

# Agent-Gym: Procedural Environment Generation and Hybrid Verifiers for Scaling Open-Weights SWE Agents

**Anonymous authors**

Paper under double-blind review

## Abstract

1 Improving open-source models on real-world SWE tasks (solving GITHUB  
 2 issues) faces two key challenges: 1) scalable curation of execution environ-  
 3 ments to train these models, and 2) optimal scaling of test-time compute.  
 4 We introduce Agent-Gym, the largest procedurally-curated executable gym  
 5 environment for training real-world SWE-agents, consisting of more than  
 6 8.7K tasks. Agent-Gym is powered by two main contributions: 1) SYNGEN:  
 7 a synthetic data curation recipe that enables scalable curation of executable  
 8 environments using test-generation and back-translation directly from com-  
 9 mits, thereby reducing reliance on human-written issues or unit tests. We  
 10 show that this enables more scalable training leading to PASS@1 of 34.4% on  
 11 SWEBENCH-VERIFIED benchmark with our 32B model. 2) Hybrid Test-time  
 12 Scaling: we next provide an in-depth analysis of two test-time scaling axes;  
 13 execution-based and execution-free verifiers, demonstrating that they ex-  
 14 hibit complementary strengths and limitations. Test-based verifiers suffer  
 15 from low distinguishability, while execution-free verifiers are biased and of-  
 16 ten rely on stylistic features. Surprisingly, we find that while each approach  
 17 individually saturates around 42-43%, significantly higher gains can be ob-  
 18 tained by leveraging their complementary strengths. Overall, our approach  
 19 achieves 51% on the SWE-Bench Verified benchmark, reflecting a new state-  
 20 of-the-art for open-weight SWE agents and for first time being competitive  
 21 with proprietary models such as o1 and sonnet-3.5-v2 w/ tools. We will  
 22 open-source our environments, models, and agent trajectories.

## 1 Introduction

24 Autonomous software engineering (SWE), aiming to solve real-world software engineering  
 25 problems such as GITHUB issues, has made significant progress in recent times (Wang et al.,  
 26 2024; Yang et al., 2024b). While LLM-based SWE-Agents have demonstrated remarkable  
 27 improvements, state-of-the-art performance is largely driven by proprietary models (An-  
 28 thropic, 2025; Jaech et al., 2024) — with open-models lagging behind (Xie et al., 2025).

29 Addressing this gap requires solving two fundamental challenges: First, scalable curation of  
 30 high-quality execution environments to train these models; and second, developing efficient  
 31 aggregation strategies to maximize test-time performance. While several benchmarks for  
 32 evaluating SWE-agents on GITHUB issues exist (Jimenez et al., 2023; Zhao et al., 2024),  
 33 scalable curation of high-quality training environments remains a challenging problem. For  
 34 instance, while the training split from SWE-Bench (Jimenez et al., 2023) contains output  
 35 patches, it lacks executable environments. Jain et al. (2024b) and Pan et al. (2024) collect  
 36 executable test environments, but the sample size and diversity of curated data are limited  
 37 by their reliance on human-written issues and test cases.

38 In this paper, we introduce AGENTGYM, the largest procedurally curated environment for  
 39 training real-world SWE-agents — consisting of more than 8.7K problems, with executable  
 40 gym environments, unit tests, and natural-language task descriptions (§2). AGENTGYM  
 41 addresses both key challenges through two primary contributions (Figures 1a and 1b):

42 **Synthetic Data Enables More Scalable Training.** We propose SYNGEN — a novel synthetic  
 43 data curation recipe that enables collection of a large number of executable training environ-

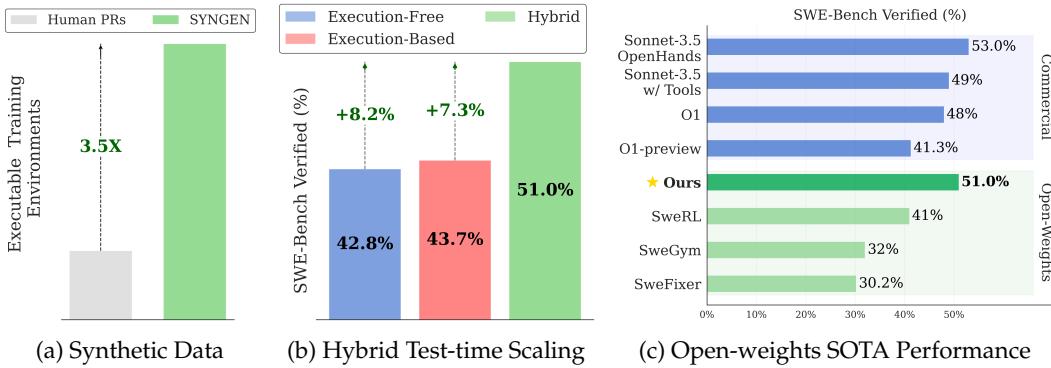


Figure 1: **Overview.** In this paper, we introduce Agent-Gym, the largest gym environment and training framework for training open-weight SWE agents. Agent-Gym is powered by two main contributions: (a) SYNGEN: a synthetic data curation recipe for curating executable training environments w/o relying on human tests and issues (§2). (b) Hybrid Inference Time Scaling: showing that while both execution-based and execution-free verifiers elicit inference-time gains; significantly better performance can be achieved by leveraging the strengths of both (§4). (c) Overall, the final approach reflects SOTA performance for open-weight SWE-agents, while also being competitive with some proprietary model baselines.

44     ments without reliance on human-written pull requests (PRs) or unit tests. We show that  
 45     instead of using human-written PRs, good-quality execution environments can directly be  
 46     curated from *commits* through backtranslation (Li et al., 2023; Wei et al., 2023) and automated  
 47     test generation (§2). Compared to PR-based data collection (Pan et al., 2024), this approach  
 48     enables more scalable data curation and agent-training, resulting in a PASS@1 performance  
 49     of 34.4% on the challenging SWEBENCH-VERIFIED benchmark (Jimenez et al., 2023).

50     **Hybrid Inference Time Scaling.** We next leverage AGENTGYM to investigate two complementary axes for scaling test-time compute (§4): 1) Execution-based verifiers that evaluate  
 51     patches through test cases (Xia et al., 2024b), and 2) Execution-free verifiers that assess  
 52     trajectories through learned models (Pan et al., 2024). While prior works have studied  
 53     these approaches in isolation, they lack a comprehensive analysis of their relative strengths  
 54     and weaknesses. We first present a unique and in-depth analysis of their working mechanisms,  
 55     demonstrating that execution-free and execution-based methods actually exhibit  
 56     complementary strengths and weaknesses. We find two key insights (studied in §4.2): a)  
 57     Execution-based methods provide direct signals for patch correctness but struggle with dis-  
 58     criminating between solutions , and b) Execution-free verifiers provide better discrimination  
 59     but can be biased by other heuristics (*e.g.*, agent thoughts) over the final patch. Based on  
 60     the above insights, we propose a hybrid scaling approach leveraging the strengths of both  
 61     methods. Surprisingly, while the performance of both execution-based and execution-free  
 62     methods plateaus around 42-43%, the hybrid approach yields significantly higher gains,  
 63     achieving a final performance of 51% on SWEBENCH-VERIFIED (Figure 1b and §4.3).

64     The key contributions of this paper are: 1) We introduce AGENTGYM, the largest procedurally  
 65     curated environment for training real-world SWE-agents, increasing the number of  
 66     executable environments by over 3 times. 2) We provide an in-depth analysis demonstrating  
 67     that execution-based and execution-free axes for scaling test-time compute exhibit comple-  
 68     mentary strengths and weaknesses. 3) Based on the above insights, we propose a *hybrid*  
 69     *scaling* approach that leverages the strengths of both methods, significantly improving  
 70     test-time performance. 4) Finally, we release an open-weights 32B model that achieves  
 71     51% on SWEBENCH-VERIFIED, reflecting a new state-of-the-art for open-weight SWE-agents,  
 72     while also for the first time demonstrating competitive or better performance compared to  
 73     commercial models (Fig. 1c), e.g., o1 (Jaech et al., 2024) and sonnet-3.5-v2 (Anthropic, 2024).

## 75     2 AGENTGYM: Procedural Synthetic Data Generation

76     **Overview.** SWE task collection methods (Jimenez et al., 2023) rely on human-written issues  
 77     and unit tests for problem statements and evaluation functions. However, this presents a

| Dataset (split)                          | Repo-Level | Executable | # Instances (total) | # Instances (train) |
|--|------------|------------|---------------------|---------------------|
| APPS (Hendrycks et al., 2021)            | ✗          | ✓          | 10,000              | 5,000               |
| HumanEval (Chen et al., 2021)            | ✗          | ✓          | 164                 | 0                   |
| R2E (Jain et al., 2024b)                 | ✓          | ✓          | 246                 | 0                   |
| SWE-Bench (train) (Jimenez et al., 2023) | ✓          | ✗          | 19,008              | 19,008              |
| SWE-Gym Raw (Pan et al., 2024)           | ✓          | ✗          | 66,894              | 66,894              |
| SWE-Bench (test) (Jimenez et al., 2023)  | ✓          | ✓          | 2,294               | 0                   |
| SWE-Gym (Pan et al., 2024)               | ✓          | ✓          | 2,438               | 2,438               |
| Agent-Gym-Lite (Ours)                    | ✓          | ✓          | 4,578               | 4,578               |
| Agent-Gym (Ours)                         | ✓          | ✓          | 8,702               | 8,702               |

Table 1: **Dataset Statistics.** Comparing statistics across different datasets curating executable training environments for SWE-agent training. Agent-Gym-lite refers to a filtered subset of tasks, with non-overlapping repositories with SWE-Bench.

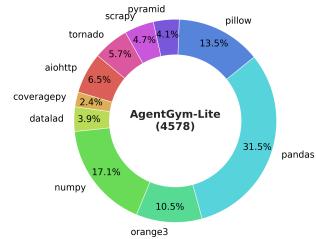


Figure 2: Repo distribution on Agent-Gym-lite subset used for training (§3).

challenge for scaling data curation as size is limited by human-written PRs. To overcome this limitation, we propose SYNGEN — a synthetic data curation recipe using backtranslation and test generation. We procedurally generate environments using only commits from GITHUB repositories, reducing reliance on both human-written issues and test cases.

**Repository and Commit Curation.** We use SEART GITHUB search<sup>1</sup> to identify PYTHON repositories with a large number of commits. Next, we extract commit history and associated code changes for each repository. We filter relevant commits using a combination of rule-based and LLM-based heuristics, identifying *interesting* code changes. For each relevant commit, we next collect build scripts by semi-manually searching across dependency pins. We expand our set of heuristics and installation procedure further in the Appendix A.

**Test-Validation and Generation for Environment Collection.** Following Jimenez et al. (2023), we use the existing test cases in the curated commits to identify Fail→Pass (F2P) test cases, i.e. test cases that fail in the original buggy commit and pass in the fixed commit. However, often, curated commits do not have associated tests, limiting the ability to use them for training environments. We supplement such commits with automatically generated Fail→Pass test-cases. Appendix A expands our test generation approach.

**Backtranslation: Non-reliance on GITHUB Issues.** Using the above steps, we collect a large number of commits, associated build environments and F2P (Fail→Pass) test cases. Now, we need to collect the problem statements associated with the commits. Prior works (Jimenez et al., 2023; Pan et al., 2024) use human-written GITHUB issues as problem statements. This inevitably cannot use the entire commit history since human-written issues are not available for all commits. Here, following Li et al. (2023); Wei et al. (2023) we propose a backtranslation approach to collect the problem statements associated with the commits.

However, naively back-translating code changes is quite noisy as models often generate generic problem statements that do not capture the essence of the code changes. Instead, we identify that human-written issues often contain failing tests and execution traces as part of bug reports. We use this observation to collect high-quality problem statements by using the F2P test-cases as part of the backtranslation prompt. Similar to existing works (Jain et al., 2024b; Zhuo et al., 2024), we find that using test execution information allows generating precise and directed problem statements. Please find prompts and examples in Appendix.

We collect over 8.7K problem statements using this approach and decontaminate the dataset to remove repository overlap with SWE-Bench repositories, obtaining 4578 problems. Tab. 1 shows the statistics of different datasets, and Fig. 2 shows the distribution of the repositories.

### 3 Training SWE-Agents using AGENTGYM Environments

**Agent Scaffolding.** We design a minimal scaffold on top of OPENHANDS (Wang et al., 2024) to experiment with agents for diverse SWE tasks. It uses a traditional REACT framework (Yao et al., 2022) without any specialized workflow; equipping the LLM with only a bash terminal, file editor, and search tool. Figure 14 depicts an example code editing trajectory.

<sup>1</sup><https://seart-ghs.si.usi.ch/>

Table 2: Resolve Rate (%) Comparison on SWE-Bench Verified and Lite benchmarks. We observe that synthetic data curation (SYNGEN) allows our approach to scale better across different model sizes. All experiments use the Qwen-2.5-Coder series for base-models.

| Model Size | SWE-Bench Lite    |                    |                            |          | SWE-Bench Verified |                    |                           |          |
|------------|-------------------|--------------------|----------------------------|----------|--------------------|--------------------|---------------------------|----------|
|            | Base-model        | SWE-Gym            | Ours                       | $\Delta$ | Base-model         | SWE-Gym            | Ours                      | $\Delta$ |
| 7B         | 1.0 ( $\pm 1.0$ ) | 10.0 ( $\pm 2.4$ ) | <b>11.0</b> ( $\pm 0.8$ )  | +1.0     | 1.8 ( $\pm 1.3$ )  | 10.6 ( $\pm 2.1$ ) | <b>19.0</b> ( $\pm 1.0$ ) | +8.4     |
| 14B        | 2.7 ( $\pm 1.9$ ) | 12.7 ( $\pm 2.3$ ) | <b>20.67</b> ( $\pm 0.7$ ) | +7.97    | 4.0 ( $\pm 1.6$ )  | 16.4 ( $\pm 2.0$ ) | <b>26.8</b> ( $\pm 1.4$ ) | +10.4    |
| 32B        | 3.0 ( $\pm 1.4$ ) | 15.3 ( $\pm 2.5$ ) | <b>23.77</b> ( $\pm 0.8$ ) | +8.47    | 7.0 ( $\pm 1.3$ )  | 20.6 ( $\pm 2.1$ ) | <b>34.4</b> ( $\pm 1.2$ ) | +13.8    |

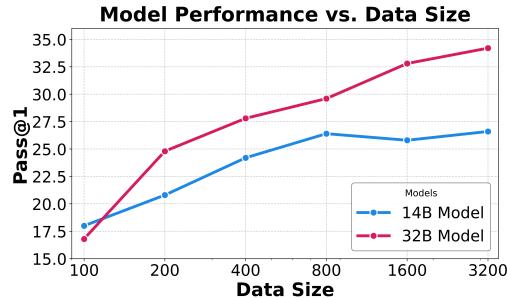


Figure 3: **PASS@1 scaling curve with increasing number of training samples.** Performance improvement with more training samples, enabled by SYNGEN approach.

| Ablation           | Config    | Pass@1 (%) |
|--------------------|-----------|------------|
| Adding Thoughts    | With      | 34.4       |
|                    | Without   | 30.4       |
| Real vs. Synthetic | Real      | 28.0       |
|                    | Synthetic | 27.8       |

Figure 4: **Top.** Use thoughts in REACT agent trajectories leads to significant performance improvements. **Bottom.** Using SYNGEN synthetic generated issues and test cases achieves similar performance as real-world issues (400 trajectories for both real & synthetic in above) while providing better scalability during data collection.

116 **Trajectory Collection and SFT Training.** We next collect SFT trajectories using from Agent-  
117 Gym environments. To avoid contamination, we only use a subset of Agent-Gym consisting  
118 of repos with no overlap with the SWE-Bench dataset. The resulting subset (Agent-Gym-  
119 lite) consists of 4578 executable environments across 10 repositories (Figure 2). For each  
120 task environment, we use CLAUDE-3.5-SONNET-V2 with our agent scaffold and collect the  
121 successful agent trajectories. Through this process, we collect 3321 trajectories from 2048  
122 unique task environments. We then use these trajectories to train our agent via supervised  
123 fine-tuning on agent thoughts and actions. For training, we use LLaMA-Factory (Zheng  
124 et al., 2024) and Qwen-2.5-Coder models (7B, 14B, 32B) as our base models. For detailed  
125 experiment configuration and hyperparameters, please refer to Appendix B.

### 126 3.1 Results and Analysis

127 **Comparison to open-weight SWE-Agents across Model Scales.** We report PASS@1 of  
128 Agent-Gym trained models on the SWEBENCH-VERIFIED and SWEBENCH-LITE benchmarks  
129 in Table 2. We also report comparisons with recently proposed SWE-Gym (Pan et al., 2024),  
130 which is most closest to our work. As seen in Tab. 2, we find that our approach enables better  
131 scaling for training SWE-agents across all model sizes. For instance, on SWE-Bench-Verified,  
132 for the same base-model type and scale, our 32B model significantly improves the pass@1  
133 performance by 14%; pushing the final performance from 20.6 (SWE-Gym) to 34.4%.

134 **Scaling with Number of Trajectories.** We investigate the relationship between training  
135 samplesize (number of trajectories) and agent performance in Fig. 3. We evaluate 14B and  
136 32B models trained with trajectory counts ranging from 100 to 3,200. Our findings indicate  
137 that performance improves with increasing trajectory count, though with diminishing  
138 returns for both models. Notably, the 14B model begins to saturate at approximately 800  
139 samples, while the 32B model still shows improvements, likely due to its larger capacity.  
140 These results extend the findings of Pan et al. (2024), who studied dataset scaling up to  $\sim 500$   
141 samples. Our analysis demonstrates that while performance does improve with increasing  
142 samplesize, the rate of improvement diminishes or even plateaus for smaller models.

143 **Real vs Synthetic Problem Statements.** The Agent-Gym approach enables us to generate  
144 problem statements without relying on human-written descriptions and test cases, offering  
145 greater scalability. We compare the performance of models trained on real GitHub issues  
146 versus our synthetic problem statements (collecting 400 trajectories from both sets). Remark-

147 ably, models trained on synthetic data achieve nearly identical performance (27.8% pass@1)  
 148 to those trained on real data (28.0%). This finding validates the efficacy of our synthetic  
 149 data generation methodology, demonstrating that procedurally generated environments  
 150 can match the training value of real-world examples while providing scalability.

151 **Explicit Thought Traces are Important.** During SFT we use both the agent’s thought  
 152 processes and actions as training targets. Models trained with thought demonstrations  
 153 achieve significantly better performance compared to those trained without (34.2% vs 30.4%  
 154 in Table 4). This suggests that exposing the model to step-by-step reasoning processes is  
 155 necessary for reliable problem-solving in complex environments.

## 156 4 Efficient Inference Time Scaling With Hybrid Verifiers

157 We utilize Agent-Gym (§2) for inference-time scaling experiments with coding agents. In  
 158 §4.1, we explore different axes for scaling test-time compute, focusing on two distinct  
 159 approaches: 1) Execution-based Verifiers and 2) Execution-free Verifiers. We analyze the  
 160 relative strengths and weaknesses of each approach, demonstrating their complementary  
 161 nature (§4.2). Based on this insight, we propose a hybrid approach that leverages the  
 162 strengths of both, significantly improving test-time performance (§4.3). Finally, we provide  
 163 detailed ablations and analysis, examining critical design choices for our approach (§4.4).

### 164 4.1 Exploring Different Axes for Training Verifiers

165 Given an input task description  $\mathcal{D}$ , a set of agent trajectories  $\{\mathcal{T}_i\}_{i=1}^K$  and candidate patch  
 166 outputs  $\{\mathcal{P}_i\}_{i=1}^K$ , our objective is to build a verifier that assigns scores  $\mathbf{S} = \{s_i\}_{i=1}^K$  to rank  
 167 the outputs. To this end, we investigate two types of verifiers:

168 **Execution-Based Verifiers.** We train a specialized *testing-agent* that generates reproduction  
 169 test cases to determine whether a candidate patch resolves the issue (i.e., whether the patch  
 170 passes the generated test suite). Additionally, following Xia et al. (2024b), we leverage  
 171 existing regression tests to filter out patches that fail to maintain backward compatibility.  
 172 Our execution-based (EB) verifier thus comprises two components: 1) a *testing-agent* that  
 173 generates targeted tests to evaluate bug fixes, and 2) a regression test filter that eliminates  
 174 patches that compromise existing functionality. Specifically, we train the testing-agent (using  
 175 Qwen-Coder-32B as base-model) to generate a comprehensive test script containing  $M = 10$   
 176 diverse tests that cover various inputs, corner cases, etc.. See Appendix D for example tests.  
 177 The final execution-based score  $s_k^{EB}$  for each patch  $\mathcal{P}_k$  is then computed as,

$$s_k^{EB} = \begin{cases} \text{TestScore}_k, & \text{if } RS_k = \max_{j \in [1, K]} RS_j, \\ 0, & \text{otherwise,} \end{cases}; \text{ where } \text{TestScore}_k = \sum_i \text{Pass}(\mathcal{P}_k, \text{Test}_i) \quad (1)$$

178 where  $RS_k$  refers to regression test score for  $k^{th}$  patch and helps select the patches with  
 179 highest regression test scores (Xia et al., 2024b).  $\text{TestScore}_k$  is simply the sum of the number  
 180 of passing tests for each patch  $\mathcal{P}_k$ . Please refer Appendix §C for further test-agent details.

181 Notably, unlike zero-shot test generation with Agentless (Xia et al., 2024b), our testing  
 182 agent interacts with the environment to examine existing test cases and generates new  
 183 tests informed by these examples with execution feedback. We demonstrate that this  
 184 environment-aware approach provides additional benefits over zero-shot methods in §4.4.

185 **Execution-free Verifiers.** We next train execution-free (EF) verifiers for selecting the best  
 186 trajectory from a set of sampled trajectories from the code-editing agent (§3). In particular,  
 187 following (Pan et al., 2024), given task description  $\mathcal{D}$ , agent-trajectory  $\mathcal{T}$  (sequence of  
 188 thought, action, and observations) and output patch  $\mathcal{P}$ , we finetune a Qwen2.5-Coder-14B  
 189 model to predict YES and NO tokens to determine correctness of a trajectory using SFT on  
 190 correct and incorrect trajectories. The execution-free score is then computed by normalizing  
 191 the relative probability of YES token as  $s^{EF} = P(\text{YES}) / (P(\text{YES}) + P(\text{NO}))$ , where  $P(\text{YES})$  and  
 192  $P(\text{NO})$  are estimated through log-probabilities of corresponding token predictions.

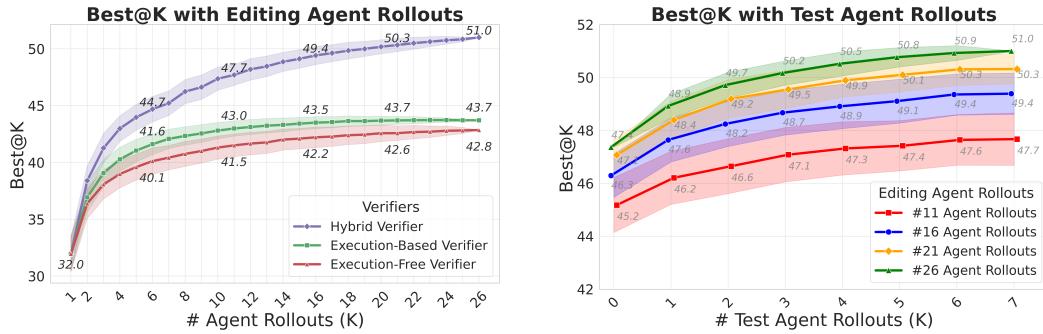


Figure 5: **Left.** BEST@K with increasing number of editing-agent rollouts. Inference-time scaling improves final performance for both execution-based and execution-free verifiers. Hybrid Verifier combining execution-based and execution-free verifiers provides significantly superior scaling. **Right.** BEST@K with increasing number of testing-agent rollouts. Increasing test-agent rollouts also improves final performance and can provide more compute efficient scaling than naively increasing only editing-agent rollouts.

Table 3: Performance of various models/methods on SWE-Bench Verified.

| Method   | Model              | Type     | Verified    |
|--|--------------------|----------|-------------|
| <b>Proprietary Models</b>                            |                    |          |             |
| Agentless-1.5 (Xia et al., 2024b)                    | GPT-4o             | Pipeline | 34.0        |
| Agentless (Xia et al., 2024b)                        | O1                 | Pipeline | 48.0        |
| Claude + Tools                                       | Claude-3.6-Sonnet  | Agent    | 49.0        |
| Agentless-1.5 (Xia et al., 2024b)                    | Claude-3.6-Sonnet  | Pipeline | 50.8        |
| OpenHands (Wang et al., 2024)                        | Claude-3.6-Sonnet  | Agent    | 53.0        |
| Claude + Tools                                       | Claude-3.7-Sonnet  | Agent    | 62.3        |
| Claude + Tools (Best@Any)                            | Claude-3.7-Sonnet  | Agent    | 70.3        |
| <b>Open-source Models</b>                            |                    |          |             |
| SWE-SynInfer (Ma et al., 2024)                       | Lingma-SWE-GPT-72B | Agent    | 30.2        |
| SWE-Fixer (Xie et al., 2025)                         | SWE-Fixer-72B      | Pipeline | 30.2        |
| SWE-Gym (BEST@16 w/ Verifier) (Pan et al., 2024)     | SWE-Gym-32B        | Agent    | 32.0        |
| SWE-RL (BEST@500 w/ Tests) (Wei et al., 2025)        | SWE-RL-70B         | Pipeline | 41.0        |
| Agentless (Xia et al., 2024b)                        | DeepSeek-R1        | Pipeline | 49.2        |
| <b>Agent-Gym-Preview (Ours)</b> (PASS@1)             | AgentGym-32B       | Agent    | <b>34.4</b> |
| <b>Agent-Gym-Preview (Ours)</b> (BEST@16 w / Hybrid) | AgentGym-32B       | Agent    | <b>49.4</b> |
| <b>Agent-Gym-Preview (Ours)</b> (BEST@26 w / Hybrid) | AgentGym-32B       | Agent    | <b>51.0</b> |

## 193 4.2 Comparative Analysis of Execution-Based and Execution-Free Verifiers

194 **Experimental Methodology.** We evaluate verifier performance using the BEST@K metric,  
195 which quantifies each verifier’s ability to identify correct patches from multiple candidates.  
196 Specifically, given  $K$  trajectories, the BEST@K metric represents the percentage of problems  
197 where the verifier successfully selects the correct patch using its scoring mechanism. For our  
198 experiments, we sample one trajectory at temperature  $T = 0$  and twenty-five trajectories at  
199 temperature  $T = 0.9$  from our 32B model for each SWEBENCH-VERIFIED benchmark problem.

200 **Both verifiers elicit inference time gains.** Figure 5 illustrates the BEST@K performance of  
201 both verifier types on the SWEBENCH-VERIFIED benchmark as a function of number of editing  
202 agent rollouts. Both execution-based and execution-free verifiers demonstrate substantial  
203 performance improvements with increased number of rollouts. However, BEST@K rate  
204 quickly plateaus for both methods, converging similarly to 43.7% and 42.8% respectively.

205 **Limited Distinguishability in Execution-Based Verifiers.** Recall that these verifiers output  
206 scores based on test pass counts and thus cannot differentiate between patches with identical  
207 pass rates, limiting their discriminative capacity. We analyze tests generated by our 32B  
208 testing agent, prompted CLAUCHE-3.5-SONNET-v2 model, and Agentless-1.5 reproduction tests  
209 (Xia et al., 2024b)<sup>2</sup>. Figure 6 (left) presents the problem density distribution for distinguisha-

<sup>2</sup>We utilize test cases from the official artifacts repository (Xia et al., 2024a).

bility rate, i.e., the proportion of tests that successfully differentiate between top-ranked correct and incorrect patches. The results demonstrate that for the majority of problems, less than 20% of tests provide discriminative signal, constraining the re-ranking. Figure 7 additionally depicts that most generated tests either do not reproduce the bug (Pass→Pass) or do not pass ground truth patches (Fail→Fail) due to bugs in the generated test cases.

**Vulnerability to Test Toxicity.** Following (Chen et al., 2022), we examine the prevalence of toxic tests, i.e., tests that pass incorrect patches but fail correct patches. Figure 6 (mid) illustrates the distribution of toxic test rates across different test generation approaches. While toxic tests are generally rare, we find that for a small but significant subset of problems, testing agents generate toxic tests (up to 10% of total tests) that can erroneously rank incorrect patches above correct ones, undermining the reliability of execution-based verification.

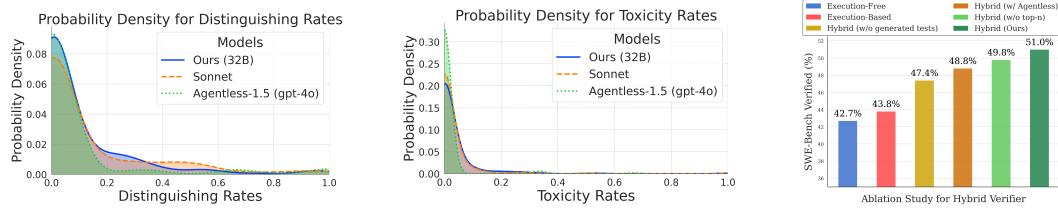


Figure 6: **Left.** Problem Probability Distributions for distinguishability rates depicting weak discrimination capabilities of tests. **Middle.** Distributions for toxicity rates showing (rare) generation of toxic tests. **Right.** Impact of design-decisions for the Hybrid Verifier.

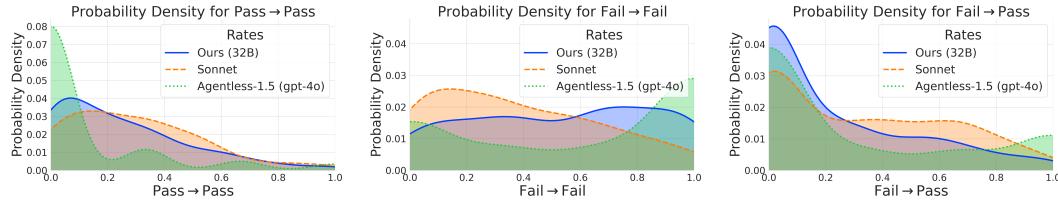


Figure 7: Problem Probability Distributions for Pass→Pass, Fail→Fail, and Fail→Pass rates for various reproduc. test generation approaches. We identify a large fraction of generated tests either do not reproduce the bug (left) or do not even pass the correct solution (middle).

**Execution-Free Verifiers can rely on heuristics.** We next study the workings and limitations of execution-free verifiers. Figure 8 (left) presents an ablation study quantifying the contribution of various trajectory components (e.g., output patch, agent thoughts) to verifier performance. We find that agent thoughts play a considerable role in determining the verifier performance. Surprisingly, the final BEST@26 drops from 42.8% to 37.6% when we remove the trajectory from the verifier input (i.e., only use the final patches). This means that while patch alone is responsible for determining the correctness, execution-free verifiers heavily rely on trajectory features, such as agent thoughts, to make predictions. To further investigate this phenomenon, we visualize attention patterns across sliding windows of the input trajectory to identify the most influential components in verifier decision-making (measured by attention allocation for the answer token). Figure 8 (right) illustrates the top two windows receiving the highest attention scores, demonstrating that verifiers disproportionately attend to agent thoughts and messages — such as expressions of confidence or uncertainty—potentially using these sentiment signals as proxies for patch correctness rather than evaluating the technical merits of the solution.

### 4.3 Hybrid Inference Time Scaling

**Combining the verifier strengths.** Given the analysis from §4.2, we can summarize two key insights: 1) Execution-based approach provides direct signal for patch correctness through execution but suffers from lack of distinguishing tests 2) Execution-free approach offers better distinguishability between patches through a continuous reward score  $s^{EF}$  but can be biased to pay more attention to heuristics (e.g., agent thoughts) over final output patch.

| Method              | Accuracy (%) | Best@26 (%) |
|---------------------|--------------|-------------|
| Final Patch + Traj. | 71.82        | 42.8        |
| Patch Only          | 68.01        | 37.6        |
| Traj. - Thoughts    | 68.77        | 41.4        |

1. Successfully reproduced the issue
2. Implemented a fix [...]
4. Ensured edge cases are handled
5. Maintained backward compatibility [...] <function=finish>submit</function> [...]

**Great!** The fix **works**. Let's see what we did to fix the issue:  
 1. We identified that the original code was failing because it was trying to use the `.inverse()` method directly on permutations, which [...]

(a) **Impact of Patch & Thoughts** on execution-free verifier. Patch alone reduces performance, indicating that model relies on other heuristics (e.g., agent thoughts) for reranking; which can be misleading (see part-b: right).

Figure 8: Quantitative and qualitative analysis on limitations of execution-free verifiers.

242 Given the above insights, we thus propose a hybrid verifier that leverages the strengths of  
 243 both approaches. Particularly, we define the hybrid verifier with score  $s_k^H$  as,

$$s_k^H = \text{Top}_n(s_k^{EF}) + s_k^{EB}, \text{ where } \text{Top}_n(s_k^{EF}) = \begin{cases} s_k^{EF}, & \text{if } s_k^{EF} \text{ is among the top } n \text{ scores,} \\ -\infty, & \text{otherwise.} \end{cases} \quad (2)$$

244 where  $s_k^{EB}$  provides execution-feedback while  $\text{Top}_n(s_k^{EF})$  provides distinguishability in case  
 245 of a tie with execution-based test scores (recall that  $s_k^{EF}$  provides a continuous score).

246 **Main Results.** Results are shown in Tab. 3 and Fig. 5. While both execution-based and  
 247 execution-free methods rapidly reach performance plateaus with increasing agent rollouts  
 248 (saturating at  $\sim 43\%$ ), our hybrid approach demonstrates substantially superior scaling  
 249 properties, yielding significant performance improvements ( $> 7\%-8\%$ ); achieving a BEST@26  
 250 performance of 51% on the challenging SWEBENCH-VERIFIED benchmark.

251 **Comparison to Open Systems.** The proposed approach significantly outperforms other  
 252 open-weight alternatives; reflecting a new state-of-the-art in this domain. Among other  
 253 generalist-agent methods, SWE-Gym (Pan et al., 2024) recently achieves a BEST@16 per-  
 254 formance of 32.0%. Similarly, concurrent work (Wei et al., 2025) recently achieve 41.0%  
 255 using RL and BEST@500 (using Agentless). In contrast, despite mainly relying on supervised  
 256 fine-tuning for training, our proposed approach achieves a PASS@1 itself of 34.4% with  
 257 BEST@26 performance of 51.0% — achieving strong performance improvements through  
 258 simply more scalable data curation (§2) and better test-time scaling (Figure 5).

#### 259 4.4 Ablation Studies on Hybrid Verification Design

260 **Variation with Test-Agent Rollouts.** As in 4.2, execution-based test generation can suffer  
 261 from a lack of distinguishing tests. One approach to address this, is to sample more test-  
 262 agent rollouts. We quantify this effect in Figure 5 (right). We observe that increasing number  
 263 of test-agent rollouts consistently helps improve performance with our hybrid approach.

264 **Compute-Efficient Rollouts.** Figure 5 (right) illustrates the BEST@K performance as a  
 265 function of both test-agent and code-editing agent rollout counts. Interestingly, we find that  
 266 sampling more test-agent rollouts can provide more compute optimized inference-scaling  
 267 over naively sampling more editing-agent rollouts. For instance, increasing the number of  
 268 editing-agent rollouts from 16 to 21 improves the BEST@K performance from 47.6% to 48.4%.  
 269 In contrast, simply sampling 5 more test-rollouts can yield better gains (BEST@K 49.3%).<sup>3</sup>

270 **Regression Tests Alone are Insufficient.** Our execution-based verification framework inte-  
 271 grates both regression and generated reproduction tests. Figure 6 (right) isolates the impact  
 272 of regression tests alone on the final performance. While regression tests alone improve  
 273 performance from 42.9% to 47.4%, using generated tests further enhances performance to  
 274 51.0%, demonstrating that both test types provide essential and complimentary signals.

275 **Agentic vs Agentless Tests.** A distinguishing feature of our approach is to train a specialized  
 276 agent for test-generation; instead of the zero-shot approach from Xia et al. (2024b). To eval-  
 277 uate this design choice, we conducted a controlled comparison using official Agentless tests  
 278 from their released artifact (Xia et al., 2024a) within our hybrid verification framework on the  
 279 SWEBENCH-VERIFIED benchmark. Figure 6 (right) demonstrates that while Agentless tests

<sup>3</sup>Note that test-agent rollouts are also usually considerably cheaper than editing-agent rollouts.

280 provide meaningful performance improvements, our agent-generated tests yield superior  
 281 results (51.0% versus 48.8%), validating our agent-based approach to test generation.

282 **Role of Top<sub>n</sub>.** We evaluate the impact of the Top<sub>n</sub> filtering mechanism introduced in Equation  
 283 (2). Figure 6 (right) shows that this selective application strategy improves performance  
 284 from 49.8% to 51.0%. This improvement likely stems from mitigating the impact of toxic tests  
 285 (§4.2) by restricting their application to higher-quality patches (identified via execution-free  
 286 reward scores  $s_k^{EF}$ ), thereby enhancing the reliability of the verification process.

## 287 5 Related Work

288 **Programming Agents.** Recent work on GITHUB issue resolution includes SWE-agent (Yang  
 289 et al., 2024b), Autocoderover (Zhang et al., 2024b), OpenHands (Wang et al., 2024), Agent-  
 290 Less (Xia et al., 2024b), Moatless Orwall (2024). All of them rely on proprietary models due  
 291 to a lack of datasets and open-weight models — a gap our work addresses.

292 **Agent Training Environments.** Existing SWE agent environments have key limitations:  
 293 SWE-Bench (Jimenez et al., 2023) lacks executable training environments, R2E (Jain et al.,  
 294 2024b) offers only 246 instances with function completion. SWE-Gym (Pan et al., 2024)  
 295 collects executable GITHUB environments similar to us but rely on human-written issues  
 296 and test cases. Synthetic data generation has been studied in various domains but our  
 297 work is the first to apply it for executable GITHUB environment collection. We use back-  
 298 translation (Li et al., 2024) and test-generation in SYNGEN approach. Please see Long et al.  
 299 (2024) for a comprehensive survey on synthetic data generation methods.

300 **SWE-Agent Training.** Ma et al. (2024) and Xie et al. (2025) train on synthetic code editing  
 301 tasks. Pan et al. (2024) study SFT on agent trajectories and inference scaling similar to our  
 302 work. Wei et al. (2025) explores reinforcement learning on large scale data collected from  
 303 real-world GITHUB issues without execution feedback.

304 **Verifiers for SWE-Coding Tasks.** Various works have explored use of verifiers for SWE tasks.  
 305 AgentLess (Xia et al., 2024b) used majority voting to select the best patch from multiple  
 306 agents. Agentless-1.5 relied on reproduction and regression tests to verify the correctness  
 307 of generated patches. Zhang et al. (2024a) proposed multi-agent committee-review (LLM  
 308 judge) to select the best patch from multiple agents. Pan et al. (2024) proposed trajectory  
 309 verifiers to re-rank the generated patches based on LLM score.

310 **Verifiers for General Coding Tasks.** Various works have explored the use of verifiers for  
 311 general coding tasks on isolated puzzles (HumanEval (Chen et al., 2021)), interviews (Jain  
 312 et al., 2024a), and competition or olympiad problems (Hendrycks et al., 2021; Li et al.,  
 313 2022) Gu et al. (2024) showed that LLM judges perform poorly on checking correctness of  
 314 generated code. Chen et al. (2022); Ridnik et al. (2024); Key et al. (2022); Zhang et al. (2023a)  
 315 study how test generation can be used to re-rank the generated code samples. Inala et al.  
 316 (2022); Zhang et al. (2023b); Ni et al. (2023) employ neural code re-ranker models.

317 In this work, we extend these lines of work by first presenting **novel insights on challenges**  
 318 **and opportunities for both execution-based and execution-free approaches in SWE-**  
 319 **Coding. Using these insights, we also propose a novel hybrid approach that effectively**  
 320 **combines their strengths** to achieve better performance (51.0% on SWEBENCH-VERIFIED).

## 321 6 Conclusion

322 In this paper, we introduce Agent-Gym, the largest gym environment and training frame-  
 323 work for scaling open-weight SWE agents. We share two key insights: 1) Synthetic data  
 324 curation can enable more scalable training on SWE tasks. 2) Hybrid-test time scaling: differ-  
 325 ent axis for test-time scaling (execution-based testing agents and execution-free verifiers)  
 326 exhibit complementary strengths; which can be leveraged to achieve significantly higher  
 327 test-time gains. Overall, our final approach achieves 51% on SWE-Bench Verified, reflect-  
 328 ing a new state-of-the-art for open-weight SWE agents, while also for first-time showing  
 329 competitive performance with some proprietary models. We hope that our work can offer  
 330 unique insights for scaling open-source SWE-agents on real-world applications.

331 **References**

- 332 Reem Aleithan, Haoran Xue, Mohammad Mahdi Mohajer, Elijah Nnorom, Gias Uddin,  
 333 and Song Wang. Swe-bench+: Enhanced coding benchmark for llms. *arXiv preprint*  
 334 *arXiv:2410.06992*, 2024. 17
- 335 Anthropic. Raising the bar on SWE-bench Verified with Claude 3.5 Sonnet. <https://www.anthropic.com/research/swe-bench-sonnet>, 2024. 2
- 336 Anthropic. Claude 3.7 sonnet. <https://www.anthropic.com/news/clause-3-7-sonnet>,  
 338 February 2025. 1
- 339 Bei Chen, Fengji Zhang, Anh Nguyen, Daoguang Zan, Zeqi Lin, Jian-Guang Lou,  
 340 and Weizhu Chen. Codet: Code generation with generated tests. *arXiv preprint*  
 341 *arXiv:2207.10397*, 2022. 7, 9
- 342 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto,  
 343 Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evalu-  
 344 ating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021. 3,  
 345 9
- 346 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,  
 347 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reichiro Nakano, et al. Training verifiers  
 348 to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021. 18
- 349 Alex Gu, Wen-Ding Li, Naman Jain, Theo X Olausson, Celine Lee, Koushik Sen, and  
 350 Armando Solar-Lezama. The counterfeit conundrum: Can code language models grasp  
 351 the nuances of their incorrect generations? *arXiv preprint arXiv:2402.19475*, 2024. 9
- 352 Dan Hendrycks, Steven Basart, Saurav Kadavath, Mantas Mazeika, Akul Arora, Ethan Guo,  
 353 Collin Burns, Samir Puranik, Horace He, Dawn Song, and Jacob Steinhardt. Measuring  
 354 coding challenge competence with apps. *NeurIPS*, 2021. 3, 9
- 355 Jeevana Priya Inala, Chenglong Wang, Mei Yang, Andres Codas, Mark Encarnación, Shu-  
 356 vendu Lahiri, Madanlal Musuvathi, and Jianfeng Gao. Fault-aware neural code rankers.  
 357 *Advances in Neural Information Processing Systems*, 35:13419–13432, 2022. 9
- 358 Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low,  
 359 Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. Openai o1 system card.  
 360 *arXiv preprint arXiv:2412.16720*, 2024. 1, 2
- 361 Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Ar-  
 362 mando Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contami-  
 363 nation free evaluation of large language models for code. *arXiv preprint arXiv:2403.07974*,  
 364 2024a. 9
- 365 Naman Jain, Manish Shetty, Tianjun Zhang, King Han, Koushik Sen, and Ion Stoica. R2e:  
 366 Turning any github repository into a programming agent environment. In *ICML 2024*,  
 367 2024b. 1, 3, 9
- 368 Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and  
 369 Karthik Narasimhan. Swe-bench: Can language models resolve real-world github issues?  
 370 *arXiv preprint arXiv:2310.06770*, 2023. 1, 2, 3, 9
- 371 Darren Key, Wen-Ding Li, and Kevin Ellis. I speak, you verify: Toward trustworthy neural  
 372 program synthesis. *arXiv preprint arXiv:2210.00848*, 2022. 9
- 373 Xian Li, Ping Yu, Chunting Zhou, Timo Schick, Omer Levy, Luke Zettlemoyer, Jason Weston,  
 374 and Mike Lewis. Self-alignment with instruction backtranslation. *arXiv preprint*  
 375 *arXiv:2308.06259*, 2023. 2, 3
- 376 Xian Li, Ping Yu, Chunting Zhou, Timo Schick, Omer Levy, Luke Zettlemoyer, Jason E  
 377 Weston, and Mike Lewis. Self-alignment with instruction backtranslation. In *The Twelfth  
 378 International Conference on Learning Representations*, 2024. URL [https://openreview.net/  
 379 forum?id=1oijHJBRsT](https://openreview.net/forum?id=1oijHJBRsT). 9

- 380 Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond,  
 381 Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, et al. Competition-level code  
 382 generation with alphacode. *Science*, 378(6624):1092–1097, 2022. 9
- 383 Lin Long, Rui Wang, Ruixuan Xiao, Junbo Zhao, Xiao Ding, Gang Chen, and Haobo Wang.  
 384 On llms-driven synthetic data generation, curation, and evaluation: A survey. *arXiv*  
 385 preprint [arXiv:2406.15126](https://arxiv.org/abs/2406.15126), 2024. 9
- 386 Yingwei Ma, Rongyu Cao, Yongchang Cao, Yue Zhang, Jue Chen, Yibo Liu, Yuchen  
 387 Liu, Binhua Li, Fei Huang, and Yongbin Li. Lingma swe-gpt: An open development-  
 388 process-centric language model for automated software improvement. *arXiv preprint*  
 389 [arXiv:2411.00622](https://arxiv.org/abs/2411.00622), 2024. 6, 9
- 390 Ansong Ni, Srini Iyer, Dragomir Radev, Veselin Stoyanov, Wen-tau Yih, Sida Wang, and  
 391 Xi Victoria Lin. Lever: Learning to verify language-to-code generation with execution. In  
 392 *International Conference on Machine Learning*, pp. 26106–26128. PMLR, 2023. 9
- 393 A. Orwall. Moatless tool. <https://github.com/aorwall/moatless-tools>, 2024. Accessed:  
 394 2024-10-22. 9, 17
- 395 Jiayi Pan, Xingyao Wang, Graham Neubig, Navdeep Jaitly, Heng Ji, Alane Suhr, and Yizhe  
 396 Zhang. Training software engineering agents and verifiers with swe-gym, 2024. URL  
 397 <https://arxiv.org/abs/2412.21139>. 1, 2, 3, 4, 5, 6, 8, 9, 18, 19
- 398 Tal Ridnik, Dedy Kredo, and Itamar Friedman. Code generation with alphacodium: From  
 399 prompt engineering to flow engineering. *arXiv preprint arXiv:2401.08500*, 2024. 9
- 400 Xingyao Wang, Boxuan Li, Yufan Song, Frank F Xu, Xiangru Tang, Mingchen Zhuge, Jiayi  
 401 Pan, Yueqi Song, Bowen Li, Jaskirat Singh, et al. Openhands: An open platform for ai  
 402 software developers as generalist agents. *arXiv preprint arXiv:2407.16741*, 2024. 1, 3, 6, 9,  
 403 17
- 404 Yuxiang Wei, Zhe Wang, Jiawei Liu, Yifeng Ding, and Lingming Zhang. Magicoder: Source  
 405 code is all you need. *arXiv preprint arXiv:2312.02120*, 2023. 2, 3
- 406 Yuxiang Wei, Olivier Duchenne, Jade Copet, Quentin Carbonneaux, Lingming Zhang,  
 407 Daniel Fried, Gabriel Synnaeve, Rishabh Singh, and Sida I. Wang. Swe-rl: Advancing  
 408 llm reasoning via reinforcement learning on open software evolution. *arXiv preprint*  
 409 [arXiv:2502.18449](https://arxiv.org/abs/2502.18449), 2025. 6, 8, 9
- 410 Chunqiu Steven Xia, Yinlin Deng, Soren Dunn, and Lingming Zhang. Agentless: De-  
 411 mystifying llm-based software engineering agents. <https://github.com/OpenAutoCoder/>  
 412 Agentless, 2024a. 6, 8
- 413 Chunqiu Steven Xia, Yinlin Deng, Soren Dunn, and Lingming Zhang. Agentless: Demysti-  
 414 fying llm-based software engineering agents. *arXiv preprint arXiv:2407.01489*, 2024b. 2, 5,  
 415 6, 8, 9
- 416 Chengxing Xie, Bowen Li, Chang Gao, He Du, Wai Lam, Difan Zou, and Kai Chen. Swe-fixer:  
 417 Training open-source llms for effective and efficient github issue resolution. *arXiv preprint*  
 418 [arXiv:2501.05040](https://arxiv.org/abs/2501.05040), 2025. 1, 6, 9
- 419 An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li,  
 420 Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong  
 421 Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jin Xu, Jingren  
 422 Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin  
 423 Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize  
 424 Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu  
 425 Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin  
 426 Wei, Xuancheng Ren, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong  
 427 Liu, Zeyu Cui, Zhenru Zhang, and Zhihao Fan. Qwen2 technical report. *arXiv preprint*  
 428 [arXiv:2407.10671](https://arxiv.org/abs/2407.10671), 2024a. 18

- 429 John Yang, Carlos E Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik  
 430 Narasimhan, and Ofir Press. Swe-agent: Agent-computer interfaces enable automated  
 431 software engineering. *arXiv preprint arXiv:2405.15793*, 2024b. 1, 9
- 432 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and  
 433 Yuan Cao. React: Synergizing reasoning and acting in language models. *arXiv preprint*  
 434 *arXiv:2210.03629*, 2022. 3, 17
- 435 Kexun Zhang, Danqing Wang, Jingtao Xia, William Yang Wang, and Lei Li. Algo: Synthesizing  
 436 algorithmic programs with generated oracle verifiers. *arXiv preprint arXiv:2305.14591*,  
 437 2023a. 9
- 438 Kexun Zhang, Weiran Yao, Zuxin Liu, Yihao Feng, Zhiwei Liu, Rithesh RN, Tian Lan, Lei  
 439 Li, Renze Lou, Jiacheng Xu, et al. Diversity empowers intelligence: Integrating expertise  
 440 of software engineering agents. In *The Thirteenth International Conference on Learning*  
 441 *Representations*, 2024a. 9
- 442 Tiansi Zhang, Tao Yu, Tatsunori Hashimoto, Mike Lewis, Wen-tau Yih, Daniel Fried, and  
 443 Sida Wang. Coder reviewer reranking for code generation. In *International Conference on*  
 444 *Machine Learning*, pp. 41832–41846. PMLR, 2023b. 9
- 445 Yuntong Zhang, Haifeng Ruan, Zhiyu Fan, and Abhik Roychoudhury. Autocoderover:  
 446 Autonomous program improvement. In *Proceedings of the 33rd ACM SIGSOFT International*  
 447 *Symposium on Software Testing and Analysis*, pp. 1592–1604, 2024b. 9
- 448 Wenting Zhao, Nan Jiang, Celine Lee, Justin T Chiu, Claire Cardie, Matthias Gallé,  
 449 and Alexander M Rush. Commit0: Library generation from scratch. *arXiv preprint*  
 450 *arXiv:2412.01769*, 2024. 1
- 451 Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyuan Luo, Zhangchi Feng, and  
 452 Yongqiang Ma. Llamafactory: Unified efficient fine-tuning of 100+ language models. In  
 453 *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume*  
 454 *3: System Demonstrations)*, Bangkok, Thailand, 2024. Association for Computational  
 455 Linguistics. URL <http://arxiv.org/abs/2403.13372>. 4, 18, 19
- 456 Terry Yue Zhuo, Minh Chien Vu, Jenny Chim, Han Hu, Wenhao Yu, Ratnadira Widayasa,  
 457 Imam Nur Bani Yusuf, Haolan Zhan, Junda He, Indraneil Paul, et al. Bigcodebench:  
 458 Benchmarking code generation with diverse function calls and complex instructions.  
 459 *arXiv preprint arXiv:2406.15877*, 2024. 3

---

460 **A Dataset Details**

461 **Commit Filtering Heuristics.** Our commit filtering approach employs multiple heuristics  
 462 to identify high-quality bug fixes and improvements suitable for training data. We filter for  
 463 manageable changes by constraining the scope of modifications, focusing on commits with  
 464 limited but meaningful edits. We distinguish between test and non-test files, prioritizing  
 465 commits where modified components have corresponding test changes—a strong signal for  
 466 intentional bug fixes rather than simple refactoring. Our heuristics also identify substantive  
 467 modifications beyond documentation updates, ensuring we capture genuine functional  
 468 improvements. We further analyze entity-level changes to verify that commits represent  
 469 focused, well-tested modifications rather than large-scale refactoring or feature additions.  
 470 These filtering strategies collectively enable the identification of precise, verifiable code  
 471 changes that represent meaningful bug fixes or improvements, providing high-quality  
 472 examples for training software engineering agents. To ensure consistency and quality, we  
 473 employ specific thresholds in our filtering process:

- 474     • Maximum of 5 non-test files modified in a single commit
- 475     • Maximum of 100 edited lines across all non-test files
- 476     • Maximum patch length of 2000 characters to ensure focused changes
- 477     • No more than 1 deleted entity in non-test files
- 478     • Maximum of 3 added entities in non-test files
- 479     • Maximum of 3 edited entities in non-test files
- 480     • No more than 10 statement-level changes to maintain tractability

481 These constraints help ensure that the selected commits represent focused, manageable  
 482 changes that are suitable for training software engineering agents. Additionally, we use  
 483 LLM as a judge filter to further refine our dataset.

484 **Repository Installation.** Installing historical commits from GitHub repositories presents  
 485 significant challenges due to evolving dependency requirements and API changes. We use  
 486 a Docker-based approach with a search-based dependency resolution strategy to create  
 487 reproducible environments for each commit. Our installation process follows these steps:

- 488     1. Extract dependency information from the repository (requirements.txt, setup.py,  
      etc.)
- 489     2. Identify potential version conflicts and compatibility issues
- 490     3. Generate multiple candidate dependency configurations
- 491     4. Test each configuration until a working environment is found

493 This process is semi-manual and challenging to scale and we aim to rely more on LLMs in  
 494 the future. The pseudocode below illustrates our search-based hole-filling approach:

```
495 # Search-based dependency resolution with hole-filling
496 function resolve_dependencies(repository, commit_hash):
497     # Extract dependency constraints from repository
498     constraints = extract_constraints(repository, commit_hash)
499
500     # Identify "holes" - unspecified or conflicting version requirements
501     holes = find_dependency_holes(constraints)
502
503     # Generate candidate configurations by filling holes
504     candidates = []
505     for hole in holes:
506         # Generate multiple version candidates for each hole
507         candidates.append(generate_version_candidates(hole)) # semi-
508                         # manual
509
510     # Try configurations until one works
```

```

512     for config in generate_configurations(candidates):
513         docker_env = create_docker_environment(config)
514         if test_installation(docker_env, repository, commit_hash):
515             return docker_env
516
517     # If all configurations fail, try fallback strategies
518     return apply_fallback_strategies(repository, commit_hash)
519
520 # Example of hole-filling for Python dependencies
521 function generate_version_candidates(dependency):
522     if dependency == "numpy":
523         return ["1.17.*", "1.20.*", "1.26.*"]
524     elif dependency == "python":
525         return ["3.7", "3.8", "3.10"]
526     elif dependency == "cython":
527         return ["<0.30", "==3.0.5"]
528     # Add more dependencies as needed
529
530     # Use LLM to suggest versions for unknown dependencies
531     return query_llm_for_version_candidates(dependency)

```

Listing 1: Repository installation search procedure

533 Example installation scripts test multiple dependency combinations sequentially, exiting on  
534 the first successful build:

```

535 # Attempt with first configuration
536 if build_and_check_pandas "3.7" "1.17.*" "<0.30" "62.*" "0.23"; then
537     echo "[INFO]_First_combo_succeeded._Exiting."
538     exit 0
539 fi
540
541 # Attempt with second configuration
542 if build_and_check_pandas "3.8" "1.20.*" "<0.30" "62.*" "0.23"; then
543     echo "[INFO]_Second_combo_succeeded._Exiting."
544     exit 0
545 fi
546
547 # Attempt with third configuration
548 if build_and_check_pandas "3.10" "1.26.*" "==3.0.5" "62.*" "0.23"; then
549     echo "[INFO]_Third_combo_succeeded._Exiting."
550     exit 0
551 fi

```

Listing 2: Example installation script excerpt

554 This approach allows us to create working environments for historical commits, enabling  
555 execution-based validation of our dataset.

556 **Test Generation.** We use an Agentless-like reproduction test generation approach. A key  
557 difference is that we use the ground truth patch as context when generating the tests.

558 **Issue Generation.** As discussed in the main paper, we use backtranslation to generate  
559 synthetic issues for commits that lack human-written GitHub issues. Our approach lever-  
560 ages both the code changes in the commit and the test execution results to create realistic,  
561 informative issue descriptions. The issue generation process follows these steps:

- 562 1. Extract failing test functions from the execution results
- 563 2. Analyze test outputs to identify error messages and expected behaviors
- 564 3. Provide the LLM with commit hash, commit message, code patch, and test execution  
565 results
- 566 4. Guide the LLM to generate a concise, informative issue that describes the bug  
567 without revealing the solution

568 For each commit, we extract and utilize specific components:

- 569 • **Commit metadata:** Hash and commit message provide context about the change
- 570 • **Code patches:** We separate non-test file changes (showing what was fixed) from
- 571 test file changes (showing how to verify the fix)
- 572 • **Test execution:** We include both old commit (failing) and new commit (passing)
- 573 execution outputs
- 574 • **Test functions:** We extract up to three relevant test functions that demonstrate the
- 575 bug
- 576 • **Assertion failures:** We extract and format the failing assertions from the old commit
- 577 to show error details

578 The prompt construction carefully organizes these components to give the LLM sufficient  
 579 context while focusing attention on the most relevant information for issue generation. We  
 580 carefully design our prompting strategy to ensure the generated issues resemble human-  
 581 written ones, focusing on clarity, naturalness, and providing sufficient information for  
 582 understanding the bug. The implementation of our issue generation pipeline involves  
 583 several key functions working together:

```
584 # Build the complete prompt with all components
585 def get_prompt(commit, execution_result, issues=None):
586     # Include commit hash and message
587     # Include commit patch (non-test files)
588     # Include test file changes
589     # Include execution results from old and new commits
590     # Include improved test functions
591     # Include test function code
592     # Include assertion failures
593     # Include example issues and instructions
```

Listing 3: Issue generation code structure

596 The template below shows our structured approach to issue generation:

597 As you are trying to generate synthetic issues, you will follow these  
 598 guidelines:  
 599

- 600 1. Keep the issue concise and informative.
- 601 2. Describe the failing test, including the input that causes the failure  
     , the nature of the failure, and the expected behavior. Do NOT  
     mention test functions or files directly.
- 602 3. Do not reveal the solution to the problem in the issue. Only describe  
     the bug and the expected behavior.
- 603 4. If there are multiple failing tests, focus on the most informative one  
     or a subset that best describes the general nature of the failure.
- 604 5. Describe the expected output of the failing test:  
     - For errors, describe the error message.  
     - For failing tests, mention what is supposed to happen.
- 605 6. Write the issue as a human would, using simple language without  
     excessive formatting.
- 606 7. Use concrete terms to describe the nature of the failure. Avoid vague  
     terms like "specific output" or "certain data".
- 607 8. INCLUDE test code to describe the bug but keep it brief and relevant.  
     Truncate or simplify tests longer than 5-6 lines.
- 608 9. Do not mention external files unless absolutely necessary.
- 609 10. Format code snippets using triple backticks.

610 The issue should include:  
 611 1. A clear and concise title  
 612 2. A description of the problem with detailed example buggy code  
 613 3. Expected behavior  
 614 4. Actual behavior or error message

626

Listing 4: Issue generation template

627 This approach enables us to generate high-quality synthetic issues that provide clear problem  
 628 statements for our training data, even for commits that lack human-written issues. Below  
 629 are examples of synthetic issues generated using our approach:

```

630 **Title:** Calling `load()` Before `draft()` Causes `draft()` to Fail for
631   JPEG Images
632
633 **Description:** When generating a thumbnail for a JPEG image using the `thumbnail()`
634   method, the method calls `load()` before `draft()`. This sequence
635   results in the `draft()` method returning `None`, which prevents the
636   thumbnail from being properly optimized.
637
638 **Example Code:**  

639 ```python  
from PIL import Image  
  
with Image.open("Tests/images/hopper.jpg") as im:  
    im.thumbnail((64, 64))  
```  

640
641 **Expected Behavior:** The `thumbnail()` method should utilize the `draft()` method to optimize
642   the image size before loading, ensuring that the thumbnail is resized
643   correctly and efficiently.
644
645 **Actual Behavior:** The `draft()` method returns `None` because `load()` is invoked before it
646   . This prevents the thumbnail from being optimized, potentially
647   leading to incorrect thumbnail sizes or unnecessary memory usage.
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
```

Listing 5: Example synthetic issue for a PIL image thumbnail bug

```

658 **Title:** Unable to Register Route with Names Containing Both Dots and
659   Colons
660
661 **Description:** After merging branch '0.18', attempting to register a route with a name
662   that includes both dots (`.`) and colons (`:`) results in a `ValueError`. The recent changes were intended to allow route names to
663   be a sequence of Python identifiers separated by dots or colons, but
664   this combination is still causing issues.
665
666 **Example Code:**  

667 ```python  
from aiohttp.web import UrlDispatcher, PlainRoute  
  
def handler(request):  
    return 'Hello'  
  
router = UrlDispatcher()  
  
# Attempting to register a route with both dots and colons in the name  
route = PlainRoute('GET', handler, 'test.test:test', '/handler/to/path')  
router.register_route(route)  
```  

668
669 **Expected Behavior:** Registering a route with a name like `test.test:test` should succeed
670   without errors, as the name follows the updated rules allowing
671   multiple identifiers separated by dots or colons.
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
```

```

687
688 **Actual Behavior:**  

689 A `ValueError` is raised with the message:  

690 ``  

691 ValueError: Incorrect route name value, Route name should be a sequence  

692     of python identifiers separated by dot or colon  

693 ``  

694 This prevents the registration of route names that include both dots and  

695     colons, contrary to the intended flexibility introduced in the recent  

696     commit.  

697

```

Listing 6: Example synthetic issue for a route name validation bug

698 **Patch Minimization.** We identify that the ground-truth patches often contain irrelevant  
 699 code changes that are not required to fix the bug, often making modifications to style and  
 700 structure of the programs. We implement a patch-minimization approach to identify the  
 701 minimal set of code changes required to fix the bug by iteratively removing the code changes  
 702 and checking whether the tests still pass. This allows us to collect fine-grained signal for  
 703 evaluating localization capabilities of LLMs.

## 704 B SFT Training

### 705 Agent Details.

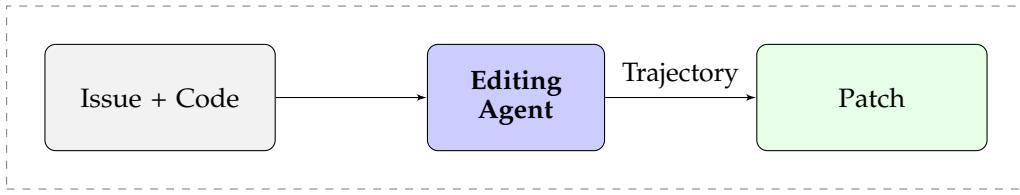


Figure 9: Code-editing agent architecture: The agent takes an issue description and codebase as input and produces a patch that fixes the issue.

706 We use Agent-Gym to train a general-purpose prompting agent. In particular, we train  
 707 our code-editing agent on tasks from Agent-Gym, where given an executable environment  
 708  $\mathcal{E}$  and problem description  $\mathcal{D}$ , the agent is asked to solve the provided issue using any  
 709 means necessary. Particularly, unlike (Orwall, 2024), we do not rely on the use of specialized  
 710 workflows. The agent is tasked to solve the entire task end-to-end, including writing its  
 711 own reproduction scripts, finding the bug, proposing a fix and then testing its correctness.  
 712 Similar to (Wang et al., 2024), the agent is also provided with a finish tool, allowing it to  
 713 submit a solution if it thinks it has completed the task.

714 **Agent and Tools.** Similar to (Aleithan et al., 2024; Wang et al., 2024), we adopt the traditional  
 715 REACT format (Yao et al., 2022) for agent-design. For AGENTHUB, we use a minimalistic set  
 716 of four tools to enable the agent to perform diverse SWE tasks: 1) file\_editor: for viewing  
 717 and editing files, 2) search\_tool: for searching a relevant term in a given file or folder, 3)  
 718 execute\_bash: allowing execution of non-interactive bash commands (e.g., for running test  
 719 scripts), 4) submit: for ending the current trajectory while returning expected outputs. No  
 720 internet or browser access is provided to the agent during the training process.

721 **Data Curation.** For training, we use supervised finetuning with rejection sampling using  
 722 trajectories from sonnet-3.5 model for supervision. To avoid contamination, we only use  
 723 a subset of Agent-Gym consisting of repos with no overlap with the SWE-Benchdataset.  
 724 The resulting subset (Agent-Gym-lite) consists of 4538 executable environments across 10  
 725 repositories (Figure 2). Overall, we collect a total of 3321 successful trajectories from 2048  
 726 unique test environments. For rejection sampling we use the unit tests from Agent-Gym  
 727 environments (both synthetic and existing). For each trajectory, we use a maximum of  
 728  $N = 40$  steps. Also, we limit the number of tokens per-trajectory to 32K max tokens. Finally,

729 we also use a maximum timeout of 10-min for the overall trajectory and 90 seconds for  
 730 each action execution, in order to avoid cases where the agent launches a long-running  
 731 background process. We collect all training data using a temperature of 0.2.

732 **Training Setup and Hyperparameters.** For training, we use the Qwen-2.5-Coder 7B, 14B  
 733 and 32B series as the base model for training SWE-agents on Agent-Gym. For training  
 734 we perform full SFT using the above collected trajectories using LLaMA-Factory (Zheng  
 735 et al., 2024). We train the overall model for a total of 2 epochs, batch size as 8 while using a  
 736 learning rate of  $1e - 5$ . The warmup ratio for training was set to 0.1. Due to computational  
 737 constraints, a maximum context length of 20K was used for training the agent. In future,  
 738 the use of context-parallelism can enable us to further push the performance when training  
 739 SWE-agents on more complex tasks requiring larger-context lengths.

## 740 C Inference Time Scaling

### 741 C.1 Execution-Based Testing Agents

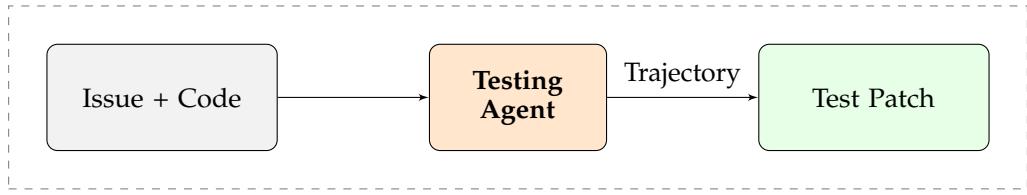


Figure 10: Testing agent architecture: The agent generates comprehensive test cases to verify if a candidate patch resolves the issue.

742 **Agent Details.** We train a specialized *testing-agent* that generates reproduction test cases  
 743 to determine whether a candidate patch resolves the issue (i.e., whether the patch passes  
 744 the generated test suite). Specifically, we train the testing-agent (using QWEN-CODER-32B  
 745 as base-model) to generate a comprehensive test script containing  $M = 10$  diverse tests  
 746 that cover various inputs, corner cases, etc. We use the same agent scaffold from Sec. 3 for  
 747 training the testing agent.

748 **Data Curation.** For training, we use supervised finetuning using trajectories from  
 749 sonnet-3.5 model for supervision. Overall, we collect a total of 2203 test-generation tra-  
 750 jectories from sonnet (both positive and negative without any rejection sampling). For  
 751 each trajectory, we use a maximum of  $N = 40$  steps. Also, we limit the number of tokens  
 752 per-trajectory to 20K max tokens. Finally, we also use a maximum timeout of 5-min for the  
 753 overall trajectory and 60 seconds for each action execution, in order to avoid cases where  
 754 the agent launches a long-running background process.

755 **Training Setup and Hyperparameters.** For training, we use the QWEN-CODER-32B model as  
 756 the base model. We then use the above collected training SFT trajectories to perform full  
 757 finetuning with the QWEN-CODER-32B model using LLaMA-Factory (Zheng et al., 2024). We  
 758 train the overall model for a total of 2 epochs, batch size as 8 while using a learning rate of  
 759  $1e - 5$ . A maximum context length of 20K was used for training the agent. The warmup  
 760 ratio for training was set to 0.1.

### 761 C.2 Execution-Free Verifiers

762 **Verifier Details.** In addition to the execution-based “testing agents”, we also explore the  
 763 execution-free outcome-supervised reward models (a.k.a verifiers) (Cobbe et al., 2021). In  
 764 particular, given a problem statement  $\mathcal{D}$ , agent-trajectory  $\mathcal{T} = \{a_1, o_1, a_2, o_2, \dots, a_n, o_n\}$  and  
 765 output patch  $\mathcal{O}$  from the code-editing agent on the Agent-Gym environments, we train  
 766 a Qwen2.5-Coder-14B model (Yang et al., 2024a) to output a scalar score value  $s^{EF} \in [0, 1]$   
 767 predicting the probability of output patch being correct. Specifically, following (Pan et al.,  
 768 2024) we output the correctness of each patch through output tokens YES (correct) and NO

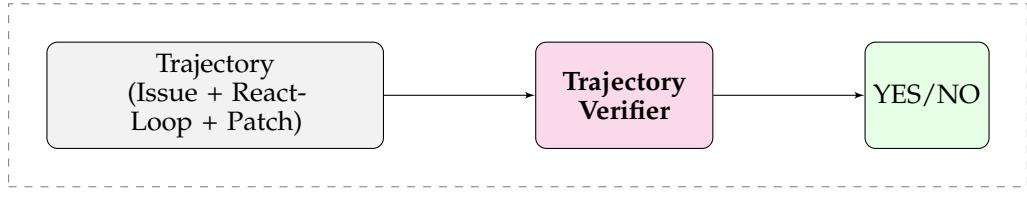


Figure 11: Execution-free verifier architecture: The verifier predicts whether a patch is correct based on the full trajectory without executing the code.

(incorrect). The overall reward score is then computed by normalizing the relative probability of YES token as  $r = P(\text{YES})/(P(\text{YES}) + P(\text{NO}))$ , where  $P(\text{YES})$  and  $P(\text{NO})$  are estimated through the log-probabilities of the corresponding token predictions.

**Training Data.** We first use the trajectories collected for code-editing agent training §3 in order to obtain a collection of positive and negative samples for verifier training. Following the best configuration from (Pan et al., 2024), we also generate on-policy trajectories using our trained 32B model. We then filter the collected samples to have an equal number of positive and negative samples. The overall dataset consists of 5700 total trajectories including both positive and negative samples. For training, we follow the template from (Pan et al., 2024), asking the LLM model to predict the output as YES for positive and NO for negative trajectories.

**Training Setup and Hyperparameters.** For training, we use the QWEN-CODER-14B model as the base model. We then use the above collected training SFT trajectories to perform finetuning using LLaMA-Factory (Zheng et al., 2024). Similar to (Pan et al., 2024), we perform LORA finetuning using a rank of 64. We train the overall model for a total of 2 epochs, batch size of 8 while using a learning rate of  $1e - 5$ . A maximum context length of 32K was used for training the agent. The warmup ratio for training was set to 0.1.

### 786 C.3 Execution-Based Analysis

787 In our analysis of execution-based testing agents, we focus on two key metrics: distinguishability and toxicity of generated tests. These metrics help us understand the effectiveness 788 and limitations of execution-based verification.

790 **Distinguishability Rate.** The distinguishability rate measures a test's ability to differentiate 791 between correct and incorrect patches. A test is considered "distinguishing" if it behaves 792 differently when applied to correct patches versus incorrect patches. In practical terms, this 793 means the test can help us identify which patches are correct and which are not.

794 For example, consider a test that passes for all correct patches but fails for all incorrect 795 patches—this test has perfect distinguishability. Conversely, a test that passes (or fails) for 796 both correct and incorrect patches provides no useful signal for distinguishing between 797 them.

798 Mathematically, for a given test  $t$  and a set of patches  $P$  divided into correct patches  $P_c$  and 799 incorrect patches  $P_i$ , the distinguishability is defined as:

$$\text{Distinguish}(t) = \mathbb{1} \left[ \max_{p \in P_i} \text{Pass}(p, t) \neq \max_{p \in P_c} \text{Pass}(p, t) \right] \quad (3)$$

800 where  $\text{Pass}(p, t)$  indicates whether patch  $p$  passes test  $t$ , and  $\mathbb{1}[\cdot]$  is the indicator function. 801 This formula checks whether the best-performing incorrect patch behaves differently on the 802 test compared to the best-performing correct patch. The distinguishability rate for a set of 803 tests  $T$  is then the average distinguishability across all tests:

$$\text{DistinguishRate}(T) = \frac{1}{|T|} \sum_{t \in T} \text{Distinguish}(t) \quad (4)$$

804 In our analysis, we found that most generated tests have low distinguishability  
 805 rates—typically less than 20% of tests can effectively differentiate between correct and  
 806 incorrect patches. This limitation significantly impacts the ability of execution-based verifi-  
 807 cation to identify the best patches, especially as the number of candidate patches increases.

808 **Toxicity Rate.** We define toxic tests as those that incorrectly favor incorrect patches over  
 809 correct ones. The toxicity rate is the proportion of tests that exhibit this behavior. Mathemati-  
 810 cally:

$$\text{Toxic}(t) = \mathbb{1} \left[ \max_{p \in P_i} \text{Pass}(p, t) > \max_{p \in P_c} \text{Pass}(p, t) \right] \quad (5)$$

811 The toxicity rate for a set of tests  $T$  is:

$$\text{ToxicityRate}(T) = \frac{1}{|T|} \sum_{t \in T} \text{Toxic}(t) \quad (6)$$

812 While toxic tests are generally rare, they can significantly impact verification reliability  
 813 when present, with toxicity rates reaching up to 10% for some problems. These findings  
 814 highlight the importance of generating diverse, high-quality tests and the value of combining  
 815 execution-based verification with other approaches, such as execution-free verifiers, to  
 816 achieve more robust results.

#### 817 C.4 Execution-Free Analysis

818 Figure 12 shows the limitations of the execution-free verifier.

## 819 D Example Testing Agent Outputs

820 This section provides examples of test cases generated by our approach.

### 821 D.1 Example 1: SymPy Relational Parsing Tests

822 The following example shows a truncated test suite for validating relational parsing in  
 823 SymPy, demonstrating our approach’s ability to generate multiple test cases. This test was  
 824 generated to address the issue in [SymPy PR #24661](#), which fixes relational parsing in the  
 825 SymPy library.

```
826 from sympy import Lt, Gt, Le, Ge, Eq, Ne
827
828 def test_relational_parsing():
829     # Test case 1: Basic less than operation
830     try:
831         result = parse_expr('1 < 2', evaluate=False)
832         expected = Lt(1, 2, evaluate=False)
833         if str(result) == str(expected):
834             print("Test_Case_1:_Issue_resolved")
835         else:
836             print("Test_Case_1:_Issue_reproduced")
837     except Exception as e:
838         print("Test_Case_1:_Other_issues")
839
840     # Test case 2: Greater than operation
841     try:
842         result = parse_expr('3 > 2', evaluate=False)
843         expected = Gt(3, 2, evaluate=False)
844         if str(result) == str(expected):
845             print("Test_Case_2:_Issue_resolved")
846         else:
```

| Metric (Execution-free scaling)   | [Final Patch + Traj.] | [Final Patch Only]  | [Final Patch + Traj. - Thoughts] |
|---|-----------------------|---|----------------------------------|
| Reward Model Accuracy (%)   | 71.82%                | 68.01%  | 68.77%                           |
| Aggregation Performance (Pass@26)   | 42.8%                 | 37.6%   | 41.4%                            |
| Since we've:  |                       | [ASSISTANT]   |                                  |
| 1. Successfully reproduced the issue<br>2. Implemented a fix<br>3. Verified the fix works<br>4. Ensured edge cases are handled<br>5. Maintained backward compatibility  |                       | Great! The fix works. Let's analyze what we did to fix the issue:   |                                  |
| I believe we can now submit this as the solution:   |                       | 1. We identified that the original code was failing because it was trying to use the `~.inverse()` method directly on permutations, which doesn't exist.  |                                  |
| <function=finish><br><parameter=command>submit</parameter><br><parameter=result>Fixed the `~check_homomorphism` function in `homomorphisms.py` to properly handle inverted generators in permutation ...  |                       | 2. We fixed the issue by:<br>- Using the correct operator `**-1` for getting inverses which works consistently for both Permutations and FreeGroupElements<br>- Splitting the logic into two clear cases: |                                  |
| [ASSISTANT]<br>The changes look good. I've made the following improvements:<br>1. Simplified the logic to explicitly check for both the generator and its inverse<br>2. Used the `**-1` operator which works for both Permutation and FreeGroupElement<br>3. Maintained the same functionality but with more robust handling of inverted generators |                       | Let's run the reproduce script again to verify the fix:   |                                  |
| Let's run the reproduce script again to verify the fix:<br><function=execute_bash><br><parameter=cmd>python<br>→ reproduce_issue.py</parameter><br></function><br>[/ASSISTANT]  |                       | <function=execute_bash><br><parameter=cmd>python<br>→ reproduce_issue.py</parameter><br></function><br>[/ASSISTANT]   |                                  |
| [USER]<br>Exit code: 0<br>Execution output of [execute_bash]:<br>[STDOUT]   |                       | Testing with DihedralGroup(3):<br>SUCCESS: Homomorphism created successfully  |                                  |
|   |                       | Testing with SymmetricGroup(3):<br>SUCCESS: Homomorphism created successfully   |                                  |

Figure 12: **Analyzing limitations of execution-free verifiers.** Top: Quantitative ablations on the impact of output patch on reward model (RM) performance. Down: Qualitative visualization analyzing top  $k = 4$  sliding windows over agent trajectory (refer §4.2) with highest mean attention score (sympy–sympy-24443: SWE-Bench). We observe that the RM can be biased by the agent thought / actions, instead of relying on the final output patch.

```

848     print("Test_Case_2:_Issue_reproduced")
849 except Exception as e:
850     print("Test_Case_2:_Other_issues")
851
852 # ... [6 more test cases omitted for brevity] ...
853
854 # Test case 9: Chained comparisons
855 try:
856     result = parse_expr('1 < x < 2', evaluate=False)
857     if isinstance(result, bool):
858         print("Test_Case_9:_Issue_reproduced")
859     else:
860         print("Test_Case_9:_Issue_resolved")
861 except Exception as e:
862     print("Test_Case_9:_Other_issues")

```

Listing 7: Test cases for SymPy relational parsing (truncated). Successfully detects incorrect code from correct code.

## 864 D.2 Example 2: Django Model Choice Field Validation Tests

865 The following truncated example demonstrates tests for Django's ModelChoiceField validation with proper setup and teardown code. This test was generated to address the issue in  
 866 [Django PR #13933](#), which improves error messages in ModelChoiceField validation.  
 867

```

868
869 from django.forms import ModelChoiceField, ModelMultipleChoiceField
870 from django.db import models
871 from django.core.exceptions import ValidationError
872 from django.test import TestCase
873 from django.db import connection
874
875 # Create a simple test model
876 class TestModel(models.Model):
877     name = models.CharField(max_length=100)
878
879     def __str__(self):
880         return self.name
881
882     class Meta:
883         app_label = 'test_app'
884
885 def setup_test_environment():
886     # Create the test table
887     with connection.schema_editor() as schema_editor:
888         schema_editor.create_model(TestModel)
889
890     # Create some test data
891     TestModel.objects.create(name="Option_1")
892     TestModel.objects.create(name="Option_2")
893     TestModel.objects.create(name="Option_3")
894
895 def cleanup_test_environment():
896     # Drop the test table
897     with connection.schema_editor() as schema_editor:
898         schema_editor.delete_model(TestModel)
899
900 def test_modelchoice_validation():
901     try:
902         # Test Case 1: Basic invalid choice for ModelChoiceField
903         field = ModelChoiceField(queryset=TestModel.objects.all())
904         try:
905             field.clean(999) # Non-existent ID
906             print("Test_Case_1:_Issue_resolved")
907         except ValidationError as e:
908             if "999" not in str(e): # Value should be in error message
909                 print("Test_Case_1:_Issue_reproduced")
910             else:
911                 print("Test_Case_1:_Issue_resolved")
912
913         # ... [4 more test cases omitted for brevity] ...
914
915         # Test Case 6: Valid choice
916         obj = TestModel.objects.first()
917         try:
918             result = field.clean(obj.id)
919             if result == obj:
920                 print("Test_Case_6:_Issue_resolved")
921             else:
922                 print("Test_Case_6:_Issue_reproduced")
923         except ValidationError:
924             print("Test_Case_6:_Issue_reproduced")
925     except Exception as e:
926         print(f"Unexpected_error:{e}")

```

927

Listing 8: Test cases for Django ModelChoiceField validation (truncated). Most test cases error due to unhandled exceptions and do not distinguish.

## 928 E Agent Trajectory Example

929 This section provides a visual example of an agent's trajectory while solving a software  
 930 engineering task. The sequence shows the step-by-step process from problem statement to  
 931 solution, demonstrating how our agent approaches and solves real-world programming  
 932 issues.

### Problem Statement

```
'PolyElement.as_expr()' not accepting symbols
The method `PolyElement.as_expr()`

https://github.com/sympy/sympy/blob/193e3825645d93c73e31cdceb6d742cc6919624d/sympy/polys/rings.py#L618-L624

is supposed to let you set the symbols you want to use, but, as it stands, either you pass the wrong number
of symbols, and get an error message, or you pass the right number of symbols, and it ignores them, using
`self.ring.symbols` instead:

```python
>>> from sympy import ring, ZZ, symbols
>>> R, x, y, z = ring("x,y,z", ZZ)
>>> f = 3*x**2*y - x*y*z + 7*z**3 + 1
>>> U, V, W = symbols("u,v,w")
>>> f.as_expr(U, V, W)
3*x**2*y - x*y*z + 7*z**3 + 1
```

```

Figure 13: Problem statement presented to the agent, describing the issue that needs to be resolved.

933 This trajectory example illustrates several key aspects of our agent's problem-solving ap-  
 934 proach:

- 935 • **Systematic Exploration:** The agent methodically explores the codebase to under-  
 936 stand the context and locate the issue.
- 937 • **Root Cause Analysis:** Rather than addressing symptoms, the agent identifies the  
 938 underlying cause of the problem using `test_issue.py`.
- 939 • **Solution Development:** The agent formulates a clear plan before implementing  
 940 changes.

941 These capabilities enable our agent to effectively tackle complex software engineering tasks  
 942 that require deep understanding of code structure, programming language semantics, and  
 943 software design principles.

**Step 1**

**Thought:**

I'll help you fix the issue with 'PolyElement.as\_expr()' not accepting symbols. Let's follow the steps systematically.

1. First, let's explore the repository to understand its structure:

**Action:**

```
<function>
<parameter>search_term=PolyElement</parameter>
</function>
```

**Observation:**

```
Execution output of [search]:
Found 196 matches for "PolyElement" in /testbed:
./sympy/utilities/tests/test_monicring.py (1 matches)
./sympy/printing/tests/test_rings.py (1 matches)
./sympy/printing/rings.py (1 matches)
./sympy/printing/repr.py (1 matches)
./sympy/printing/str.py (1 matches)
./sympy/printing/tests/test_rings.py (1 matches)
./sympy/printing/tests/test_repr.py (2 matches)
./sympy/printing/tests/test_str.py (1 matches)
./sympy/polys/rings.py (17 matches)
./sympy/polys/modular_ugcd.py (63 matches)
```

(a) Step 1: Initial analysis and exploration

**Step 2**

**Thought:**

Great! I can see that 'PolyElement' is primarily defined in '/testbed/sympy/polys/rings.py'. Let's look at the specific implementation:

**Action:**

```
<function>
<parameter>command=view</parameter>
<parameter>path=/testbed/sympy/polys/rings.py</parameter>
<parameter>view_range={610, 636}</parameter>
</function>
```

**Observation:**

```
Execution output of [file_editor]:
Here's the code of running 'cat -n' on the file: /testbed/sympy/polys/rings.py:
 610     if self.ring == new_ring:
 611         return self
 612     elif self.ring.symbols == new_ring.symbols:
 613         terms = list(zip(dict_reorder(self, self.ring.symbols), new_ring.symbols))
 614         return new_ring._from_terms(terms, self.ring.domain)
 615     else:
 616         return new_ring._from_dict(self, self.ring.domain)
 617
 618     def as_expr(self, symbols):
 619         if symbols and len(symbols) != self.ring.ngens:
 620             raise ValueError("not enough symbols, expected %s" % (self.ring.ngens,
 621                                         len(symbols)))
```

(b) Step 2: Detailed Exploration

**Step 3**

**Thought:**

Now I understand the issue. The 'as\_expr()' method is not properly using the symbols passed to it. Let's first create a script to reproduce the issue:

**Action:**

```
<function>
<parameter>command=create</parameter>
<parameter>path=/testbed/reproduce_issue.py</parameter>
<parameter>file_text=<from sympy import Ring, ZZ, symbols
import signal
import sys

def timeout_handler(signum, frame):
    print("Script timed out after 20 seconds")
    sys.exit()

# Set timeout of 20 seconds
signal.signal(signal.SIGALRM, timeout_handler)
signal.alarm(20)

def test_as_expr():
    # Create polynomial ring and polynomial
    R, x, y, z = Ring("x,y,z", ZZ)
    f = 3*x**2*y - x*y + 7*x**3 + 1

    # Create new symbols
    U, V, W = symbols("u,v,w")>>>
```

(c) Step 3: Reproducing the issue

**Step 4**

**Thought:**

Let's run the reproduction script to confirm the issue:

**Action:**

```
<function>
<parameter>cmd=python /testbed/reproduce_issue.py</parameter>
</function>
```

**Observation:**

```
Exit code: 0
Execution output of [execute_bash]:
[STDOUT]
[STDERR]

Correctly raised ValueError for wrong number of symbols
Expression with new symbols: 3*x**2*y - x*y + 7*x**3 + 1
ERROR: Original symbols are still in the expression
Tests failed!
[STDERR]
```

(d) Step 4: Running reproduction tests

**Step 5**

**Thought:**

The issue is confirmed! The script shows that:

1. The error for wrong number of symbols works correctly
2. When the correct number of symbols is provided, they are not being used in the output expression

Looking at the code in 'rings.py', the issue is in the 'as\_expr()' method. The problem is that when symbols are provided, they are correctly validated for number but then immediately overwritten with 'self.ring.symbols' on line 622. This is clearly a bug.

Let's fix the code:

**Action:**

```
<function>
<parameter>path=/testbed/reproduce_issue.py</parameter>
<parameter>cmd=python /testbed/reproduce_issue.py</parameter>
<parameter>id_start as_expr(self, symbols):
    if symbols and len(symbols) != self.ring.ngens:
        raise ValueError("not enough symbols, expected %s" % (self.ring.ngens,
len(symbols)))
    else:
        symbols = self.ring.symbols

    return expr_from_dict(self.as_expr_dict(), symbols)</parameter>
<parameter>name_asr_id_start as_expr(self, symbols):
    if symbols:
        if len(symbols) != self.ring.ngens:
            raise ValueError("not enough symbols, expected %s" % (self.ring.ngens,
```

(e) Step 5: Implementing the fix

**Step 6**

**Thought:**

The changes look good. The key changes are:

1. Changed the logic to only check symbol count if symbols are provided
2. Only use 'self.ring.symbols' as a fallback when no symbols are provided
3. Maintained the same error checking for wrong number of symbols

Let's run the reproduction script again to verify the fix:

**Action:**

```
<function>
<parameter>cmd=python /testbed/reproduce_issue.py</parameter>
</function>
```

**Observation:**

```
Exit code: 0
Execution output of [execute_bash]:
[STDOUT]
[STDERR]

Correctly raised ValueError for wrong number of symbols
Expression with new symbols: 3*x**2*y - x*y + 7*x**3 + 1
SUCCESS: New symbols are correctly used in the expression
All tests passed!
[STDERR]
```

(f) Step 6: Verifying the fix

Figure 14: Short successful agent trajectory (using our 32B model) showing the step-by-step process of solving a software engineering task. The agent analyzes the problem, identifies the root cause, implements a solution, tests it, and verifies that it resolves the issue.