

SCALABLE SUPERVISING SOFTWARE AGENTS WITH PATCH REASONER

Anonymous authors

Paper under double-blind review

ABSTRACT

While large language model agents have advanced software engineering tasks, the unscalable nature of existing test-based supervision is limiting the potential improvement of data scaling. The reason is twofold: (1) building and running test sandbox is rather heavy and fragile, and (2) data with high-coverage tests is naturally rare and threatened by test hacking via edge cases. In this paper, we propose R4P, a patch verifier model to provide scalable rewards for training and testing SWE agents via reasoning. We consider that patch verification is fundamentally a reasoning task, mirroring how human repository maintainers review patches without writing and running new reproduction tests. To obtain sufficient reference and reduce the risk of reward hacking, R4P uses a group-wise objective for RL training, enabling it to verify multiple patches against each other’s modification and gain a dense reward for stable training. R4P achieves 72.2% Acc. for verifying patches from SWE-bench-verified, surpassing OpenAI o3. To demonstrate R4P’s practicality, we design and train a lite scaffold, Mini-SE, with *pure* reinforcement learning where all rewards are derived from R4P. As a result, Mini-SE achieves 26.2% Pass@1 on SWE-bench-verified, showing a 10.0% improvement over the original Qwen3-32B. This can be further improved to 33.8% with R4P for test-time scaling. The stable scaling curves in both RL test rewards and test-time accuracy reflect R4P’s practical utility for scalable supervision on software agents.

1 INTRODUCTION

Large language models (LLMs) agents have made notable progress in software engineering (SWE) thanks to rule-based reinforcement learning (RL) and test-time scaling (TTS) (Jimenez et al., 2023; Luo et al., 2025a). However, unlike tasks with extensive publicly available and easy-to-check answers (e.g., math), SWE tasks naturally have *non-unique* answers (patches) that are hard to *formally verify* (ter Beek et al., 2024). Thus, *testing* was adopted as a proxy approach (Tihanyi et al., 2025) for verifying patches in SWE’s RL and TTS. Yet, testing is *heavy*, hindering data scaling in SWE. Specifically, building test environment (e.g., Docker images) for each data instance is highly labor-intensive (Pan et al., 2024; Jain et al., 2025), and maintaining sandbox instances (e.g., Docker containers) for parallel testing is unstable due to surge workload (Luo et al., 2025a). Furthermore, the inherent *test coverage problem* (Yu et al., 2025) also harms the reliability of test results, making tests being easily *hacked* by edge cases when expanding data to less maintained GitHub projects. We found that only 28.11% of resolved issues have corresponding tests on projects with over 500 stars (Appendix A). This hinders the utilization of vast, unlabeled data from the open-source community.

We identify the core challenge is *how to acquire test-free supervision for patches at scale*. Ideally, a reward model (Minghao Yang, 2024; Shiwen et al., 2024; Chen et al., 2025) could serve as an alternative patch verifier. However, they are trained to judge the *relative* preference of solutions than their *absolute* correctness. Since SWE is a hard reasoning task, all LLM rollout answers may be incorrect and should not be rewarded, making such relative assessment inadequate. Furthermore, without test execution, for an LLM to reliably identify subtle incompatibilities, it needs an exhaustive, agentic search across the repository’s call graph and dependencies. However, this approach makes the verification process as hard and slow as the patch generation process itself (Meng et al., 2024), making it impractical for scalable supervision. To mitigate this, some test-time verifiers (Pan et al., 2024; Jain et al., 2025; Luo et al., 2025a) use certain agent’s trajectories as an additional input, trained to

maximize the log probability of token YES or NO from a process-oriented perspective. However, such task formulation creates a strong coupling between the model and a specific interaction style, hindering its ability to generalize to other agents.

In this paper, we introduce *R4P*, a *reasoning* patch verifier model for SWE agents that do not rely on any golden test, developer patch, agent trajectory, or run-time sandbox. It enables efficient data scaling for real-world issues that are untested or have low test coverage, offering potential for continual learning beyond the constraints of test-based sandbox data. We consider patch verification is fundamentally a reasoning task which can be improved via RL. This mirrors how real-world repository maintainers evaluate pull request (PR) patches based on their understanding of the repository and a detailed analysis of the patch, rather than relying on writing and running test scripts to judge correctness (Baum et al., 2016). To avoid reward hacking and enable effective test-free verification, R4P adopts a group-wise training objective. Specifically, for a given issue and a set of candidate patches, R4P assesses each patch by comparing it against others in the group for mutual contextual information, compensating for the absence of tests and facilitating the detection of subtle errors. Moreover, by expanding the verifier’s solution space beyond binary classification, R4P reduces the risk of reward hacking and offers a denser reward compared to the original sparse bipolar reward. This leads to more stable training and improved convergence. We evaluate R4P on patch generated by four different agents from SWE-bench-verified Experiments. When built upon Qwen2.5-Coder-32B-Instruct, R4P achieves 72.2% accuracy, surpassing advanced models like OpenAI o3.

To further demonstrate R4P’s practicality, we developed *Mini-SE*, a lite, execution-free agentic scaffold with issue-resolving-oriented code *search* and *edit* capabilities. We then use R4P for supervision and train Mini-SE based on Qwen3-32B model on R2E-Gym issues without using its sandbox. While Mini-SE does not test its generated patches during rollout, its verification burden is transferred to R4P for patch selection for test-time scaling. The Pass@1 resolution rate on SWE-bench-verified steadily improves with more training data and finally reaches 26.2%, outperforming Lingma Agent + Lingma SWE-GPT-72B (Ma et al., 2024). With R4P for patch selection, it can be further improved to 33.8%. In addition, R4P verifies each patch within a second on average, much faster than the minute-level time cost of testing. These illustrates the practicality of reasoning-based verification.

2 PRELIMINARIES

Reinforcement Learning: Reinforcement learning (RL) aims to learn an optimal policy π_θ that maximizes the expected cumulative discounted reward r when interacting with an environment. In the context of auto-regressive LLM, state at step t is the concatenation of prompt x and current response $y_{<t}$, and the action is the t -th token y_t . Empirically, policy gradient methods are widely used to directly optimizes this objective J and updates model parameters θ :

$$\nabla_\theta J(\theta) = \mathbb{E}_{x \sim D, y \sim \pi_\theta} \left[\sum_{t=0}^T \nabla_\theta \log \pi_\theta(y_t | x, y_{<t}) A_t \right] \quad (1)$$

For scalable rewarding, reward models are trained to provide a dense, normalized scalar reward (e.g., $r_s \in [0, 1]$), reflecting human preference rankings of answers. They have shown strength in human-AI alignment tasks like writing and safety. More recently, sparse verifiable rewards (e.g., $r_v = \{0, 1\}$) are proven to be effective on reasoning tasks with ground-truth correctness like math (Shao et al., 2024) and coding (Luo et al., 2025b), giving rise to a new domain, Reinforcement Learning with Verifiable Rewards (RLVR) (Schulman et al., 2017; Guo et al., 2025).

Test-Time Scaling: Test-time scaling (TTS) refers to approaches of utilizing test-time compute to enhance the answer quality. Given a certain test-time compute budge N , TTS aims to maximize the expectation E of the number of model output $y^*(x)$ equals to the ground-truth y , which is notated as indicator function $\mathbb{I}_{y=y^*(x)}$ for a given prompt x , which equals 1 if the predicted label y_i for patch p_i matches the ground-truth label $y^*(p_i)$, and 0 otherwise (Snell et al., 2024). Assume the output tokens y follows the model distribution $D(\theta|N, x)$, an ideal strategy θ^* should be:

$$\theta^*(N) = \operatorname{argmax}_\theta \left(\mathbb{E}_{y \sim D(\theta|N, x)} [\mathbb{I}_{y=y^*(x)}] \right) \quad (2)$$

Common TTS approaches include expanding multi-turn searching or revisions during rollout (Madaan et al., 2023; Wang et al., 2024) and sampling multiple candidates for filtering or selec-

tion after rollout (Wang et al., 2022; Xia et al., 2024) against verifier. This verifier could also be rules or reward models, as they requires similar supervision signals like rewards in RL. For example, Snell et al. train a process reward model to enable Monte Carlo tree search (Sutton et al., 1998) during inference, while Xia et al. use an LLM for scaling tests to filter flawed answers after rollout.

LLM for Issue Resolving: Software Engineering (SWE) tasks have increasingly gained attention, especially for real-world repository-level issue resolving tasks (Jimenez et al., 2023). Given an issue description in natural language (e.g., a feature request or a bug report), the models need to scan the whole repository to understand the problem, locate the place for editing, and finally submit a patch as the answer. The golden tests (i.e., the unit tests submitted in the developer patches) and corresponding sandbox environments (typically based on Docker) will then be applied to verify the patch. While such test-based verification is effective for benchmarking, it cannot provide scalable supervision for improving model’s SWE capability. We identified two major problems:

P.1: Testing is heavy: Even with the help of LLM, building a single issue’s Docker image takes about 7 human minutes (Jain et al., 2025). Thus, existing datasets contain less than 5,000 real-world executable issue instances (Pan et al., 2024). While *issue synthesis* allows for data scaling in a reusable sandbox via commit backtranslation (Yang et al., 2025), they are proven to be much less effective than real-world issues (Luo et al., 2025a). Furthermore, maintaining hundreds of Docker containers (e.g., 64 batch size \times 8 rollouts) for parallel online rewarding is highly unstable, as containers are prone to crashing under heavy peer load (Luo et al., 2025a). In addition, while a timeout threshold is necessary to avoid process blocking from infinite loop of flawed patch, it can also introduce false negative rewards on instances running heavy test suites.

P.2: Test coverage problem: In practice, developer’s unit tests often fail to cover all edge cases, allowing tricky patches to pass without genuinely solving the issue (Mockus et al., 2009). This problem persists even in widely recognized high quality dataset, SWE-bench (Yu et al., 2025), and will be eventually more exacerbated in less actively maintained repositories for RL scaling. Such test coverage problem will unavoidable results in reward hacking after running out of existing real-world high quality issues. Thus, it is challenging to continue exploiting existing tests for data scaling in-the-wild. Furthermore, our study shows that 28.11% issues from projects with more than 500 stars even do not contain *any* tests, let along the test coverage problems (see Appendix A). This prevents the exploitation of vast, unlabeled data from the open-source community (e.g., GitHub, Jira).

Challenges of Model-based Verification: A way to scale supervision is to train a reward model (RM) as patch verifier, which judges LLM responses without relying on tests and environments. However, it is usually infeasible in practice to incorporate RMs into SWE due to two challenges:

C.1 Lack of sufficient reference: In addition to scalar and trajectory RMs discussed in Sec. 1, we explored the capability of advanced LLMs as RMs for patch verification (see Sec. 5.3). Our analysis of failure cases reveals the reason is due to the lack of dependency context for the patch, preventing the model from identifying subtle incompatibilities or bugs with the existing codebase. While agentic search is a potential solution, precisely identifying all issue-relevant context is both difficult and inefficient, as an issue in a large project can be far-reaching (Meng et al., 2024).

C.2 Sparsity of outcome space: Since the output space of patch verification is binary (Pass/Fail), the supervision signal for RM training is very sparse, making naively mapping input issue-patch pairs to a Pass/Fail label very challenging. In practice, diverse patches can solve the same issue without a dense, relative relationship between each other, which further precludes the use of dense training objectives like Bradley-Terry model (Kendall & Smith, 1940). In addition, simple binary classification SFT offers limited benefits to open-source models. Furthermore, the binary outcome reward is very easy to hack, making the training unstable.

3 R4P

To provide scalable patch verification for SWE tasks, we propose R4P, a reasoning reward model for patch verification trained via pure rule-based RL. We consider patch verification is naturally a reasoning task, as human developers review pull requests (PRs) by reasoning about the static code changes, rather than designing new, issue-specific golden tests for every submission. R4P enables efficient data scaling for unlabeled or low-test-coverage real-world issues, thereby supporting continuous learning beyond the limited test-based sandbox data.

Task formulation To overcome the challenges of reference deficiency (C.1) and sparse rewards (C.2), R4P introduces a group-wise task formulation. Given a SWE issue I in issue description and a group of patches $P = [p_1, p_2, \dots, p_N]$ trying to resolve the issue, R4P is expected to generate a sequence of tokens, which include the reasoning tokens and the ID of the correct patches y_i :

You are a software expert. You will be given a software issue and some patch candidates in user query. You need to judge which patch(es) can resolve the issue. Carefully review, critic, and compare the given candidates. You need to first think about the reasoning process in the mind until you get the final answer. Finally, put the ID(s) of correct patch candidates in `\boxed{\}`.

[ISSUE] {user issue I } [/ISSUE]
 [PATCH 1] {agent patch P_1 } [/PATCH 1]
 ...
 [PATCH N] {agent patch P_N } [/PATCH N]

This group-level verification fully leverages the non-unique and diverse characteristics of patches, since the models are tended to modify various positions with different code changes, where each patch provides sufficient information for judging other patches. As a result, R4P could cross-reference all patches' edit locations and content, inferring potential context and identifying subtle errors in one patch by observing others. Furthermore, this group-wise approach transforms the binary outcome space into a much denser one (e.g., $\sum C_N^i$ possibilities for random selecting correct patches from a group $P = [p_1, \dots, p_N]$), which provides a richer supervision signal and significantly mitigates the reward hacking risk inherent in simple binary classification tasks. Specifically, during RL training, the reward of R4P $r(P)$ is calculated by the indicator function $\mathbb{I}_{y_i=y^*(p_i)}$, which equals 1 if the predicted label y_i for patch p_i matches the ground-truth label $y^*(p_i)$, and 0 otherwise.

$$r(P) = \begin{cases} \frac{1}{N} \sum_{i=1}^N \mathbb{I}_{y_i=y^*(p_i)} & \text{if answer is boxed} \\ 0 & \text{if answe is unboxed} \end{cases} \quad (3)$$

Reward modeling To avoid reward hacking caused by the model's random guesses on a patch's binary correctness when uncertain, we adopt a continuous group-wise reward requiring the model to verify as many correct patches as possible (i.e., group-wise accuracy). The intuition is: a reasoning that correctly verifies more patches should be considered better than one that verifies fewer, even though certain random guessing is inevitable. Thus, the reward clearly reflect how much better one reasoning result is compared to another, rather than treating real correct reasoning as equivalent to a randomly guessed correct answer. Note that when the number of patches to be verified (N_v) exceeds the group size (N_t), the set of patches is partitioned into smaller groups for verification, with the total number of inference calls being $\lceil N_v/N_t \rceil$. Conversely, if the number of patches is less than N_t , the group is padded with empty patches to match the input format used during training, ensuring consistent model behavior.

Data sampling To collect a dataset of patches that closely resemble those generated by contemporary SWE agents, we use Claude-3.7-sonnet on OpenHands (Wang et al., 2024) scaffold to resolve issues in SWE-Gym (Pan et al., 2024). For each of the 2,438 issue instances in SWE-Gym, we sampled 6 patch candidates, resulting in a total of 14,628 verified patches. Since the patch accuracy of this process is only 30%, the original dataset was highly imbalanced with a majority of incorrect patches. Thus, we filtered a subset of the incorrect patches to create a more balanced label distribution. These data is then used for RL training of R4P via GRPO (Guo et al., 2025) (Appendix B.1).

4 MINI-SE

To validate R4P's practicality, we aim to train an agent and evaluate its performance with all supervision from R4P. To achieve it, a straightforward approach is to adapt existing agent scaffolds (Yang et al., 2024; Wang et al., 2024). However, this faces two major challenges. First, their tool designs are often overly general. The bash-like or python-like commands often leads to excessively long agent trajectories due to atomic actions like `cd` or `ls`. Furthermore, some complex actions like

grep may have a vast argument space, which is error-prone for LLMs, involving multiple interactions of trial-and-errors. This makes the RL rollout and backward processes very inefficient for validating R4P’s RL supervision capability. Second, they typically follow a “generate-then-verify” workflow, attempting to generate and execute tests to validate their own patches during rollout. However, since existing methods typically adopts a separate and more powerful TTS verifier for final patch selection (Jain et al., 2025), this in-loop self-verification introduces redundant interactions and further hampers training efficiency.

To enable an efficient validation of R4P’s practicality, we introduce Mini-SE, a lightweight scaffold that focuses on issue resolving without overlong trial-and-error iterations and avoid dependency on executable environment. It contains:

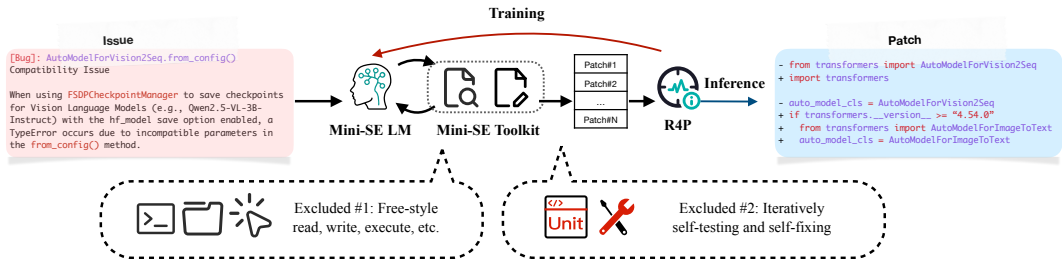


Figure 1: Mini-SE adopts an test-free issue-resolution-oriented tool design strategy: **Search**: input an *entity* name, output a file path and *code* snippet. **Edit**: input a file path, old code and new code snippet, output a diff *patch*. This design prevents the inefficient, redundant iterations caused by free-style exploration during training rollouts. It also removes self-testing and fixing when generating patches since R4P could take this role for penalizing wrong ones during training and selecting the correct one during inference. This further improves efficiency and avoid relying on sandbox.

Code Search A tool that maps an input entity’s simple name to corresponding source code and file path. Upon issue input, it initialize a code graph database for entire repository, extracting all *class*, *method*, and *function* entities with static analyzer. To reduce the interaction turns, this tool returns the source code of *all* entities with that simple name (e.g., `func`). This *fuzzy matching* also reduces tool call failures due to incorrectly spelled fully qualified names (e.g., `file.Class.func`).

Code Edit A tool that replace an old string to a new string in the input file path, which could work as *insert*, *delete*, and *modify* actions. Upon issue input, it checkouts the specific commit and creates a shallow clone of the target repository as a workspace. All modifications are validated using *syntax checker*, and any change with syntax errors is *discarded*. This check reduces failures from the LLM repeatedly trying to fix syntax errors introduced in previous wrong fixes.

Rollout process of Mini-SE Mini-SE focuses solely on issue resolution, generally following a targeted “*entity* → *code* → *patch*” paradigm. Specifically, it begins with “symptom” entities mentioned in the issue description, examines its code, and iteratively traces through the dependency chain to identify the “root cause” entities. Once located, it modifies the corresponding code and submits a patch. The absence of trial-and-error interactions due to incorrectly using general-purpose tools enables Mini-SE swiftly generating a large volume of candidate patches. In addition, Mini-SE decouples the verification process from the generation process. During rollout, it directly leverages the reward signal provided by R4P to reinforce correct patches and penalize incorrect ones without self-test-self-fix, leading to highly efficient training. The saved token budget can be reallocated to patch verifier during test-time patch selection. All of these features makes Mini-SE an efficient scaffold for validating R4P’s practicality.

5 EXPERIMENT

5.1 EXPERIMENTAL SETUP

Implementation We implement R4P using `Qwen-2.5-Coder-Instruct-32B` with group size $N = 4$, as it offers a balanced trade-off between patch verification capability and computational

Table 1: Comparison with general models.
 (“*” mark means point-wise patch verification)

Model	Acc	F1	EM
claude-3.7-sonnet	68.1	50.5	26.8
claude-4-sonnet	68.4	50.0	25.7
gemini-2.5-pro	72.7	56.6	34.6
o4-mini	68.5	52.0	29.6
gpt-4o	61.2	43.0	17.3
o3	71.5	57.4	36.4
qwen-2.5-coder-7b*	54.0	N/A	N/A
qwen-2.5-coder-7b	60.0	43.3	17.9
qwen-2.5-coder-32b*	55.9	N/A	N/A
qwen-2.5-coder-32b	61.9	44.5	18.5
R4P	72.2	63.3	41.8

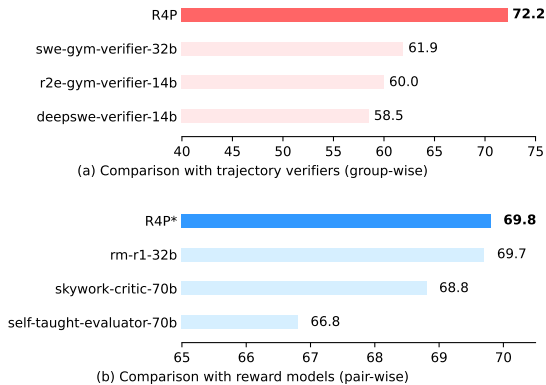


Figure 2: Comparison with specialized models.
 (“*” mark means pair-wise patch selection)

cost. We do not adopt Qwen-3 series for R4P because enabling its thinking mode results in excessively long reasoning chains (around 32K tokens), which significantly increases computational overhead during both training and inference. Disabling thinking mode, however, leads to substantial performance degradation. In contrast, for Mini-SE, we base it on Qwen-3-32B because the agentic interaction pattern involves large volumes of tool-returned tokens rather than model-generated ones, making its training more efficient than that of R4P. Additionally, it demonstrates better tool-use instruction compliance compared to Qwen 2.5 series. Mini-SE’s training details are in Appendix B.2.

Datasets We collect patches generated by four different agents and LLMs from submissions in the SWE-bench-verified experiments, with a resolution rate of around 50% to ensure label balance (Appendix B.3). To ensure a fair comparison, we adopt the same group-wise formulation for general models as used in R4P. To maintain consistent task difficulty, all empty patches are removed. The final dataset comprises 1,340 patches, forming 335 group-wise instances. To further assess the practicality of R4P, we use real-world issues from R2E-Gym along with rewards generated by R4P to train Mini-SE via RL. This dataset contains 4,578 issues from 10 repositories, where we exclude any issues that overlap with R4P’s training data. We evaluate Mini-SE on SWE-bench-verified, which consists of 500 issues with sandbox testing.

5.2 MAIN RESULTS

Effectiveness on patch verification As detailed in Table 1, we compare R4P with various advanced generative and reasoning LLMs. R4P achieves a verification accuracy of 72.2%, surpassing strong proprietary models like OpenAI o3, despite R4P having a significantly smaller parameter scale. Since group-wise verification is a multi-choice retrieval task, we also report the F1-score and Exact Match ratio (details of all evaluation metrics are in Appendix B.3.). On these metrics, R4P achieved 63.3% F1 and 41.8% EM, outperforming all baselines. We also compare R4P with trajectory verifiers and general reward models. For trajectory verifiers, we use their original prompts and apply controlled decoding to compare the probabilities of YES and NO for judging patch correctness. As shown in Fig. 2 (a), R4P largely outperforms them, as it provides a scaffold-ignorant verification, which is much more generalizable. For reward models, since they can only provide relative preference on answers, we construct a pair-wise subset from our evaluation data, where each instance contains one correct and one incorrect patch from the same issue. While this formulation is unoptimized for R4P, it still achieves a 69.8% accuracy, slightly outperforming the state-of-the-art, as shown in Fig. 2 (b). This result underscores that R4P’s core verification capabilities are robust and can be effectively generalized to other patch evaluation scenarios.

Practicality on RL To validate if R4P’s supervision could support RL scaling, we trained Mini-SE via *pure RL*. Notably, unlike normal SWE agents, Mini-SE is lightweight and does not test its patches during rollout. Thus, it naturally achieves a lower Pass@1 score in the absence of an

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

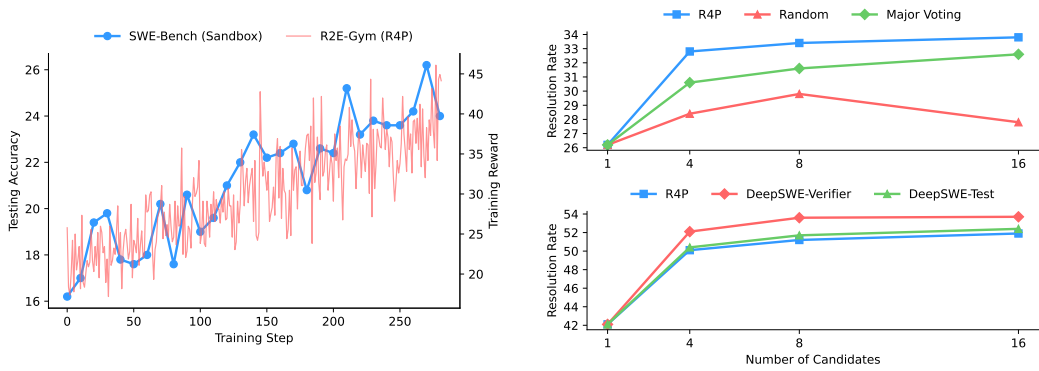


Figure 3: **Left:** Mini-SE’s training rewards and testing accuracy in RL process. **Upper Right:** Mini-SE’s resolution rates w.r.t. different TTS strategy. **Lower Right:** R4P’s performance on selecting DeepSWE’s patches. The scaling curves illustrates R4P’s practicality and generalizability.

extra test-time patch verifier. Nevertheless, as shown in Fig. 3, Mini-SE achieved a 26.2% Pass@1 within 300 steps, boosting the base Qwen3-32B model performance by 10.0%. Furthermore, the test accuracy increases steadily with training rewards, indicating that R4P provides stable supervision signals without leading to reward collapse. This feature enables R4P to further supervise agents on vast, untested issues in open-source community, illustrating its practicality for agentic RL training.

Practicality on TTS To validate R4P’s effectiveness for TTS, we sample 16 patches from Mini-SE on the SWE-bench-verified. We employed a two-round process: the patches are grouped for R4P verification, and those predicted as incorrect patches will be filtered out. Then we prompt R4P to output the single most likely correct patch from all remaining candidates as output. As shown in Fig. 3, Mini-SE’s resolution rate increases steadily with the number of candidate patches evaluated. When all 16 patches are considered, it reaches 33.8%. To further show its generalizability, we employ R4P on DeepSWE’s patches for selection. While both of the test agent and verifier model of DeepSWE are optimized for its own trajectory, R4P achieves a comparable results of the test agent.

Efficiency on patch verification Fig. 4 shows the distribution of the average time required to verify each patch per step when R4P is deployed on 2x80GB GPUs via vLLM. We compare it against the time taken to test the golden patches for 500 instances from swe-bench-verified using a CPU server with 64 cores and 128GB of RAM. The average time per instance for R4P does not exceed 1 second, significantly outperforming the average testing time per instance of 50 seconds. Furthermore, patches generated by LLMs often contain algorithmic inefficiencies or infinite loops. During the training of Mini-SE, we observed that 59.2% (29/49) of validation instances reached the 30-minute timeout of SWE-bench harness. This demonstrates the superior efficiency of R4P’s group-wise reasoning-based verification compared to testing-based verification.

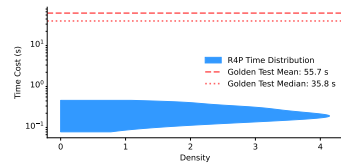


Figure 4: Time cost per patch.

5.3 ANALYSIS AND DISCUSSION

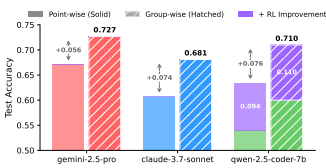


Figure 5: Task formulation

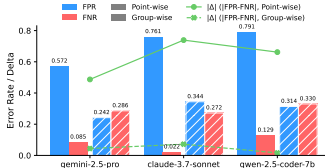


Figure 6: Bias of FPR/FNR

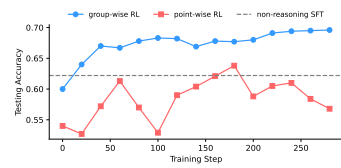


Figure 7: RL convergence

432 approximately 25% higher. This suggests that although some noisy supervision creates conflicting
433 gradients, the overall update direction remains correct. Given R4P’s ability to scale data efficiently,
434 this noisy yet stable supervision is an acceptable trade-off.
435

436 6 RELATED WORK 437

438 **SWE data scaling** Data scaling has become a central challenge of SWE. Typically, people mines
439 issues from GitHub and creating sandbox for testing (Pan et al., 2024), which is very labor-intensive.
440 To facilitate it, one way is LLM-assisted sandbox construction (Hu et al., 2025; Yang et al., 2025).
441 However, effort is still required for checking sandbox correctness (Hu et al., 2025). Another way
442 is issue generation and commit backtranslation to mutate numerous instances from several seed
443 instances (Jain et al., 2025; Yang et al., 2025). However, they often differ from real-world issues
444 and empirically lead to suboptimal training performance (Luo et al., 2025a). A third approach relies
445 on similarity metrics to supervising SWE sub-tasks separately Ma et al. (2025); Wei et al. (2025);
446 Xie et al. (2025). However, since correct patches are often non-unique, it can incorrectly penalize
447 novel yet valid solutions, introducing noise and reducing sample efficiency. In contrast to these data-
448 centric methods, we focus on expanding supervision rather than data, as a large number of diverse
449 untested issues remain under-exploited in the whole open-source community.
450

451 **Reward models and verifiers** Reward models are initially designed to predict an answer’s relative
452 quality (Minghao Yang, 2024; Chen et al., 2025) for AI-human preference alignment. With the
453 emergence of Reinforcement Learning with Verifiable Rewards (RLVR) and Monte Carlo-based
454 RL algorithms, e.g., GRPO (Guo et al., 2025; Shao et al., 2024), rule-based rewards have gained
455 prominence in reasoning tasks where absolute correctness can be determined. In such settings,
456 reward models are often repurposed as verifiers during test-time selection (Pan et al., 2024; Jain et al.,
457 2025; Luo et al., 2025a). However, SWE tasks lack easily accessible high-quality rule-based rewards
458 from testing. Thus, we explore using reward models with absolute patch verification capability to
459 provide scalable, deterministic supervision.
460

460 7 LIMITATION AND CONCLUSION 461

462 **Limitation** *First*, R4P’s weights remain fixed after training. As the agent’s policy improves, the
463 static reward model may become misaligned with true answer quality (Gao et al., 2023). In future
464 work, we plan to explore periodic calibration using issues with high-quality sandbox during extended
465 RL training, where discrepancies between R4P’s predictions and actual test results could guide
466 weight updates. *Second*, as a patch verifier, R4P are designed to supervise patch generation process
467 and cannot solely provide end-to-end supervision for complex agents that perform in-loop self-
468 verification via test generation and execution (Gao et al., 2025; Luo et al., 2025a). Future work
469 will explore mixed training strategies, e.g., scaling normal issues with R4P to train straightforward
470 patch generation while leveraging challenging issues with high-quality test sandbox to train patch
471 generation with in-loop verification. *Third*, R4P and Mini-SE are Python-centric. To generalize to
472 other languages, it is required to collect patches from their datasets and adapt the code search tool to
473 their static analyzer. *Fourth*, R4P’s performance correlates with multiple factors like group diversity
474 as discussed in Sec. 5.3. To mitigate it, it might be helpful to follow our best practice suggestions.
475

476 **Conclusion** This work explores a model-based approach to address the scalability bottleneck in
477 supervising software engineering agents due to reliance on testing. It introduces R4P, a reasoning-
478 based patch verifier. R4P utilizes group-wise training objective to effectively verify patches against
479 each other and mitigate instability during training its patch reasoning capability. Empirical results
480 demonstrate that R4P achieves 72.2% verification accuracy, outperforming strong reasoning models
481 like OpenAI o3. We further showcased its practical value by training Mini-SE, a fully test-free agent
482 that leverages R4P for both reinforcement learning and test-time scaling. Mini-SE achieves 26.2%
483 Pass@1 and 33.8% Best@16 on SWE-bench-verified, showing that reasoning-based supervision is
484 a viable and powerful alternative to traditional test suites.
485

486 ETHICS STATEMENT
487

488 This work adheres to the ICLR Code of Ethics. It uses only publicly available datasets and involves
489 no human subjects or sensitive personal data. We considered fairness, privacy, and security and
490 found no ethical concerns; the authors take full responsibility.
491

492 REPRODUCIBILITY STATEMENT
493

494 All datasets are public, and preprocessing details appear in the Appendix. We release
495 anonymized code and scripts—with configs, hyperparameters, and evaluation protocols—to
496 fully reproduce our experiments and results in [https://anonymous.4open.science/r/
497 pacth-verifier-4F53/README.md](https://anonymous.4open.science/r/pacth-verifier-4F53/README.md).
498

499 LLM USE DISCLOSURE
500

501 LLM was used only for language polishing; it did not assist with the research content, and the
502 authors are fully responsible for the manuscript.
503
504

505 REFERENCES
506

- 507 Tobias Baum, Olga Liskin, Kai Niklas, and Kurt Schneider. A faceted classification scheme for
508 change-based industrial code review processes. In *2016 IEEE International conference on soft-
509 ware quality, reliability and security (QRS)*, pp. 74–85. IEEE, 2016.
- 510 Xiusi Chen, Gaotang Li, Ziqi Wang, Bowen Jin, Cheng Qian, Yu Wang, Hongru Wang, Yu Zhang,
511 Denghui Zhang, Tong Zhang, et al. Rm-r1: Reward modeling as reasoning. *arXiv preprint
512 arXiv:2505.02387*, 2025.
- 513 Leo Gao, John Schulman, and Jacob Hilton. Scaling laws for reward model overoptimization. In
514 *International Conference on Machine Learning*, pp. 10835–10866. PMLR, 2023.
- 515 Pengfei Gao, Zhao Tian, Xiangxin Meng, Xinchun Wang, Ruida Hu, Yuanan Xiao, Yizhou Liu, Zhao
516 Zhang, Junjie Chen, Cuiyun Gao, et al. Trae agent: An llm-based agent for software engineering
517 with test-time scaling. *arXiv preprint arXiv:2507.23370*, 2025.
- 518 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,
519 Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms
520 via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- 521 Ruida Hu, Chao Peng, Xinchun Wang, Junjielong Xu, and Cuiyun Gao. Repo2run: Automated
522 building executable environment for code repository at scale. *arXiv preprint arXiv:2502.13681*,
523 2025.
- 524 Naman Jain, Jaskirat Singh, Manish Shetty, Liang Zheng, Koushik Sen, and Ion Stoica. R2e-gym:
525 Procedural environments and hybrid verifiers for scaling open-weights swe agents. *arXiv preprint
526 arXiv:2504.07164*, 2025.
- 527 Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik
528 Narasimhan. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint
529 arXiv:2310.06770*, 2023.
- 530 M. G. Kendall and B. Babington Smith. On the method of paired comparisons. *Biometrika*, 1940.
531 URL <https://www.jstor.org/stable/2332613>.
532
- 533 Michael Luo, Naman Jain, Jaskirat Singh, Sijun Tan, Ameen Patel, Qingyang Wu, Alpav Ariyak,
534 Colin Cai, Tarun Venkat, Shang Zhu, Ben Athiwaratkun, Manan Roongta, Ce Zhang, Li Erran Li,
535 Raluca Ada Popa, Koushik Sen, and Ion Stoica. DeepSWE: Training a fully open-sourced, state-
536 of-the-art coding agent by scaling RL. Blog, July 2025a. URL [https://www.together.
537 ai/blog/deepswe](https://www.together.ai/blog/deepswe). Accessed: [Insert Date of Access].
538
539

- 540 Michael Luo, Sijun Tan, Roy Huang, Ameen Patel, Alpay Ariyak, Qingyang Wu, Xiaoxiang Shi,
541 Rachel Xin, Colin Cai, Maurice Weber, Ce Zhang, Li Erran Li, Raluca Ada Popa, and Ion Stoica.
542 DeepCoder: A fully open-source 14b coder at o3-mini level. Blog, April 2025b. URL <https://www.together.ai/blog/deepcoder>. Accessed: [Insert Date of Access].
543
- 544 Yingwei Ma, Rongyu Cao, Yongchang Cao, Yue Zhang, Jue Chen, Yibo Liu, Yuchen Liu, Binhua
545 Li, Fei Huang, and Yongbin Li. Lingma swe-gpt: An open development-process-centric language
546 model for automated software improvement. *arXiv preprint arXiv:2411.00622*, 2024.
547
- 548 Zexiong Ma, Chao Peng, Pengfei Gao, Xiangxin Meng, Yanzhen Zou, and Bing Xie. Sorft: Issue
549 resolving with subtask-oriented reinforced fine-tuning. *arXiv preprint arXiv:2502.20127*, 2025.
550
- 551 Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri
552 Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. Self-refine: Iterative refinement
553 with self-feedback. *Advances in Neural Information Processing Systems*, 36:46534–46594, 2023.
- 554 Xiangxin Meng, Zexiong Ma, Pengfei Gao, and Chao Peng. An empirical study on llm-based agents
555 for automated bug fixing. *arXiv preprint arXiv:2411.10213*, 2024.
- 556 Xiaoyu Tan Minghao Yang, Chao Qu. Inf-orm-llama3.1-70b, 2024. URL [<https://huggingface.co/infly/INF-ORM-Llama3.1-70B>] (<https://huggingface.co/infly/INF-ORM-Llama3.1-70B>).
- 557
558
559
- 560 Audris Mockus, Nachiappan Nagappan, and Trung T Dinh-Trong. Test coverage and post-
561 verification defects: A multiple case study. In *2009 3rd international symposium on empirical
562 software engineering and measurement*, pp. 291–301. IEEE, 2009.
- 563 Jiayi Pan, Xingyao Wang, Graham Neubig, Navdeep Jaitly, Heng Ji, Alane Suhr, and Yizhe
564 Zhang. Training software engineering agents and verifiers with swe-gym. *arXiv preprint
565 arXiv:2412.21139*, 2024.
- 566
567
- 568 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
569 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 570 Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang,
571 Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathemati-
572 cal reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- 573
574
- 575 Tu Shiwen, Zhao Liang, Chris Yuhao Liu, Liang Zeng, and Yang Liu. Skywork critic model
576 series. <https://huggingface.co/Skywork>, September 2024. URL <https://huggingface.co/Skywork>.
- 577
578
- 579 Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling llm test-time compute optimally
580 can be more effective than scaling model parameters. *arXiv preprint arXiv:2408.03314*, 2024.
- 581
582
- 583 Richard S Sutton, Andrew G Barto, et al. *Reinforcement learning: An introduction*, volume 1. MIT
584 press Cambridge, 1998.
- 585
586
- 587 Maurice H ter Beek, Rod Chapman, Rance Cleaveland, Hubert Garavel, Rong Gu, Ivo Ter Horst,
588 Jeroen JA Keiren, Thierry Lecomte, Michael Leuschel, Kristin Yvonne Rozier, et al. Formal
589 methods in industry. *Formal Aspects of Computing*, 37(1):1–38, 2024.
- 590
591
- 592 Norbert Tihanyi, Tamas Bisztray, Mohamed Amine Ferrag, Bilel Cherif, Richard A Dubniczky,
593 Ridhi Jain, and Lucas C Cordeiro. Vulnerability detection: from formal verification to
594 large language models and hybrid approaches: a comprehensive overview. *arXiv preprint
595 arXiv:2503.10784*, 2025.
- 596
597
- 598 Xingyao Wang, Boxuan Li, Yufan Song, Frank F Xu, Xiangru Tang, Mingchen Zhuge, Jiayi Pan,
599 Yueqi Song, Bowen Li, Jaskirat Singh, et al. Openhands: An open platform for ai software
600 developers as generalist agents. *arXiv preprint arXiv:2407.16741*, 2024.
- 601
602
- 603 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdh-
604 ery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models.
605 *arXiv preprint arXiv:2203.11171*, 2022.

594 Yuxiang Wei, Olivier Duchenne, Jade Copet, Quentin Carbonneaux, Lingming Zhang, Daniel Fried,
595 Gabriel Synnaeve, Rishabh Singh, and Sida I Wang. Swe-rl: Advancing llm reasoning via rein-
596 forcement learning on open software evolution. *arXiv preprint arXiv:2502.18449*, 2025.
597

598 Chunqiu Steven Xia, Yinlin Deng, Soren Dunn, and Lingming Zhang. Agentless: Demystifying
599 llm-based software engineering agents. *arXiv preprint arXiv:2407.01489*, 2024.

600 Chengxing Xie, Bowen Li, Chang Gao, He Du, Wai Lam, Difan Zou, and Kai Chen. Swe-fixer:
601 Training open-source llms for effective and efficient github issue resolution. *arXiv preprint*
602 *arXiv:2501.05040*, 2025.
603

604 John Yang, Carlos E Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik Narasimhan,
605 and Ofir Press. Swe-agent: Agent-computer interfaces enable automated software engineering.
606 *Advances in Neural Information Processing Systems*, 37:50528–50652, 2024.

607 John Yang, Kilian Leret, Carlos E Jimenez, Alexander Wettig, Kabir Khandpur, Yanzhe Zhang,
608 Binyuan Hui, Ofir Press, Ludwig Schmidt, and Diyi Yang. Swe-smith: Scaling data for software
609 engineering agents. *arXiv preprint arXiv:2504.21798*, 2025.

610 Boxi Yu, Yuxuan Zhu, Pinjia He, and Daniel Kang. Utboost: Rigorous evaluation of coding agents
611 on swe-bench. *arXiv preprint arXiv:2506.09289*, 2025.
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

A PRELIMINARY STUDY OF TESTED ISSUES

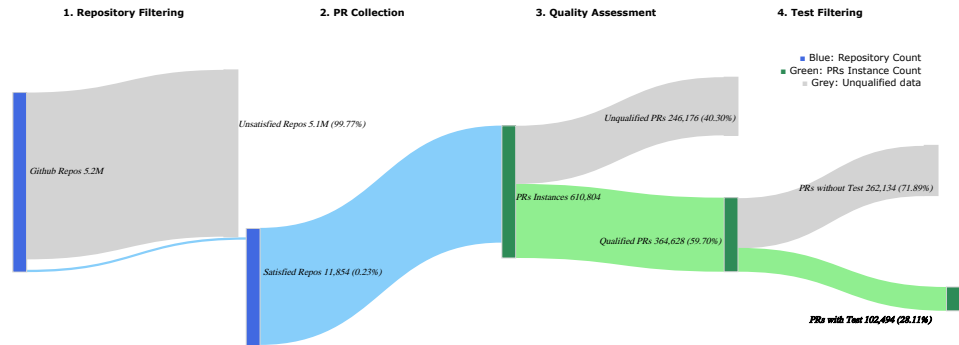


Figure 12: Visualization of data flow across four hierarchical steps in our GitHub issue curation pipeline: (1) *Repository Filtering*, (2) *PRs Collection*, (3) *Quality Assessment*, and (4) *Test Filtering*. The text in **Blue** indicates repository counts, the text in **Green** indicates issue–instance counts, and **Grey** flows denote data filtered out as unsatisfied or unqualified at each stage. Numbers and percentages shown in the diagram are computed within each step.

This study measures the prevalence of associated tests for issues and pull requests in open-source communities. We begin from the universe of 5.2 M public GitHub repositories and pre-filter to Python projects with at least 500 stars, excluding any repository already included in *SWE-Bench Verified*. The repository filtering leaves 11,854 eligible projects (0.23%). From these repositories, we harvest 610,804 PRs. We then apply a series of quality filters to the candidate issues: (i) the problem description must fall within 100–4,000 characters to align with the length distribution observed in *SWE-Bench Verified*; (ii) descriptions containing non-ASCII content (e.g., Chinese characters) are discarded; (iii) associated fixes must not modify non-Python files; (iv) patches must be small and self-contained, touching 1–5 files; and (v) issues embedding images in the description are removed to avoid non-textual ambiguity. Finally, in qualified PRs, we filter the patch without corresponding test changes, leaving only 102,494 qualified PRs with tests. As summarized in Fig. 12, this pipeline yields only **28.11%** PRs have corresponding tests.

B EXPERIMENT DETAILS

B.1 TRAINING R4P

Training settings We adopt GRPO algorithm for training R4P. Based on the insights from recent DAPO report, we set the clip ratio high to 0.28 and do not employ a KL loss in objective function. With a batch size of 128, we trained Qwen-2.5-Coder-Instruct for 225 steps. To accelerate the convergence process, R4P generates 16 samples and performs 16 gradient backward in each step. Throughout training, the temperature is maintained at 1.0 and the learning rate is fixed at 1e-6. To implement the training, we adopt VeRL framework with vLLM engine for rollout.

Dataset construction As mentioned in Sec. 3, we sampled 6 patches for each issue in the SWE-gym sandbox. During training data construction, the group size was set to 4. To improve the sample efficiency of the patches, we created multiple instances for each issue by forming different combinations of patches, i.e., constructing C_6^4 combinations for each issue. However, since the average solve rate of Claude-3.7 with OpenHands on SWE-gym is only about 30%, directly use such label imbalanced data will lead to significant under-estimate tendency of LLM. Thus, we balanced the combinations such that each group contained a roughly equal number of instances with correct patches ranging from 0 to 4—each category accounting for approximately 50% of the total instances. Ultimately, this process yielded 7,599 training samples with a group size of 4.

Prompt template We provide the system and user prompt in our R4P implementation, where the verification group size $N=4$:

System prompt

702 You are a software expert. You will be given a software issue and some
 703 patch candidates in user query. You need to judge which patch(es) can
 704 resolve the issue. Carefully review, critic, and compare the given
 705 candidates. You need to first think about the reasoning process in
 706 the mind until you get the final answer. Finally, put the ID(s) of
 707 correct patch candidates within `\boxed{}`, e.g., `\boxed{1}`, `\boxed{2,`
 708 `4}`, `\boxed{1, 2, 3, 4}` (all correct), `\boxed{}` (all wrong).

709 *User prompt*

```
710 <issue>
711 {problem_statement}
712 </issue>
713 <patch-1>
714 {generated_patch_1}
715 <patch-1>
716 <patch-2>
717 {generated_patch_2}
718 <patch-2>
719 <patch-3>
720 {generated_patch_3}
721 <patch-3>
722 <patch-4>
723 {generated_patch_4}
724 <patch-4>
```

725 B.2 TRAINING MINI-SE

726 **Training settings** We mainly follow the best practice provided by DeepSWE, including using
 727 leave-one-out and remove reward standard deviation for advantage estimation, filter out overlong
 728 trajectory without loss calculation, and normalize policy loss by sequence length. During training,
 729 we adopt fully on-policy training, with batch size of 64 and sampling 8 rollout per step. We also
 730 fix the temperature at 1.0 and the learning rate at 1e-6. To implement agentic RL, we adopt VeRL’s
 731 AgentLoop framework for training, and use vLLM for rollout. To accelerate the training process,
 732 we extract all instance’s code graph via tree-sitter in advance, which can be directly used for code
 733 search tool. We also checkout the copy of repositories in a separate directory for models to edit. All
 734 the code changes will be recorded in the final patches generated by git diff. In addition, we set a
 735 maximum limit of 50 rounds and 28K tokens during training. We deploy R4P model on a vLLM
 736 server for providing verification. [The reward for agent’s RL training is vanilla rule-based reward \(1](#)
 737 [for “correct” patch, 0 for “incorrect” patch\).](#) We finish the training process before 300 stpes.

738 **Prompt template** We provide the system prompt and tool configs of Mini-SE for reference. The
 739 user prompt is the SWE-bench or R2E-gym problem statement itself.

741 *System prompt*

```
742 You are an expert AI software engineering agent. Your primary goal is to
743 resolve a GitHub issue given in the user message. Following this
744 workflow methodically:
```

- 745 1. Understand the problem:
 - 746 - Thoroughly comprehend the issue description, identifying core
 - 747 components and expected behavior
 - 748 - Determine reproduction steps and failure conditions
- 749 2. Explore and Locate:
 - 750 - Use `'search_tool'` to explore the relevant files, entities, and test
 - 751 cases related to the bug report
 - 752 - Locate the exact root cause of the bug
- 753 3. Develop and Fix:
 - 754 - Develop minimal, targeted code modifications to address the root
 - 755 cause
 - Use `'edit_tool'` to apply surgical patch. Aim for minimal, clean
 - changes

- 756
757
758
759
760
761
762
763
764
765
766
767
768
4. Review and Revise:
 - Review the original reproduction steps to ensure the fix effectiveness
 - Review the relevant regression tests to avoid introducing any new bugs
 - Iterate using `search_tool` for review and `edit_tool` for revise until you confirm no edge cases are overlooked
 5. Submit the patch:
 - Call `patch_submission` tool to generate a unix diff patch and submit it to the user when confirming full resolution
 - Ensure the final patch is non-empty before finishing this conversation
 - All code changes persist throughout the conversation and will be included in the final patch

769 *Tool configs*

770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

```

type: "function"
function:
  name: "edit_tool"
  description: "A file edit tool that replaces an old string of text
    with a new string.\nNotes:\n1. The `old_str` parameter must match
    a segment of the file's content exactly. Pay close attention
    to whitespace, indentation, and newlines.\nThe edit will fail if `
    old_str` is not found, or if it is found multiple times in the
    file. Ensure `old_str` is unique enough to target the specific
    code block.\n3. `edit_tool` permanently modifies the actual
    repository** (changes persist to the final outcome)."
```

```

parameters:
  type: "object"
  properties:
    path:
      type: "string"
      description: "Relative file path e.g. `dir/file.py`."
    old_str:
      type: "string"
      description: "The exact string/content to be replaced in the file.
        This must be a unique match within the entire file."
    new_str:
      type: "string"
      description: "The new string that will replace `old_str`. Use an
        empty string to perform a deletion."
    required: ["path", "old_str", "new_str"]
```

```

type: "function"
function:
  name: "search_tool"
  description: "A code graph tool to show the source code of a function/
    class/class_method by its name.\nNotes:\n1. This tool is designed
    for Python code; it cannot find entities in other languages.\n2.
    If the `entity` name is not unique across the repository, the
    source code for all matching items will be displayed.\n3. This
    tool operates on a static, initial snapshot of the repository**.
    Changes made with `edit_tool` will not be visible in the
    search results."
```

```

parameters:
  type: "object"
  properties:
    construct:
      type: "string"
      description: "Type of program construct to search for."
      enum: ["function", "class", "class_method"]
    entity:
      type: "string"
```

```

810         description: "The simple name of the entity to find (e.g., `
811             my_function`, `MyClass`, `my_class_method`)."
812         required: ["construct", "entity"]
813
814     type: "function"
815     function:
816         name: "patch_submission"
817         description: "Automatically generate a diff patch for existing code
818             changes and submit it to the user."
819         parameters:
820             type: "object"
821             properties: {}
822             required: []

```

Note that we also provide a patch submission tool to enable a final check of LLM before ending the interaction and the real patch submission process.

B.3 EVALUATION DATA

Data source Our evaluation data for patch verification task is sampled from SWE-bench verified experiments, which include all the patch submissions of the agents on the leaderboard. Specifically, our training data is from the following submissions: (1) 20250520_openhands_devstral_small, (2) 20241029_openhands-codeact-2.1-sonnet-20241022, (3) 20240728_sweagent_gpt4o, (4) 20240620_sweagent_claude3.5sonnet, as they have a overall moderate accuracy for balancing patch labels. Note that we remove the empty patches and patches that has meaningless modifications (e.g., only comments or tests are chagned) to avoid model shortcut through these low-level features. We thereby formulate a final datasets with 1,340 patches from 335 issues, where the overall accuracy is 49.93% (669/1340).

Prompt for reward model We follow the official prompt tempalte of these pair-wise reward models:

System prompt

```

840     Please act as an impartial judge and evaluate the quality of the
841     responses provided by two AI assistants to the user question
842     displayed below. You should choose the assistant that follows the
843     user's instructions and answers the user's question better. Your
844     evaluation should consider factors such as the helpfulness,
845     relevance, accuracy, depth, creativity, and level of detail of
846     their responses. Begin your evaluation by comparing the two
847     responses and provide a short explanation. Avoid any position
848     biases and ensure that the order in which the responses were
849     presented does not influence your decision. Do not allow the
850     length of the responses to influence your evaluation. Do not favor
851     certain names of the assistants. Be as objective as possible.
852     After providing your explanation, output your final verdict by
853     strictly following this format: `"[A]"` if assistant A is
854     better, `"[B]"` if assistant B is better.

```

User prompt

```

855     [User Question]
856     {issue}
857
858     "[The Start of Assistant A's Answer]
859     {patch_1}
860     [The End of Assistant A's Answer]
861
862     [The Start of Assistant B's Answer]
863     {patch_2}
864     [The End of Assistant B's Answer]

```

Prompt for trajectory verifier We follow the official prompt template of these verifiers:

System prompt

You are an expert judge evaluating AI assistant interactions. Your task is to determine if the assistant successfully resolved the user’s request.

Key evaluation criteria:

1. Did the assistant complete the main task requested by the user?
2. Did the assistant handle all edge cases and requirements specified?
3. Were there any errors or issues in the final solution?
4. Did the assistant verify the solution works as intended?

Respond only with "`<judgement>YES</judgement>`" or "`<judgement>NO</judgement>`".

where the user prompt is specific agent’s trajectory.

Evaluation metrics If a dataset has N patches and M unique issues (where each issue corresponds to K patches), then these N patches can be grouped into M patch groups (each of size K , i.e., $N = M \times K$). Below are the detailed definitions:

- **Acc.:** Accuracy is a *patch-level* evaluation metric that quantifies the proportion of correctly verified patches out of the total number. Let y_i be the ground-truth label for patch p_i , and $y^*(p_i)$ be the model’s prediction of its correctness. Using the indicator function \mathbb{I} (which evaluates to 1 if the condition is true, and 0 otherwise), the accuracy is calculated as:

$$Acc = \frac{\sum_{i=1}^N \mathbb{I}_{y_i=y^*(p_i)}}{N} \quad (4)$$

- **F1** The F1 score is a *group-level* evaluation metric. In our implementation, we calculate it using the set-theoretic definition (i.e., the Dice coefficient). For a given group i :

- Let A_i be the set of truly correct patches.
- Let B_i be the set of patches predicted as correct by the model.

The F1 score for group i is calculated as:

$$F1_i = \frac{2|A_i \cap B_i|}{|A_i| + |B_i|} \quad (5)$$

This metric serves to measure the similarity between the predicted set (B_i) and the ground-truth set (A_i). Furthermore, it can provide an alternative view of a *retrieval task* i.e., *the task of "retrieving the correct patches from a set of candidate patches."* In this context, the F1 score represents the harmonic mean of Precision (P_i) and Recall (R_i):

$$F1_i = \frac{2P_iR_i}{P_i + R_i} = \frac{2TP_i}{2TP_i + FP_i + FN_i} = \frac{2|A_i \cap B_i|}{|A_i| + |B_i|} \quad (6)$$

Note that in our implementation, we calculate the overall average F1 score for each individual group:

$$F1 = \frac{\sum_{i=1}^M F1_i}{M} \quad (7)$$

To handle the boundary conditions where one or both sets might be empty, we consider:

- $F1_i = 1$ if both sets are empty, as the model correctly predicted that no patch is correct.
- $F1_i = 0$ if only one set is empty, as it is a completely incorrect prediction.

- **EM** Exact Match is a *group-level* metric that assesses the model’s ability to correctly verify all patches within an entire group simultaneously. It is formally defined as the proportion of patch groups where the predicted set of correct patches exactly matches the ground-truth set:

$$EM = \frac{\sum_{i=1}^M \mathbb{I}_{|A_i|=|B_i|}}{M} \quad (8)$$

918 C CASE STUDY

919
920 In this section, we provide a case study of `django-django-13810`. We release R4P’s reasoning
921 trajectory as well as those of proprietary models with accessible reasoning trajectory (Claude-3.7-
922 Sonnet, Claude-4.0-Sonnet) as a reference.

923
924 **Issue and patches** This issue is about Django’s *MiddlewareNotUsed* error. In the original code,
925 the handler variable is updated via `self.adapt_method_mode()` before the middleware instantiation
926 triggers the exception. Thus, even though the middleware is skipped, the handler remains perma-
927 nently modified (incorrectly adapted for subsequent middleware), leading to a sync/async mismatch.

928 To solve this issue, each patch has its own way:

- 929
- 930 • **Patch 1 & 4** (incorrect): Attempt to initialize the middleware *before* adapting the handler.
931 This fails because middleware constructors typically expect the handler to already be in the
932 specific mode (sync/async) required.
- 933 • **Patch 2** (correct): Assigns the adapted handler to a temporary variable. The main handler
934 is only updated in the *else* block of the *try/except* statement. This ensures the update only
935 occurs when no exception was raised.
- 936 • **Patch 3** (correct): Also uses a temporary variable, but performs the assignment inside the
937 *try* block immediately. Since execution halts immediately if an exception is raised, this also
938 prevents the previous bug.

939
940 In this case, R4P correctly identifies Patches 2 and 3 as correct patch, yet incorrectly consider Patch
941 4 as well. In contrast, the two Claude models only correctly identify either Patch 2 or 3.

942 Below are the raw content of the issue and patches, as well as the reasoning trajectory and corre-
943 sponding analysis of R4P, Claude-3.7-Sonnet, and Claude-4.0-Sonnet.

944 <issue>

945 MiddlewareNotUsed leaves undesired side effects when loading middleware
946 in ASGI context
947 Description

948 I experienced strange issues when working with ASGI, `django-debug-toolbar`
949 and my own small middleware. It was hard problem to debug, I
950 uploaded an example project here: [https://github.com/hbielenia/asgi-](https://github.com/hbielenia/asgi-djangotoolbar-bug)
951 `djangotoolbar-bug` (the name is misleading - I initially thought it’s
952 a bug with `django-debug-toolbar`).

953 The `SESSION_FILE_PATH` setting is intentionally broken to cause a 500
954 error. When starting the application and accessing `/admin` (any
955 location really, but I wanted to leave it at a minimum and didn’t add
956 any views) it gives `TypeError: object HttpResponse can't be used in`
957 `'await' expression`. Commenting out `asgi_djangotoolbar_bug.middleware`.
958 `DummyMiddleware` fixes the issue (in that I receive a 500
959 `ImproperlyConfigured` exception). I’m not sure about the overall role
960 of `django-debug-toolbar` here - removing it causes Daphne to return a
961 500 error page but without debug information and there’s no traceback
962 in console either. I decided to leave it since it helped me
963 approximate the causes of issue.

964 I notice that in [https://github.com/django/django/blob/3.1.4/django/core/](https://github.com/django/django/blob/3.1.4/django/core/handlers/base.py#L58)
965 `handlers/base.py#L58` while `MiddlewareNotUsed` causes the loop to skip
966 futher processing and go to next middleware, it does leave handler
967 variable overwritten with output of `self.adapt_method_mode()`. On next
968 pass, this handler is passed to next middleware instance,
969 disregarding all the previous checks for (lack of) async support.
970 This likely causes the middleware chain to be "poisoned" from this
971 point onwards, resulting in last middleware in response cycle to
972 return an `HttpResponse` as a synchronous middleware would, instead of
973 coroutine that is expected.

974 This is probably avoided by adding async support to my middleware, but
975 unless I’m missing something docs indicate it should work as it is.
976 It is my intention that it’s applied only on synchronous requests, so

```

972         I didn't make it async compatible on purpose. If it's intentional in
973         Django that every middleware needs to support async if the
974         application is run as ASGI app, the documentation should probably
975         state that clearly. Though it kinda defeats the purpose of having
976         async_capable = False flag in the first place.
977     </issue>
978
979     <patch-1>
980     diff --git a/django/core/handlers/base.py b/django/core/handlers/base.py
981     --- a/django/core/handlers/base.py
982     +++ b/django/core/handlers/base.py
983     @@ -51,11 +51,12 @@ class BaseHandler:
984         middleware_is_async = middleware_can_async
985         try:
986             # Adapt handler, if needed.
987             + # Adapt handler, if needed.
988             + mw_instance = middleware(handler)
989             handler = self.adapt_method_mode(
990                 middleware_is_async, handler, handler_is_async,
991                 debug=settings.DEBUG, name='middleware %s' %
992                 middleware_path,
993             )
994             - mw_instance = middleware(handler)
995             except MiddlewareNotUsed as exc:
996                 if settings.DEBUG:
997                     if str(exc):
998
999     </patch-1>
1000     <patch-2>
1001     diff --git a/django/core/handlers/base.py b/django/core/handlers/base.py
1002     --- a/django/core/handlers/base.py
1003     +++ b/django/core/handlers/base.py
1004     @@ -51,11 +51,11 @@ class BaseHandler:
1005         middleware_is_async = middleware_can_async
1006         try:
1007             # Adapt handler, if needed.
1008             - handler = self.adapt_method_mode(
1009             + adapted_handler = self.adapt_method_mode(
1010                 middleware_is_async, handler, handler_is_async,
1011                 debug=settings.DEBUG, name='middleware %s' %
1012                 middleware_path,
1013             )
1014             - mw_instance = middleware(handler)
1015             + mw_instance = middleware(adapted_handler)
1016             except MiddlewareNotUsed as exc:
1017                 if settings.DEBUG:
1018                     if str(exc):
1019
1020     @@ -63,6 +63,9 @@ class BaseHandler:
1021         else:
1022             logger.debug('MiddlewareNotUsed: %r', middleware_path)
1023             continue
1024         + else:
1025         + # Only update handler if middleware was not skipped
1026         + handler = adapted_handler
1027
1028         if mw_instance is None:
1029             raise ImproperlyConfigured(
1030
1031     </patch-2>
1032
1033     <patch-3>
1034     diff --git a/django/core/handlers/base.py b/django/core/handlers/base.py
1035     --- a/django/core/handlers/base.py
1036     +++ b/django/core/handlers/base.py
1037     @@ -50,12 +50,14 @@ class BaseHandler:
1038         else:
1039             middleware_is_async = middleware_can_async

```

```

1026         try:
1027     - # Adapt handler, if needed.
1028     - handler = self.adapt_method_mode(
1029     + # Only adapt handler if the middleware is used
1030     + adapted_handler = self.adapt_method_mode(
1031         middleware_is_async, handler, handler_is_async,
1032         debug=settings.DEBUG, name='middleware %s' %
1033         middleware_path,
1034     )
1035     - mw_instance = middleware(handler)
1036     + mw_instance = middleware(adapted_handler)
1037     + # If we get here, the middleware will be used, so update the handler
1038     + handler = adapted_handler
1039     except MiddlewareNotUsed as exc:
1040         if settings.DEBUG:
1041             if str(exc):
1042 </patch-3>
1043 <patch-4>
1044 diff --git a/django/core/handlers/base.py b/django/core/handlers/base.py
1045 --- a/django/core/handlers/base.py
1046 +++ b/django/core/handlers/base.py
1047 @@ -50,11 +50,7 @@ class BaseHandler:
1048     else:
1049         middleware_is_async = middleware_can_async
1050     try:
1051     - # Adapt handler, if needed.
1052     - handler = self.adapt_method_mode(
1053     - middleware_is_async, handler, handler_is_async,
1054     - debug=settings.DEBUG, name='middleware %s' % middleware_path,
1055     - )
1056     + # Create middleware instance first to check if it will be used
1057     + mw_instance = middleware(handler)
1058     except MiddlewareNotUsed as exc:
1059         if settings.DEBUG:
1060 @@ -64,6 +60,12 @@ class BaseHandler:
1061         logger.debug('MiddlewareNotUsed: %r', middleware_path)
1062         continue
1063
1064     + # Only adapt the handler if the middleware will be used
1065     + handler = self.adapt_method_mode(
1066     + middleware_is_async, handler, handler_is_async,
1067     + debug=settings.DEBUG, name='middleware %s' % middleware_path,
1068     + )
1069     +
1070     if mw_instance is None:
1071         raise ImproperlyConfigured(
1072             'Middleware factory %s returned None.' % middleware_path
1073 </patch-4>
1074
1075

```

1069 **Reasoning of R4P** R4P correctly finds that Patch 2 and Patch 3 achieve the fix with a temporary
1070 variable and conditionally updating the main handler. However, R4P hallucinates regarding Patch 4.
1071 It considers Patch 4 ensures the handler is correctly adapted before used by the middleware instance,
1072 but actually this handler is unadapted.

1073 Let's analyze each patch to determine which one(s) can resolve the issue
1074 in the 'MiddlewareNotUsed' handling in the ASGI context.

1075 **### Issue Analysis**

1076 The issue is that when 'MiddlewareNotUsed' is raised, the 'handler'
1077 variable is being overwritten with the output of 'self.
1078 adapt_method_mode()', which can cause issues in the middleware chain.
1079 Specifically, it can lead to a 'TypeError' when the last middleware
in the chain returns an 'HttpResponse' instead of a coroutine.

```
1080
1081 ### Patch Analysis
1082
1083 #### Patch-1
1084 - Adds a comment to explain the adaptation of the handler.
1085 - Moves the `mw_instance = middleware(handler)` line after the adaptation
1086   of the handler.
1087 - This patch ensures that the `handler` is adapted before being passed to
1088   the middleware instance.
1089 - This is a good change because it ensures that the `handler` is
1090   correctly adapted before being used by the middleware instance.
1091
1092 #### Patch-2
1093 - Adds a comment to explain the adaptation of the handler.
1094 - Moves the `mw_instance = middleware(handler)` line after the adaptation
1095   of the handler.
1096 - Adds an `else` block to update the `handler` only if the middleware is
1097   not skipped.
1098 - This patch ensures that the `handler` is correctly adapted before being
1099   passed to the middleware instance and only updates the `handler` if
1100   the middleware is not skipped.
1101 - This is a good change because it ensures that the `handler` is
1102   correctly adapted before being used by the middleware instance and
1103   prevents the `handler` from being updated if the middleware is
1104   skipped.
1105
1106 #### Patch-3
1107 - Adds a comment to explain the adaptation of the handler.
1108 - Adds a comment to explain that the handler is only adapted if the
1109   middleware is used.
1110 - Moves the `mw_instance = middleware(handler)` line after the adaptation
1111   of the handler.
1112 - Adds a comment to explain that the handler is updated only if the
1113   middleware is used.
1114 - This patch ensures that the `handler` is correctly adapted before being
1115   passed to the middleware instance and only updates the `handler` if
1116   the middleware is not skipped.
1117 - This is a good change because it ensures that the `handler` is
1118   correctly adapted before being used by the middleware instance and
1119   prevents the `handler` from being updated if the middleware is
1120   skipped.
1121
1122 #### Patch-4
1123 - Moves the `mw_instance = middleware(handler)` line before the
1124   adaptation of the handler.
1125 - Adds a comment to explain that the middleware instance is created first
1126   to check if it will be used.
1127 - Adds a comment to explain that the handler is only adapted if the
1128   middleware will be used.
1129 - This patch ensures that the `handler` is correctly adapted before being
1130   passed to the middleware instance and only updates the `handler` if
1131   the middleware is not skipped.
1132 - This is a good change because it ensures that the `handler` is
1133   correctly adapted before being used by the middleware instance and
1134   prevents the `handler` from being updated if the middleware is
1135   skipped.
1136
1137 ### Conclusion
1138 Patches 2, 3, and 4 are all good changes. They ensure that the `handler`
1139   is correctly adapted before being passed to the middleware instance
1140   and only updates the `handler` if the middleware is not skipped.
1141 Patch-1 is also a good change, but it doesn't add the `else` block to
1142   update the `handler` only if the middleware is not skipped.
1143
1144 Therefore, the correct patches are:
```

1134

1135 `\boxed{2, 3, 4}`

1136

1137

1138 **Reasoning of Claude-3.7-Sonnet** Claude-3.7 correctly identifies Patch but rejects Patch 3 due to
 1139 a fundamental misunderstanding of exception control flow. It states that Patch 3 “updates handler
 1140 before the except block, which means if *MiddlewareNotUsed* is raised, the change to handler will
 1141 be lost.” However, if the line `mw_instance = middleware(...)` raises an exception, the immediately
 1142 following line `handler = adapted_handler` is never executed. Thus, Patch 3 functionally equivalent
 1143 to Patch 2. The model hallucinated a side effect where the assignment occurs despite the exception
 1144 interrupting the flow.

1145 I need to analyze the issue and the proposed patches to determine which
 1146 one(s) correctly fix the problem.

1147
 1148 `## Understanding the Issue`

1149 The issue occurs when `'MiddlewareNotUsed'` is raised in an ASGI context.
 1150 Currently, the code adapts the handler first (via `'self.
 1151 adapt_method_mode()'`), and then creates the middleware instance. If `'
 1152 MiddlewareNotUsed'` is raised when creating the middleware instance,
 1153 the adapted handler still remains and is passed to the next
 1154 middleware, which causes type mismatches when mixing sync and async
 1155 middleware.

1155 The key problem is that `'handler'` is modified before we know if the
 1156 middleware will be used or not. If the middleware isn't used (raises
 1157 `'MiddlewareNotUsed'`), we've already changed `'handler'` to be sync or
 1158 async based on the middleware that ended up not being used.

1159 `## Patch Analysis`

1160
 1161 `### Patch 1`

```
1162 ```python
1163 # Adapt handler, if needed.
1164 # Adapt handler, if needed. (duplicate comment)
1165 mw_instance = middleware(handler)
1166 handler = self.adapt_method_mode(
1167     middleware_is_async, handler, handler_is_async,
1168     debug=settings.DEBUG, name='middleware %s' % middleware_path,
1169 )
1170 ```
```

1170 This patch switches the order of operations - it creates the middleware
 1171 instance first, and only then adapts the handler. This approach would
 1172 fail because the middleware constructor expects to receive a handler
 1173 that's compatible with its sync/async nature, but we're passing the
 1174 potentially incompatible handler directly.

1175 `### Patch 2`

```
1176 ```python
1177 adapted_handler = self.adapt_method_mode(
1178     middleware_is_async, handler, handler_is_async,
1179     debug=settings.DEBUG, name='middleware %s' % middleware_path,
1180 )
1181 mw_instance = middleware(adapted_handler)
1182 ...
1183 else:
1184     # Only update handler if middleware was not skipped
1185     handler = adapted_handler
1186 ```
```

1186 This patch computes the adapted handler but doesn't update the main `'
 1187 handler'` variable until we're sure the middleware isn't skipped. It
 passes the adapted handler to the middleware instance, and only if no

```

1188     `MiddlewareNotUsed` is raised does it update `handler` to be the
1189     adapted version. This correctly addresses the issue.
1190
1191     ### Patch 3
1192     ```python
1193     adapted_handler = self.adapt_method_mode(
1194         middleware_is_async, handler, handler_is_async,
1195         debug=settings.DEBUG, name='middleware %s' % middleware_path,
1196     )
1197     mw_instance = middleware(adapted_handler)
1198     # If we get here, the middleware will be used, so update the handler
1199     handler = adapted_handler
1200     ```
1201
1202     Similar to Patch 2, this uses a temporary variable for the adapted
1203     handler, but it updates `handler` before the `except` block, which
1204     means if `MiddlewareNotUsed` is raised, the change to `handler` will
1205     be lost. This doesn't fix the issue.
1206
1207     ### Patch 4
1208     ```python
1209     # Create middleware instance first to check if it will be used
1210     mw_instance = middleware(handler)
1211     ...
1212     # Only adapt the handler if the middleware will be used
1213     handler = self.adapt_method_mode(
1214         middleware_is_async, handler, handler_is_async,
1215         debug=settings.DEBUG, name='middleware %s' % middleware_path,
1216     )
1217     ```
1218
1219     This patch tries to create the middleware instance first using the
1220     unadapted handler, then only adapts the handler if the middleware
1221     will be used. However, as noted for Patch 1, this would likely fail
1222     because the middleware constructor needs a compatible handler.
1223
1224     ## Conclusion
1225
1226     After analyzing all patches, I determine that only Patch 2 correctly
1227     fixes the issue by:
1228     1. Creating a temporary adapted handler
1229     2. Using it to initialize the middleware
1230     3. Only updating the main handler variable if the middleware isn't
1231     skipped
1232
1233     The other patches either introduce new issues or don't fully address the
1234     original problem.
1235
1236     \boxed{2}
1237
1238     Reasoning of Claude-4.0-Sonnet Claude-4.0 correctly analyzes the logic of Patch 2 and 3 (but
1239     also incorrectly consider Patch 4 as correct at its initial analysis), and even explicitly acknowledges
1240     that Patch 2 “addresses the core issue” and “should work.” Then, in the final comparison of Patch
1241     2,3,4, it correctly filters out Patch 4. However, it also dismisses Patch 2 merely because the else block
1242     adds “extra nesting” and Patch 3 appeared “cleaner.” Consequently, it generated a false negative for
1243     Patch 2. Interestingly, it is aware of the final objective is to find all patch(es) at the very beginning of
1244     the reasoning, but soon forgot this goal after a very long thought, and tend to find the “best” patch.
1245     We consider it is a problem caused by long-term forgetting.
1246
1247     <think>
1248     Let me understand the issue first:
1249
1250     1. The problem occurs when using ASGI with middleware that raises `
1251     MiddlewareNotUsed` exception

```

```

1242 2. In the current code, when 'MiddlewareNotUsed' is raised, the 'handler'
1243     variable has already been modified by 'self.adapt_method_mode()' but
1244     the middleware is skipped
1245 3. This "poisoned" handler is then passed to the next middleware, causing
1246     async/sync mismatch issues
1247 4. The result is that synchronous middleware returns 'HttpResponse'
1248     instead of the expected coroutine
1249
1249 Let me examine the current code flow:
1250 ```python
1251 try:
1252     # Adapt handler, if needed.
1253     handler = self.adapt_method_mode(...) # This modifies handler
1254     mw_instance = middleware(handler) # This might raise MiddlewareNotUsed
1255 except MiddlewareNotUsed as exc:
1256     # handler is already modified at this point!
1257     continue
1258 ```
1259
1259 The problem is that 'handler' gets modified even when the middleware will
1260     be skipped due to 'MiddlewareNotUsed'.
1261
1261 Now let me analyze each patch:
1262
1262 **Patch 1:**
1263 ```python
1264 # Adapt handler, if needed.
1265 # Adapt handler, if needed. # Duplicate comment
1266 mw_instance = middleware(handler)
1267 handler = self.adapt_method_mode(...)
1268 ```
1269
1269 This moves the middleware instantiation before the handler adaptation,
1270     but it doesn't solve the core issue. If 'MiddlewareNotUsed' is raised
1271     during middleware instantiation, the handler won't be adapted at all
1272     , which could break the chain. This doesn't look right.
1273
1273 **Patch 2:**
1274 ```python
1275 adapted_handler = self.adapt_method_mode(...)
1276 mw_instance = middleware(adapted_handler)
1277 # ... exception handling ...
1278 else:
1279     # Only update handler if middleware was not skipped
1280     handler = adapted_handler
1281 ```
1282
1282 This creates a temporary 'adapted_handler', uses it to create the
1283     middleware instance, and only assigns it back to 'handler' if no
1284     exception was raised (in the 'else' block of the try/except). This
1285     looks good and addresses the core issue.
1286
1286 **Patch 3:**
1287 ```python
1288 adapted_handler = self.adapt_method_mode(...)
1289 mw_instance = middleware(adapted_handler)
1290 # If we get here, the middleware will be used, so update the handler
1291 handler = adapted_handler
1292 ```
1293
1293 This is similar to Patch 2, but updates the handler immediately after
1294     successful middleware instantiation, still within the try block. This
1295     also addresses the issue but is slightly different in timing.
1296
1296 **Patch 4:**
1297 ```python
1298 try:
1299     # Create middleware instance first to check if it will be used

```

```
1296     mw_instance = middleware(handler)
1297 except MiddlewareNotUsed as exc:
1298     # ... handle exception ...
1299     continue
1300
1301 # Only adapt the handler if the middleware will be used
1302 handler = self.adapt_method_mode(...)
1303
1304 This completely reverses the order - first creates the middleware
1305 instance with the original handler, and only adapts the handler if
1306 the middleware will be used. This could be problematic because the
1307 middleware might expect an adapted handler during instantiation.
1308
1309 Let me think about which approach is best:
1310
1311 1. **Patch 1** is clearly wrong - it moves the adaptation after
1312 middleware creation, which could break things.
1313
1314 2. **Patch 4** might have issues because middleware instantiation happens
1315 with the non-adapted handler, but then the adapted handler is used
1316 later. This could cause inconsistencies.
1317
1318 3. **Patches 2 and 3** both use the correct approach: adapt the handler
1319 first (to a temporary variable), try to create the middleware with
1320 the adapted handler, and only commit the adaptation if successful.
1321
1322 Between Patch 2 and Patch 3:
1323 - Patch 2 uses the `else` clause of try/except, which is more explicit
1324 about "only if no exception occurred"
1325 - Patch 3 updates the handler immediately after successful middleware
1326 creation, still in the try block
1327
1328 Both should work, but Patch 2 is slightly more explicit and follows the
1329 pattern of "only do this if the try block succeeded completely". The
1330 `else` clause in try/except is specifically designed for this use
1331 case.
1332
1333 However, let me double-check the logic. In both patches 2 and 3, if `
1334 MiddlewareNotUsed` is raised during `middleware(adapted_handler)`,
1335 the original `handler` remains unchanged, which is exactly what we
1336 want.
1337
1338 Patch 3 is actually more straightforward and easier to read - if we
1339 successfully create the middleware instance, we immediately update
1340 the handler. The else clause in Patch 2 adds an extra level of
1341 nesting that might not be necessary.
1342
1343 Actually, let me reconsider Patch 4. The issue might be more subtle. Let
1344 me think about when `MiddlewareNotUsed` is typically raised:
1345
1346 `MiddlewareNotUsed` is usually raised in the middleware's `__init__`
1347 method when it determines it shouldn't be used. At this point, the
1348 middleware examines the handler that's passed to it. If the handler
1349 adaptation is important for the middleware to make its decision, then
1350 Patch 4 could be problematic.
1351
1352 However, if the middleware's decision to raise `MiddlewareNotUsed` is
1353 independent of the handler's adaptation (which is likely the case in
1354 most scenarios), then Patch 4 could actually work.
1355
1356 But looking at the specific issue described, it seems like the problem is
1357 that the handler adaptation happens regardless of whether the
1358 middleware will be used, so the safer approach is Patches 2 or 3.
1359
```

1350 Between Patch 2 and 3, I think Patch 3 is cleaner and more direct. The
1351 comment also makes the intent very clear.
1352 </think>
1353
1354 Looking at this issue, I need to understand the core problem: when `MiddlewareNotUsed` is raised, the `handler` variable has already been
1355 modified by `self.adapt_method_mode()`, but the middleware is
1356 skipped. This "poisoned" handler then gets passed to subsequent
1357 middleware, causing async/sync compatibility issues.
1358
1359 Let me analyze each patch:
1360
1361 ****Patch 1****: This moves middleware instantiation before handler
1362 adaptation, but has issues:
1363 - The duplicate comment suggests it's incomplete
1364 - More importantly, it reverses the logical order - middleware typically
1365 expects an adapted handler
1366 - If `MiddlewareNotUsed` is raised, the handler won't be adapted at all,
1367 which could break the middleware chain
1368
1369 ****Patch 2****: This uses a temporary variable approach:
1370 - Creates `adapted_handler` first
1371 - Uses it for middleware instantiation
1372 - Only updates the main `handler` in the `else` clause if no exception
1373 occurred
1374 - This properly isolates the adaptation from the exception handling
1375
1376 ****Patch 3****: Similar to Patch 2 but simpler:
1377 - Creates `adapted_handler` and uses it for middleware instantiation
1378 - Updates `handler` immediately after successful middleware creation
1379 - The comment clearly explains the intent
1380 - More direct and easier to follow than Patch 2
1381
1382 ****Patch 4****: This completely reverses the order:
1383 - Creates middleware instance first with the original handler
1384 - Only adapts the handler if middleware will be used
1385 - This could be problematic because middleware might need the adapted
1386 handler during instantiation to make decisions
1387
1388 Patches 2 and 3 both solve the core issue correctly by preventing handler
1389 modification when middleware is skipped. Between them, Patch 3 is
1390 cleaner and more straightforward - it updates the handler immediately
1391 after successful middleware creation, with a clear comment
1392 explaining the logic.
1393
1394 Patch 4 risks breaking middleware that depends on receiving an adapted
1395 handler during instantiation, and Patch 1 has the wrong order and
1396 incomplete implementation.
1397
1398
1399
1400
1401
1402
1403