Anonymous authors

004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

023

025

026 027 028

029

Paper under double-blind review

ABSTRACT

As LLM-based agents exhibit exceptional capabilities in addressing complex problems, there is a growing focus on developing coding agents to tackle increasingly sophisticated tasks. Despite their promising performance, these coding agents often produce programs or modifications that contain runtime errors, which can cause code failures and are difficult for static analysis tools to detect. Enhancing the ability of coding agents to statically identify such errors could significantly improve their overall performance. In this work, we introduce *Execution-free* Runtime Error Detection for COding Agents (REDO), a method that integrates LLMs with static analysis tools to detect runtime errors for coding agents, without code execution. Additionally, we propose a benchmark task, SWE-Bench-Error-Detection (SWEDE), based on SWE-Bench (lite), to evaluate error detection in repository-level problems with complex external dependencies. Finally, through both quantitative and qualitative analyses across various error detection tasks, we demonstrate that REDO outperforms current state-of-the-art methods by achieving a 11.0% higher accuracy and 9.1% higher weighted F1 score; and provide insights into the advantages of incorporating LLMs for error detection.

1 INTRODUCTION

Large language models (LLMs) and LLM-031 based agents have exhibited significant potential in code generation, code editing, and 033 code evaluation. This progress has culmi-034 nated in the development of advanced LLMbased agents (hereafter referred to as *coding* agents) designed to address increasingly complex tasks. For example, SWE-Bench (Jimenez 037 et al., 2024a) presents a demanding benchmark comprising repository-level coding challenges. This benchmark requires coding agents to gen-040 erate a modification patch that solves a given 041 problem within a GitHub repository, based on 042 a problem statement expressed in natural lan-043 guage. To effectively navigate complex tasks 044 such as those posed by SWE-Bench, coding agents must demonstrate proficiency in the following core competencies: 1) comprehension 046 of the problem statement and retrieving relevant 047 code, 2) reasoning towards a functionally cor-048 rect solution, and 3) generation of programs free from runtime errors such as SyntaxError, AttributeError, or TypeError. 051



Figure 1: REDO outperforms current SOTA methods (Bieber et al., 2022) with respect to accuracy and weighted F1 (W.F1) on different tasks.

While the majority of coding agents across different tasks focus on enhancing comprehension, re trieval and reasoning capabilities, the systematic detection of runtime errors has received compara tively limited attention. However, ensuring that generated code is free from runtime errors is as criti-

054 cal as the aforementioned capabilities. For example, an AttributeError can cause the modified code 055 to fail, irrespective of the agent's comprehension and reasoning processes. Indeed, coding agents are 056 not immune to runtime errors. On SWE-Bench-lite (Jimenez et al., 2024b), the top six coding agents 057 as of August 2024 (CodeStory, Mentatbot, Marscode, Lingma, Droid, AutoCodeRover) produce, on 058 average, 1.8 SyntaxErrors, 25.8 TypeErrors, 5.2 NameErrors, and 11.3 AttributeErrors. Moreover, SWE-Bench (lite) is deliberately curated to include only straightforward problems, suggesting that the incidence of runtime errors could be substantially higher on the full SWE-Bench dataset. Ad-060 ditionally, detecting these runtime errors enables coding agents to rapidly iterate solutions, thereby 061 reducing both time and cost. 062

063 Recent coding agents do often incorporate modules for detecting runtime errors. For example, SWE-064 agent (Yang et al., 2024) and AutoCoderRover (Zhang et al., 2024) employ static analysis tools such as Pylint (Foundation, 2024c), while CodeR (Chen et al., 2024) utilizes dynamic analysis by 065 generating unit tests and executing the modified code. However, these methods have three critical 066 limitations. First, because runtime errors are often triggered by specific inputs that are typically 067 unpredictable from the code alone, static analysis approaches struggle to detect errors such as many 068 TypeErrors. Second, although dynamic analysis techniques may identify these runtime errors, they 069 require execution of the underlying code, which is problematic for three reasons. First, execution 070 ceases upon encountering the first runtime error, allowing only one error to be detected at a time, 071 thereby increasing the cost of detection. Secondly, dynamically configuring execution environments 072 for diverse project setups and testing frameworks poses significant challenges and might trigger 073 technical, legal or privacy concerns (Puddu et al., 2022; Lacoste & Lefebvre, 2023). Furthermore, 074 although these mechanisms are commonly employed as separate modules, their performance has 075 never been comprehensively evaluated. This lack of modular evaluation makes it difficult to assess the extent to which these error detection modules enhance the performance of coding agents. 076

077 In this work, we introduce an execution-free runtime error detection method, termed REDO, which 078 integrates static analysis tools, such as Pyflakes (Foundation, 2024b) or PyRight (pyright, 2024), 079 with a large language model (LLM). This approach extends the capabilities of static analysis tools to detect a broader range of errors without the need for code execution. The combination of these 081 tools is specifically designed to maximize the advantages of both tools, achieving a balanced tradeoff between reliability and the breadth of detectable errors. Similar to other works (Jimenez et al., 2024a; Bieber et al., 2022), we specifically focus on the runtime error detection in Python reposito-083 ries; but our method could be straightforwardly extended to other languages. Moreover, we present 084 a challenging and practical benchmark, SWE-Bench-Error-Detection (SWEDE), which is the first 085 repository-level error detection in the presence of complex external dependencies. Beyond SWEDE, 086 we perform a suite of experiments encompassing various tasks to further evaluate error detection al-087 gorithms. As shown in Figure 1, REDO significantly outperforms previous methods, obtaining 088 SOTA performance across diverse scenarios; and through qualitative analysis, we provide insights 089 into the underlying mechanisms of REDO.

090 091 092

093 094

2 RELATED WORK

2.1 LLM-BASED CODE GENERATION, CODING AGENTS, AND SWE-BENCH

095 LLMs (Ouyang et al., 2022; et al., 2024b; 2023b; 2024a; Anthropic, 2024) have been increasingly 096 leveraged for automatic code generation (Nijkamp et al., 2023b; et al., 2021a; Chai et al., 2023; 097 et al., 2024c; Nijkamp et al., 2023a; et al., 2023a; Gunasekar et al., 2023) and code repair (Xia 098 et al., 2023; Shypula et al., 2024; Gunasekar et al., 2023; Prenner & Robbes, 2023; Huang et al., 2023). The advent of LLM-based agents has further expanded the scope of problem-solving in 100 complex coding tasks, leading to the development of specialized coding agents. For example, SWE-101 Agent (Yang et al., 2024) is built on the ReAct framework (Yao et al., 2023) and incorporates a 102 custom Agent-Computer Interface (ACI), enhancing the agent's ability to interact with the envi-103 ronment effectively. Similarly, AutoCodeRover (Zhang et al., 2024) provides an integrated set of 104 search tools and includes *Pylint* within its toolkit, which serves to statically detect errors. Another 105 significant contribution is *CodeR* (Chen et al., 2024), a multi-agent framework that facilitates code editing by coordinating multiple agents through structured graphs. Within this framework, specific 106 agents like the *Reproducer* and *Verifier* are designed to generate unit tests and validate modified 107 implementations, respectively. These coding agents are frequently evaluated using the SWE-Bench



Figure 2: REDO employes a two-step process. When a repository-level modification is given, REDO first applies differential analysis to compare the original and modified implementations. If runtime errors are identified, REDO triggers an alert and rejects the patch. Conversely, if no errors are detected, the patch is forwarded to the LLM-based detection. The final determination regarding the patch's safety is made based on the detection provided by the LLM.

benchmark Jimenez et al. (2024a), which is composed of coding issues extracted from public GitHub
 repositories. This benchmark poses significant challenges by requiring coding agents to comprehend
 issues, localize relevant code segments, and produce correct modifications.

135 2.2 STATIC ANALYSIS TOOLS AND RUNTIME ERROR PREDICTION

136 Static analysis, a technique for examining computer programs without execution, is particularly 137 valuable in contexts where executing the program might lead to legal, privacy, or computational 138 concerns. Due to its non-executive nature, static analysis has found widespread application in error 139 detection (Zheng et al., 2006; Dillig et al., 2007; Chow et al., 2024), bug identification (Ayewah 140 et al., 2008; Mashhadi et al., 2024), and vulnerability discovery Charoenwet et al. (2024); Son-141 nekalb et al. (2023); Esposito et al. (2024); Chess & McGraw (2004); Livshits & Lam (2005); Evans 142 & Larochelle (2002). In the context of Python programming, various professional static analysis 143 tools have been developed to enhance code quality. For instance, *Pylint* (Foundation, 2024c) and *Pyflakes* (Foundation, 2024b) are designed to identify errors, while *Bandit* PyCOA (2024) focuses 144 on detecting common security vulnerabilities. Additionally, tools such as *PyRight* (pyright, 2024) 145 and MyPy (Foundation, 2024a) perform type checking, contributing to more robust software devel-146 opment. Although the integration of LLMs with static analysis is still in its infancy, some recent 147 studies have proposed combining these technologies to enhance bug detection in complex systems, 148 such as the Linux kernel (Li et al., 2024a), and to identify security vulnerabilities (Li et al., 2024b). 149 However, the exploration of LLMs in conjunction with static analysis for runtime error detection 150 remains limited. Additionally, Bieber et al. (2022) presents a notable effort in this domain by lever-151 aging Graph Neural Networks (GNNs) for predicting runtime errors, along with proposing a dataset, 152 referred to as STA, to evaluate their approach's efficacy.

153 154

155

134

3 EXECUTION-FREE RUNTIME ERROR DETECTION FOR CODING AGENTS

In this study, we introduce REDO, which serves to check the safety of modification patches. Here,
"Unsafe" instances are those that might crash due to runtime errors; and "Safe" instances are those that can be successfully run. REDO operates through a two-phase process: differential analysis and LLM-based detection.

161 The differential analysis component employs a static analysis tool, which provides a dependable method for detecting runtime errors. However, its detection capabilities are generally constrained to

SyntaxError, AttributeError, and NameError. To address this limitation, the LLM-based detection is incorporated to reason about the input contexts. It extends REDO's detection capabilities to errors such as TypeError and ValueError, which are typically beyond the scope of static analysis.

By integrating these mechanisms, REDO achieves a balanced trade-off between reliability and the breadth of error detection. An overview of REDO's architecture is illustrated in Figure 2.

167 168 169

170

166

3.1 DIFFERENTIAL ANALYSIS

Static analysis tools like Pyflakes and PyRight typically ensure detection through reliable methods,
such as verifying syntax correctness and maintaining data type consistency, making them essential
for identifying runtime errors in coding agents. In this study, we employ **PyRight** as our static
analysis tool since it is fast and lightweight; however, our framework is designed to be flexible,
allowing the integration of any static analysis tool.

Despite their utility, static analysis tools are affected by two challenges. First, they are prone to 177 generating false positives, where potential vulnerabilities are incorrectly flagged (Kang et al., 2022; 178 Kharkar et al., 2022; Murali et al., 2024). For example, when PyRight is applied to original python 179 scripts containing the modified functions, which do not contain runtime errors, it falsely classifies an 180 average of 267 instances, considering 89% of all testing instances as "Unsafe" across various coding 181 agents. To mitigate this issue, particularly in code edit tasks, we introduce the concept of *differential* 182 analysis. This method involves applying static analysis tools to both the original and modified 183 implementations separately. By comparing the errors detected in the original implementation (S_{Orig}) 184 with those in the modified implementation (S_{Mod}) , we can identify any new errors introduced by 185 the modifications. If new errors are detected in the modified implementation, the patch is flagged as "Unsafe". Differential analysis effectively refines static analysis tools to focus specifically on runtime errors introduced by modifications, thereby filtering out false positives from the original 187 implementation. Notably, this removes almost all positives induced by the original implementation. 188

The second challenge with static analysis tools is their inability to detect errors that are triggered under specific inputs. In dynamically typed languages like Python, variable data types are unknown before execution and also can change in response to different inputs. Since static analysis tools cannot reason about input contexts, they often fail to detect these errors, resulting in low recall in error detection performance. We show a concrete example in Figure 4. To address this challenge, we propose the LLM-based detection described in the next section.

195 196

3.2 LLM-BASED DETECTION

197 198

199 In comparison to static analysis tools, LLMs possess the ability to comprehend both the problem 200 statement and the modification patch. This capability allows them to reason about potential input 201 contexts and anticipate runtime errors that might be overlooked by static analysis tools. However, the reasoning process of LLMs is not always as reliable as the error detection mechanisms inherent 202 in static analysis tools. To harness the complementary strengths of both approaches, we restrict 203 the application of LLMs to instances deemed 'Safe' by static analysis tools. Accordingly, we have 204 designed the LLM prompt template to identify potential runtime errors that static analysis tools may 205 have missed. 206

Specifically, for each modification patch, we provide two additional pieces of information: the problem statement and the Python script containing the original version of the modified functions. The problem statement outlines the input contexts and describes the potential verification process for the modified implementation. The original implementation provides the running context of the modified functions, including the safe utilization of variables and functions. Subsequently, we prompt the LLM to enumerate potential runtime errors that may arise due to the modification and ultimately assess the safety of the patch. A detailed prompt template is provided in Appendix E.2.

Given the cost and latency associated with LLM calls, the LLM API is restricted to a single invocation per instance in our study. For the LLM, we utilize *Claude-SONNET-3-5*, with the temperature setting fixed at zero. A pseudo-code of whole REDO framework is given in Algorithm 1.

216 4 SWE-BENCH-ERROR-DETECTION (SWEDE)

217

218 The difficulty of runtime error detection can vary drastically among different coding problems. For 219 instance, on the task proposed by Bieber et al. (2022), only one python script is considered on 220 each data point. This task is practical as it resonates the competitive programming scenario where 221 one python script should contain all functionality; and external dependencies are simple. However, as coding agents become more powerful, additional challenging and practical scenarios should be 222 considered. In this work, we propose an repository-level error detection task, which is based on 223 SWE-Bench (lite) (Jimenez et al., 2024a) and its evaluation results using different coding agents. 224 We name this task as SWE-Bench-Error-Detection or SWEDE. 225

226 SWE-Bench (lite) is a popular benchmark dataset containing instances of repository-level coding 227 problems. Taking a GitHub repository and a problem statement as inputs, SWE-Bench (lite) asks coding agents to generate a modification patch to resolve the problem. SWEDE extends SWE-Bench 228 (lite) to include generated patches and evaluation logs from SWE-Bench leaderboard (Jimenez et al., 229 2024b); but with a focus on detecting runtime errors induced by generated patches, without execut-230 ing the code. 231

232 The task is **challenging** for two reasons. First, the modified scripts usually call other python scripts. 233 referred to as *external dependencies*, within the same repository. For instance, as shown in Table 1, 234 when only one directory level above the location where the modified file resides is considered, there already are many dependencies on average. When all files in the repository are considered, the 235 number of dependencies could become even more intimidating. Furthermore, the unit tests might 236 not directly interact with the modified files. These two factors make SWEDE challenging for error 237 detection algorithms as running contexts of variables and functions are harder to infer. The task is 238 practical because, when coding agents autonomously modify repositories, detecting runtime errors 239 early offers instrumental information; and can potentially reduce cost and time. 240

Table 1: Average External Dependencies in Parent Folder on SWEDE Across Coding Agents.

Method	CodeStory	Demo	Aider	Lingma	Droid	ACR
Average dependencies	5.56	4.56	6.07	5.75	5.61	5.98

As the problem is challenging, in this work, we focus on detecting if a modification patch will induce any runtime errors, e.g., SyntaxError, AttributeError, TypeError, etc. As a result, given a patch, we label it as **positive** if it contains an runtime error; and **negative** if the patch passes all unit tests or fails the unit tests only because of wrong functionality. We propose to measure the performance of error detection algorithms on SWEDE using precision, recall and F-1 score.

251 252 253

254

255

241

247

248

249

250

5 **EXPERIMENTS**

5.1 QUANTITATIVE RESULTS

256 SWE-Bench-Error-Detection (SWEDE). We first evaluate REDO on SWEDE task. We consider 257 six State-of-the-art (SOTA) coding agents on the SWE-Bench-lite leaderboard, including CodeStory, 258 Mentatbot, Marscode, Lingma, Droid, and AutoCodeRover (ACR) Zhang et al. (2024). We compare 259 REDO to several baselines. First, we include two widely used static analysis tools, namely Pyflakes 260 and PyRight. To eliminate false positive detection, we apply differential analysis (introduced in Section 3.1. Second, we include an LLM-only method, denoted as LLM. The LLM is prompted with 261 the same template introduced in Section 3.2. Lastly, to study how the choice of static analysis tool 262 affects REDO, we include a REDO framework with Pyflakes, named as REDO-Pyflakes. Due to the 263 inherent stochasticity in generation, even with a temperature setting of zero, we generate the LLM 264 responses three times and report the mean and standard deviation of the outcomes. 265

266 As shown in Table 2, static analysis tools (Pyflakes and PyRight) usually have higher precision scores, while LLM obtains better recall scores. This justifies our analysis on the difference between 267 static analysis tools and LLM (as in Section 3). Furthermore, confusion matrices in Figure 3 show 268 that the static analysis tool is more prudent in claiming unsafe patches, and therefore misses more 269 failed patches than LLM does. By appropriately combining these two kinds of tools, REDO makes

a good trade-off between reliability and detectable score. This better trade-off contributes to supe-rior F1 scores. Second, the performance of LLM is closer to REDO because REDO utilizes LLM on more than 70% instances. However, we note that despite the similarity, REDO makes less API calls because of the differential analysis step. We show confusion matrices using REDO and differ-ent coding agents in Figure 6. Also, as can be seen on the other dataset (STA) below, when static analysis tools performs better, REDO can significantly outperform the LLM-only baseline. Third, the performance of REDO-Pyflakes and REDO are close and outperform others, demonstrating the robustness of REDO against static analysis tool choices. We further report results using Anthropic Claude OPUS in Table 4 and SONNET-3.5 with temperature being 0.5 in Table 5. We can observe that REDO and REDO-Pyflakes consistently outperform static analysis tools and LLM, demonstrat-ing the robustness of our proposed method against LLM configurations.

Coding agent	Metric			Meth	od	
County agoin	Wieute	Pyflakes	PyRight	LLM	REDO-Pyflakes	RE
	Precision	32.1	43.7	34.0 0.2	33.3 _{0.2}	33.9
CodeStory	Recall	22.8	48.1	69.2 _{0.7}	$70.5_{-0.7}$	75.5
	F1	26.7	<u>45.8</u>	$45.6_{0.4}$	$\overline{45.3}_{0.3}$	46.8
	Precision	34.1	<u>33.9</u>	30.7 0.6	31.2 0.6	31.
Demo	Recall	34.9	45.3	61.2 _{1.8}	$70.2_{1.3}$	74.4
	F1	34.5	38.8	$40.9_{\ 1.0}$	43.2 0.8	43.8
	Precision	22.9	33.8	$28.7_{1.0}$	27.9 _{0.7}	29.
Marscode	Recall	16.4	37.3	$68.7_{1.5}$	$74.6_{1.5}$	77.0
	F1	19.1	35.5	$40.5_{\ 1.2}$	40.6 1.0	42.
	Precision	<u>45.8</u>	48.3	40.5 0.1	40.3 0.2	41.
Lingma	Recall	27.8	29.9	66.3 _{0.6}	71.5 _{0.6}	71.
	F1	34.6	36.9	$50.3_{0.2}$	<u>51.5_{0.3}</u>	52.3
	Precision	53.3	42.1	41.8 _{0.3}	40.1 _{0.3}	42.0
Droid	Recall	22.0	14.7	63.6 _{0.5}	$65.4_{0.5}$	68.2
	F1	31.2	21.8	$50.4_{0.4}$	$49.7_{0.4}$	52.0
	Precision	45.3	54.7	39.2 _{0.2}	38.3 0.2	40.0
ACR	Recall	23.5	34.3	65.7 _{1.0}	<u>69.3 _{0.6}</u>	73.9
	F1	31.0	42.2	49.1 _{0.3}	$49.3_{0.0}$	52.4

Table 2: Performance metrics by method and benchmark





Ensemble patch To demonstrate how the detection errors serve to enhance coding agents, we propose a preliminary algorithm to ensemble a *base* coding agent with an *auxiliary* agent, basing on the detection results from REDO. Specifically, given an instance in SWEDE, REDO first checks the safety of the patch p_{base} from the base agent. If p_{base} is safe, it will be accepted; otherwise, REDO will check the safety of the patch p_{aux} from the auxiliary agent. If p_{aux} is safe, it will be alternatively accepted; otherwise, the algorithm falls back to the base agent and accept p_{base} . Note that, since the top agents on the SWE-Bench leaderboard are proprietary, we are unable to run their methods to generate new patches. Instead, we utilize the patches reported in the SWE-Bench repository. A 324 pseudo-code of the ensemble algorithm is given in Algorithm 2. The quantitative results and analysis 325 can be found in Section D. 326

327 **STA dataset.** To assess the general applicability of REDO, we conduct evaluations on another 328 dataset: the balanced test set introduced by Bieber et al. (2022), hereafter referred to as STA. This dataset is derived from a code generation dataset CodeNet (et al., 2021b) and comprises approximately 27,000 Python submissions for competitive programming tasks. The balanced test set was 330 designed to ensure that the number of submissions without errors is approximately equal to those 331 containing runtime errors. The dataset categorizes 26 distinct error types, including the category 332 "no error." The detailed enumeration of errors and their corresponding indices are consistent with 333 those listed in the prompt template provided in Appendix E.3. Unlike the SWEDE dataset, each 334 submission in STA consists of a single Python script that processes input via standard input (stdin) 335 and performs its functionality without relying on external dependencies. In the context of STA, error 336 detection algorithms are tasked with predicting runtime error types based on the Python submission, 337 with or without the presence of input contexts. These input contexts describe the potential values 338 that stdin may take, providing clues regarding possible runtime errors. We refer to the scenario 339 involving input context as With Context, and the scenario without input context as Without Context.

340 Given that the task in STA differs from that in SWEDE, we adapt REDO using modified prompt 341 templates. The templates for both the scenarios with and without input contexts are provided in 342 Appendix E.3 and Appendix E.4, respectively. The performance of the algorithms is assessed us-343 ing three metrics: accuracy (Acc), weighted F1 score (W.F1), and weighted error F1 score (E.F1). 344 Weighted metrics are computed by calculating the F1 score for each class independently and av-345 eraged using a weight that depends on the number of true instances for each class; the E.F1 is 346 calculated exclusively for data points containing a runtime error. We compare our methods with 347 those reported by Bieber et al. (2022). To further investigate the contributions of individual components within REDO, we also evaluate PyRight and LLM separately. Since the original dataset 348 samples do not include error-free submissions, we employed three different random seeds to sample 349 an equivalent number of submissions for a fair comparison. The means are reported in Table 3, and 350 the standard deviations are presented in Table 6. 351

352 First, as demonstrated in Table 3, REDO-PyRight achieves state-of-the-art (SOTA) performance 353 across all six evaluated metrics, underscoring the superior performance of REDO in the context of STA. Additionally, we observed discrepancies between the annotated error types and the results 354 obtained from our running results. Consequently, we also present evaluation results, enclosed in 355 brackets, based on error types identified in our runs. Under these conditions, REDO exhibits even 356 more pronounced improvements over the baselines. Second, our ablation PyRight obtains decent 357 performance in both Without and With Context tasks. Considering that the input to STA submissions 358 are usually most common cases, as opposed to corner cases in SWEDE, our result demonstrates that 359 static analysis tools could perform well if the inputs and dependencies are simple. Furthermore, 360 when comparing REDO to PyRight, we can also see that REDO obtains extra benefits by including 361 the LLM based tool. This shows the effectiveness of our proposed framework.

362 Moreover, the differing conclusions regarding static analysis tools and LLMs between SWEDE 363 and STA highlight that each excels in distinct tasks. By integrating these two tools, we can 364 leverage the strengths of both, rendering our framework more versatile and effective across a broader 365 range of tasks. 366

5.2 QUALITATIVE ANALYSIS 368

367

371

372 373

374

369 In this section, we qualitatively analyze how key design factors in REDO affects the performance. 370 Consequently, we target answering the following two questions:

- 1. When does LLM succeed and fail to help REDO?
- 2. How could the detected errors help?

375 When Does the LLM Succeed? Figure 4 illustrates a case in which the LLM successfully identifies a runtime error that PyRight overlooked. This example is on the *django-13551* instance using 376 coding agent *Lingma*. The complete LLM response and modification patch are presented in Fig-377 ure 8. Specifically, the evaluation log in Figure 4 reveals that an AttributeError is triggered when the Table 3: Quantitative results on detecting runtime errors in CodeNet dataset. Each method is eval-uated with and without the input context. Three metrics are included, namely accuracy (Acc), weighted F1 score (W. F1), and weighted error F1 score (E. F1). Bold values indicate the best performance for each metric, while underlined values represent the second-best performance.

383		W	ithout Conte	xt		With Context	t
384	Method	Acc	W. F1	E. F1	Acc	W. F1	E. F1
385	GGNN	62.8	58.9	45.8	68.3	66.5	56.8
380	TRANSFORMER	63.6	60.4	48.1	67.3	65.1	54.7
387	LSTM	66.1	61.4	48.4	68.1	66.8	58.3
388	IPA-GNN	68.3	64.8	53.8	71.4	70.1	62.2
389	E. IPS-GNN	68.7	64.9	53.3	71.6	70.9	<u>63.5</u>
390	PyRight	74.5 (76.4)	67.2 (70.1)	52.2 (54.2)	74.5 (76.4)	67.2 (70.1)	52.2 (54.2)
391	LLM	61.8 (64.1)	54.0 (57.2)	32.7 (35.8)	64.3 (66.8)	58.6 (62.1)	40.4 (43.7)
392	REDO	75.1 (76.8)	69.9 (72.6)	57.0 (59.1)	77.4 (79.3)	74.5 (77.5)	64.3 (66.6)

394	
395	LLM Prediction
396 397	I'll analyze the modification patch for potential runtime errors:
398	<runtime errors=""></runtime>
399	[AttributeError]: The modification adds `user.email` to the hash value calculation. However,
400	used that doesn't include an email field, this could raise an AttributeError.
401	
402	<conclusion> Unsafe </conclusion>
403 404	Evaluation log
405 406 407 408 409 410 411 412 413	<pre>Traceback (most recent call last): File "/opt/pyenv/versions/3.6.15/lib/python3.6/unittest/case.py", line 59, in testPartExecutor yield File "/opt/pyenv/versions/3.6.15/lib/python3.6/unittest/case.py", line 523, in subTest yield File "/opt/django_django/tests/auth_tests/test_tokens.py", line 57, in test_token_with_different_email tk1 = p0.make_token(user) File "/opt/django_django/django/contrib/auth/tokens.py", line 28, in make_token return selfmake_token_with_timestamp(user, selfnum_seconds(selfnow())) File "/opt/django_django/django/contrib/auth/tokens.py", line 70, in _make_token_with_timestamp selfmake_hash_value(user, timestamp), File "/opt/django_django/django/contrib/auth/tokens.py", line 98, in _make_hash_value return str(user.pk) + user.password + str(login_timestamp) + user.email + str(timestamp) AttributeError: 'CustomEmailEield' object bas no attribute 'email'</pre>
414	

Figure 4: Successful example

invalid attribute *email* is accessed on the variable *user*, which is an instance of a *CustomEmailField* object. PyRight failed to detect this runtime error due to the inability to infer the data type of *user* through static analysis. In contrast, the LLM, capable of reasoning about potential runtime contexts, successfully identifies the runtime error, thus marking the patch as Unsafe. This example demonstrates the advantage of leveraging LLMs to anticipate runtime errors. Another successful instance is documented in Appendix F.1.

When Does the LLM Fail? Figure 5 presents an instance where both PyRight and the LLM failed to detect a runtime error. This example pertains to the instance django-11797 using coding agent CodeStory. The complete modification patch and the LLM's response are detailed in Figure 11. The runtime error occurs when an invalid column name is referenced in a query. Given the current input context, it is nearly impossible to infer the content of the query, leading to the LLM's failure to predict this runtime error. This example highlights the limitations of REDO in inferring potential inputs with limited contextual information, a challenge we plan to address in future work. Another instance of failure is documented in Appendix F.2.

 LLM Prediction

 I Ve analyzed the modification patch and the original implementation. Here's my assessment:

 Runtime Errors>

 No potential runtime errors detected.

 Conclusion> Safe </Conclusion>

 Evaluation log

 Traceback (most recent call last):

 File "/opt/django_django/db/backends/utils.py", line 86, in _execute

 return self.cursor.execute(sql, params)

 File "/opt/django_django/db/backends/sqlite3/base.py", line 396, in execute

 return Database.Cursor.execute(self, query, params)

 sqlite3.0perationalError: no such column: U0.name

Figure 5: Failed example

How could the detected errors help? This section presents a qualitative example illustrating how detected errors can facilitate the correction of flawed patches. We introduce a preliminary patch-fixing algorithm that leverages error messages generated by REDO. Specifically, when presented with a modified patch (referred to as the *original generated patch*) and its associated detected errors, our approach first applies the patch and extracts the modified code chunks, with ten lines of code added before and after each chunk. We then prompt the LLM to refine the code chunks according to the identified errors. Finally, a fixed patch is generated based on the modifications to the code chunks. The detailed LLM prompt is provided in Appendix E.6. For instance, in the case of *Django-12308* utilizing *ACR*, REDO initially identifies a risky attribute call that could potentially lead to an AttributeError. Based on this detection, the patch-fixing algorithm generates an fixed patch that avoids invoking the risky attribute, thereby preventing the AttributeError present in the original patch. A comprehensive description of the algorithm and a detailed analysis of this example are provided in Section G.

6 CONCLUSION

In this study, we first present REDO, an innovative error detection framework that operates through a two-step process: differential analysis followed by LLM-based detection. This approach achieves a balanced trade-off between reliability and the scope of detectable errors. We also propose SWE-Bench-Error-Detection (SWEDE), a novel and challenging runtime error detection task that aligns with the increasing deployment of autonomous coding agents responsible for repository-level modifications. Furthermore, we conduct a comprehensive set of quantitative experiments to empirically demonstrate the efficacy of REDO across various tasks. In addition, our qualitative analysis offers insights into the conditions under which LLM integration proves beneficial or falls short, and how detected runtime errors using REDO could help fix the previous flawed patches.

7 LIMITATION AND FUTURE WORK

The current implementation of the LLM-based detection step constrains the number of LLM API
call on each data point to just one. Expanding the number of API calls, potentially by leveraging
agentic AI techniques, could significantly improve error detection capabilities. Additionally, the
SWEDE task is presently confined to a binary classification of 'Safe' and 'Unsafe.' As error detection algorithms become more sophisticated, it is crucial to consider expanding the classification
schema to encompass a wider spectrum of error types, similar to those addressed in STA. Finally, the
advancement of more refined ensemble and patch-fixing algorithms holds the potential to enhance the efficacy of REDO in supporting coding agents.

486 ETHICS STATEMENT.

487 488

491

492 493

494

501

502

522

523

524

539

Our work leverages results from previous methods, including publicly available sources and the 489 SWE-Bench dataset and leaderboard, as cited in the Experiment section. To the best of our knowl-490 edge, this study does not pose any risks related to harmful insights, discrimination, bias, fairness, privacy, or security concerns.

REPRODUCIBILITY STATEMENT.

495 We provide details of our experiments and implementation in Sections 3 and 5, including the models 496 used and a description of the data processing steps. Given the limited time, we are unable to wrap 497 up all the code files before the submission deadline. However, we will gladly provide them during 498 the rebuttal stage if required. Additionally, for each experiment, we used three random seeds and 499 report both the means and standard deviations. 500

REFERENCES

- 503 Anthropic. Meet claude. https://https://www.anthropic.com/claude, 2024. Accessed: 2024-08-12. 504
- 505 Nathaniel Ayewah, William Pugh, David Hovemeyer, J David Morgenthaler, and John Penix. Using 506 static analysis to find bugs. *IEEE software*, 25(5):22–29, 2008. 507
- David Bieber, Rishab Goel, Daniel Zheng, Hugo Larochelle, and Daniel Tarlow. Static prediction of 508 runtime errors by learning to execute programs with external resource descriptions, 2022. URL 509 https://arxiv.org/abs/2203.03771. 510
- 511 Yekun Chai, Shuohuan Wang, Chao Pang, Yu Sun, Hao Tian, and Hua Wu. Ernie-code: Beyond 512 english-centric cross-lingual pretraining for programming languages, 2023. URL https://arxiv. 513 org/abs/2212.06742.
- 514 Wachiraphan Charoenwet, Patanamon Thongtanunam, Van-Thuan Pham, and Christoph Treude. An 515 empirical study of static analysis tools for secure code review, 2024. URL https://arxiv.org/ 516 abs/2407.12241. 517
- 518 Dong Chen, Shaoxin Lin, Muhan Zeng, Daoguang Zan, Jian-Gang Wang, Anton Cheshkov, Jun Sun, 519 Hao Yu, Guoliang Dong, Artem Aliev, Jie Wang, Xiao Cheng, Guangtai Liang, Yuchi Ma, Pan 520 Bian, Tao Xie, and Qianxiang Wang. Coder: Issue resolving with multi-agent and task graphs, 2024. URL https://arxiv.org/abs/2406.01304. 521
 - B. Chess and G. McGraw. Static analysis for security. IEEE Security & Privacy, 2(6):76–79, 2004. doi: 10.1109/MSP.2004.111.
- Yiu Wai Chow, Luca Di Grazia, and Michael Pradel. Pyty: Repairing static type errors in python. In 525 Proceedings of the IEEE/ACM 46th International Conference on Software Engineering, pp. 1–13, 526 2024. 527
- 528 Isil Dillig, Thomas Dillig, and Alex Aiken. Static error detection using semantic inconsistency 529 inference. In Proceedings of the 28th ACM SIGPLAN Conference on Programming Language 530 Design and Implementation, pp. 435–445, 2007.
- 531 Matteo Esposito, Valentina Falaschi, and Davide Falessi. An extensive comparison of static appli-532 cation security testing tools, 2024. URL https://arxiv.org/abs/2403.09219. 533
- 534 Chen et al. Evaluating large language models trained on code, 2021a. URL https://arxiv.org/ 535 abs/2107.03374.
- 536 Dubey et al. The llama 3 herd of models, 2024a. URL https://arxiv.org/abs/2407.21783. 537
- 538 Li et al. Starcoder: may the source be with you!, 2023a. URL https://arxiv.org/abs/2305.06161.
 - OpenAI et al. Gpt-4 technical report, 2024b. URL https://arxiv.org/abs/2303.08774.

540 Puri et al. Codenet: A large-scale ai for code dataset for learning a diversity of coding tasks, 2021b. 541 URL https://arxiv.org/abs/2105.12655. 542 Rozière et al. Code llama: Open foundation models for code, 2024c. URL https://arxiv.org/abs/ 543 2308.12950. 544 Touvron et al. Llama 2: Open foundation and fine-tuned chat models, 2023b. URL https://arxiv. 546 org/abs/2307.09288. 547 David Evans and David Larochelle. Improving security using extensible lightweight static analysis. 548 *IEEE software*, 19(1):42–51, 2002. 549 550 Python Software Foundation. Mypy. https://pypi.org/project/mypy/, 2024a. 551 Python Software Foundation. Pyflakes. https://pypi.org/project/pyflakes/, 2024b. 552 553 Python Software Foundation. Pylint. https://pypi.org/project/pylint/, 2024c. 554 Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio César Teodoro Mendes, Allie Del Giorno, Sivakanth 555 Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo de Rosa, Olli Saarikivi, Adil Salim, Shital 556 Shah, Harkirat Singh Behl, Xin Wang, Sébastien Bubeck, Ronen Eldan, Adam Tauman Kalai, Yin Tat Lee, and Yuanzhi Li. Textbooks are all you need, 2023. URL https://arxiv.org/abs/ 558 2306.11644. 559 Kai Huang, Xiangxin Meng, Jian Zhang, Yang Liu, Wenjie Wang, Shuhao Li, and Yuqing Zhang. 561 An empirical study on fine-tuning large language models of code for automated program repair. 562 In 2023 38th IEEE/ACM International Conference on Automated Software Engineering (ASE), 563 pp. 1162-1174, 2023. doi: 10.1109/ASE56229.2023.00181. 564 Carlos E. Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik 565 Narasimhan. Swe-bench: Can language models resolve real-world github issues?, 2024a. URL 566 https://arxiv.org/abs/2310.06770. 567 568 Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R Narasimhan. Swe-bench leaderboard, 2024b. URL https://www.swebench.com/. Accessed: 569 2024-08-19. 570 571 Hong Jin Kang, Khai Loong Aw, and David Lo. Detecting false alarms from automatic static 572 analysis tools: how far are we? In Proceedings of the 44th International Conference on 573 Software Engineering, ICSE '22. ACM, May 2022. doi: 10.1145/3510003.3510214. URL 574 http://dx.doi.org/10.1145/3510003.3510214. 575 Anant Kharkar, Roshanak Zilouchian Moghaddam, Matthew Jin, Xiaoyu Liu, Xin Shi, Colin 576 Clement, and Neel Sundaresan. Learning to reduce false positives in analytic bug detectors. In 577 Proceedings of the 44th International Conference on Software Engineering, ICSE '22. ACM, May 578 2022. doi: 10.1145/3510003.3510153. URL http://dx.doi.org/10.1145/3510003.3510153. 579 580 Marc Lacoste and Vincent Lefebvre. Trusted execution environments for telecoms: Strengths, 581 weaknesses, opportunities, and threats. *IEEE Security & Privacy*, 21:37–46, 2023. URL https://api.semanticscholar.org/CorpusID:258069655. 582 583 Haonan Li, Yu Hao, Yizhuo Zhai, and Zhiyun Qian. Enhancing static analysis for practical bug 584 detection: An Ilm-integrated approach. Proc. ACM Program. Lang., 8(OOPSLA1), apr 2024a. 585 doi: 10.1145/3649828. URL https://doi.org/10.1145/3649828. 586 Ziyang Li, Saikat Dutta, and Mayur Naik. Llm-assisted static analysis for detecting security vulner-587 abilities. arXiv preprint arXiv:2405.17238, 2024b. 588 589 V Benjamin Livshits and Monica S Lam. Finding security vulnerabilities in java applications with 590 static analysis. In USENIX security symposium, volume 14, pp. 18–18, 2005. 591 Ehsan Mashhadi, Shaiful Chowdhury, Somayeh Modaberi, Hadi Hemmati, and Gias Uddin. An 592 empirical study on bug severity estimation using source code metrics and static analysis. Journal of Systems and Software, pp. 112179, 2024.

- 594 Aniruddhan Murali, Noble Mathews, Mahmoud Alfadel, Meiyappan Nagappan, and Meng Xu. Fuz-595 zslice: Pruning false positives in static analysis warnings through function-level fuzzing. In 596 Proceedings of the IEEE/ACM 46th International Conference on Software Engineering, ICSE 597 '24. ACM, February 2024. doi: 10.1145/3597503.3623321. URL http://dx.doi.org/10.1145/ 598 3597503.3623321.
- Erik Nijkamp, Hiroaki Hayashi, Caiming Xiong, Silvio Savarese, and Yingbo Zhou. Codegen2: 600 Lessons for training llms on programming and natural languages, 2023a. URL https://arxiv. org/abs/2305.02309. 602
- 603 Erik Nijkamp, Bo Pang, Hiroaki Hayashi, Lifu Tu, Huan Wang, Yingbo Zhou, Silvio Savarese, and Caiming Xiong. Codegen: An open large language model for code with multi-turn program 604 synthesis, 2023b. URL https://arxiv.org/abs/2203.13474. 605
- 606 Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong 607 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kel-608 ton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, 609 and Ryan Lowe. Training language models to follow instructions with human feedback, 2022. 610 URL https://arxiv.org/abs/2203.02155.
- Julian Aron Prenner and Romain Robbes. Out of context: How important is local context in neural 612 program repair?, 2023. URL https://arxiv.org/abs/2312.04986. 613
- 614 Ivan Puddu, Moritz Schneider, Daniele Lain, Stefano Boschetto, and Srdjan Čapkun. On (the lack 615 of) code confidentiality in trusted execution environments, 2022. URL https://arxiv.org/abs/ 616 2212.07899.
- 617 PyCQA. bandit. https://https://github.com/PyCQA/bandit, 2024. 618
- 619 pyright. Microsoft. https://https://github.com/microsoft/pyright, 2024. 620
- Alexander Shypula, Aman Madaan, Yimeng Zeng, Uri Alon, Jacob Gardner, Milad Hashemi, Gra-621 ham Neubig, Parthasarathy Ranganathan, Osbert Bastani, and Amir Yazdanbakhsh. Learning 622 performance-improving code edits, 2024. URL https://arxiv.org/abs/2302.07867. 623
- 624 Tim Sonnekalb, Christopher-Tobias Knaust, Bernd Gruner, Clemens-Alexander Brust, Lynn von 625 Kurnatowski, Andreas Schreiber, Thomas S. Heinze, and Patrick Mäder. A static analysis plat-626 form for investigating security trends in repositories, 2023. URL https://arxiv.org/abs/2304. 627 01725.
- Chunqiu Steven Xia, Yuxiang Wei, and Lingming Zhang. Automated program repair in the era 629 of large pre-trained language models. In 2023 IEEE/ACM 45th International Conference on 630 Software Engineering (ICSE), pp. 1482–1494, 2023. doi: 10.1109/ICSE48619.2023.00129. 631
- 632 John Yang, Carlos E. Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik Narasimhan, and Ofir Press. Swe-agent: Agent-computer interfaces enable automated software engineering, 633 2024. URL https://arxiv.org/abs/2405.15793. 634
 - Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models, 2023. URL https://arxiv.org/ abs/2210.03629.
- Yuntong Zhang, Haifeng Ruan, Zhiyu Fan, and Abhik Roychoudhury. Autocoderover: Autonomous 639 program improvement, 2024. URL https://arxiv.org/abs/2404.05427. 640
- 641 Jiang Zheng, Laurie Williams, Nachiappan Nagappan, Will Snipes, John P Hudepohl, and Mladen A 642 Vouk. On the value of static analysis for fault detection in software. *IEEE transactions on software* 643 engineering, 32(4):240-253, 2006.
- 644 645

635

636

637

638

601

- 646
- 647

А **REDO** PSEUDOCODE

52	1:	Input: original github repo, problem statement <i>i</i> , generated patch <i>p</i> , function searching method
53		f, git apply function g, git revert function r, static analysis tool t, LLM l
54	2:	Search for original python script c containing the modified functions
55	3:	Detect runtime errors S_{Orig} in the original implementation using t
56	4:	Apply the patch
57	5:	Detect runtime errors S_{Mod} in the modified implementation using t
18	6:	Revert applied patch
:0	7:	if $S_{\text{Orig}} \neq S_{\text{Mod}}$ then
59	8:	return 'Unsafe'
0	9:	else
51	10:	Prompt 1 with i , original script c , and modification patch p
52	11:	if Do not detect runtime errors then
63	12:	return 'Safe'
64	13:	else
65	14:	return 'Unsafe'
66	15:	end if
7	16:	end if

В **CONFUSION MATRICES**

С ADDITIONAL QUANTITATIVE RESULTS

Table 4: Performance metrics by method and benchmark using Claude-3 OPUS

Coding agent	Metric			Meth	nod	
County ugont	Wietite	Pyflakes	PyRight	LLM	REDO-Pyflakes	REDO
	Precision	32.1	43.7	43.1	35.6	39.2
CodeStory	Recall	22.8	48.1	39.2	46.8	62.0
	F1	26.7	45.8	41.1	40.4	48.0
	Precision	34.1	33.9	34.0	33.9	33.6
Demo	Recall	34.9	45.3	20.9	47.7	55.8
	F1	34.5	38.8	25.9	39.6	41.9
	Precision	22.9	33.8	30.8	28.7	32.8
Marscode	Recall	16.4	37.3	35.8	46.3	59.7
	F1	19.1	35.5	33.1	35.4	42.3
	Precision	45.8	48.3	48.6	45.9	49.1
Lingma	Recall	27.8	29.9	35.1	51.5	54.6
	F1	34.6	36.9	40.7	48.5	51.7
	Precision	53.3	42.1	56.0	50.6	48.2
Droid	Recall	22.0	14.7	25.7	40.4	36.7
	F1	31.2	21.8	35.2	44.9	41.7
	Precision	45.3	54.7	43.1	41.2	45.9
ACR	Recall	23.5	34.3	27.5	41.2	49.0
	F1	31.0	42.2	33.5	41.2	47.4

D **ENSEMBLE PATCHES**

For a base agent, we denote the number of instance on SWEDE switching from *failed* to *passed* as M; and from *passed* to *failed* as N. The performance enhancement E of the base agent is therefore measured by the difference in number between newly passed instances to newly failed instances, or







743

744

745

746

E = M - N. We further consider two scenarios. First, only passed instances or instances with runtime errors using the base agent are considered; second, all instances are considered. These two scenarios evaluates how well REDO detects specifically runtime errors and any errors, respectively. We report the results in Figure 7a and 7b, where the rows represent different base agents and columns represent auxiliary agents.

747 First, the enhancements E's are positive for the majority of the entries, and are especially significant 748 when ensemble a base agent with a more powerful auxiliary agent (as shown by numbers in lower 749 triangular entries). The average enhancement is **3.1** in Figure 7a and **4.6** in Figure 7b. These en-750 hancements are meaningful as more powerful agents usually involve more API calls. Switching 751 to more powerful agents only when necessary can reduce over cost and time. Second, when the 752 base agent is CodeStory, the enhancements are negative. This could result from the fact that the 753 ensemble algorithm currently does not consider the difference in the coding editing performance between base and auxiliary agents, therefore being over-confident with less powerful auxiliary agents. 754 Since the ensemble algorithm is orthogonal to our contribution, we regard improving the algorithm 755 as an important future work.

Coding agent	Metric			Meth	nod	
couning ugoin	metric	Pyflakes	PyRight	LLM	REDO-Pyflakes	REDO
	Precision	32.1	43.7	33.3	32.9	34.1
CodeStory	Recall	22.8	48.1	67.1	69.6	75.9
	F1	26.7	45.8	44.5	44.7	47.1
	Precision	34.1	33.9	31.8	31.4	32.1
Demo	Recall	34.9	45.3	65.1	70.9	77.9
	F1	34.5	38.8	42.7	43.6	45.4
	Precision	22.9	33.8	29.5	28.2	29.3
Marscode	Recall	16.4	37.3	68.7	74.6	76.1
	F1	19.1	35.5	41.3	41.0	42.3
	Precision	45.8	48.3	40.9	40.7	41.4
Lingma	Recall	27.8	29.9	69.1	74.2	74.2
-	F1	34.6	36.9	51.3	52.6	53.1
	Precision	53.3	42.1	43.7	42.3	43.0
Droid	Recall	22.0	14.7	67.0	70.6	70.6
	F1	31.2	21.8	52.9	52.9	53.5
	Precision	45.3	54.7	39.6	39.1	40.9
ACR	Recall	23.5	34.3	63.7	68.6	72.5
	F1	31.0	42.2	48.9	49.8	52.3

Table 5: Performance metrics by method and benchmark using Claude-3 SONNET and tempera ture=0.5

Table 6: Means and standard deviations using PyRight, LLM, and REDO on STA.

Metric		Meth	od				
	PyRight	LLM	REDO-PyRight				
Without context							
Accuracy	74.5 _{0.1}	61.8 _{0.0}	75.1 _{0.3}				
Running Accuracy [†]	76.4 _{0.0}	64.1 _{0.1}	$76.8_{0.3}$				
W.F1	67.2 _{0.1}	$54.0_{\ 0.0}$	69.9 _{0.2}				
Running W.F1	$70.1_{\ 0.0}$	57.2 _{0.0}	$72.6_{0.2}$				
	With con	ntext					
Accuracy	74.5 _{0.1}	64.3 _{0.1}	77.4 0.4				
Running Accuracy	76.4 _{0.0}	$66.8_{0.2}$	79.3 _{0.3}				
W.F1	67.2 _{0.1}	58.6 _{0.1}	74.5 $_{0.3}$				
Running W.F1	$70.1_{\ 0.0}$	$62.1_{\ 0.1}$	77.5 $_{0.2}$				

[†] Metrics with "Running" prefixes are those evaluated on our running results.

E PROMPT TEMPLATES

E.1 PROMPT TEMPLATE: SYSTEM PROMPT

You are an experienced program analyzer who can identify → potential runtime errors without running the programs.

E.2 PROMPT TEMPLATE: ERROR DETECTION ON CODE EDITING

A modification patch is proposed to resolve an issue with the \hookrightarrow current github repo. This modification might introduce \hookrightarrow runtime errors that cannot be captured by static analysis



876

864 {modification_patch} 865 </Modification Patch> 866 867 First, please check if there are potential runtime errors, please \hookrightarrow list their error type and reasoning in the <Runtime 868 \hookrightarrow Errors ></Runtime Errors> section , with the format 869 \hookrightarrow [ErrorType]: [Reasoning]. If there are no potential runtime 870 \hookrightarrow errors, please return 'Safe'; otherwise, please return 871 \hookrightarrow 'Unsafe'. The conclusion should be wrapped by 872 \hookrightarrow <Conclusion ></Conclusion >. 873

E.3 PROMPT TEMPLATE: ERROR DETECTION ON CODENET WITH INPUT CONTEXT

877 Given the description of input and the implemented script, please 878 \hookrightarrow check if the implementation contains runtime errors. You 879 \hookrightarrow can assume that the inputs are always valid, and relect the 880 \hookrightarrow common case. 881 882 Here is the implementation: <Implementation> 883 {implementation} 884 </Implementation> 885 886 Here is the description of the input: 887 <Input description> 888 {input} 889 </Input description> 890 891 Potential runtime errors are: 892 <Error list> 1: 'No Error', 893 2: 'Other', 894 3: 'Timeout' 895 4: 'AssertionError' 896 5: 'AttributeError', 897 6: 'decimal', 898 7: 'EOFError', 899 8: 'FileNotFoundError', 900 9: 'ImportError', 901 10: 'IndentationError', 902 11: 'IndexError', 903 12:'KeyError', 13: 'MathDomainError', 904 14: 'MemoryError', 905 15: 'ModuleNotFoundError', 906 16: 'NameError', 907 'OSError', 17: 908 18: 'OverflowError', 909 're.error', 19: 910 'RecursionError', 20: 911 21: 'RuntimeError ' 912 22: 'StopIteration', 913 23: 'SyntaxError', 914 'TabError', 24: 915 25: 'TypeError', 26: 'UnboundLocalError', 916 27: 'ValueError', 917 28: 'ZeroDivisionError',

918 29: 'numpy.AxisError' 919 </Error list> 920 921 Please explain the logic of the implementation in the 922 \hookrightarrow "Implementation" section, especially how empty strings or \hookrightarrow lists are handled. If the implementation is mostly correct 923 \hookrightarrow and should run without errors in most cases, please claim 924 \hookrightarrow "No Error"; finally, the index of the identified runtime 925 \hookrightarrow error that crashes the program in the "Conclusion" section, 926 \hookrightarrow being wrapped by <Conclusion></Conclusion>. 927 928 929 E.4 PROMPT TEMPLATE: ERROR DETECTION ON CODENET WITHOUT INPUT CONTEXT 930 931 Given the implemented script, please check if the implementation 932 \hookrightarrow contains runtime errors. Please assume that the inputs are 933 \hookrightarrow always valid; and only reflect the most common case. 934 935 Here is the implementation: 936 <Implementation> 937 {implementation} </Implementation> 938 939 Potential runtime errors are: 940 <Error list> 941 1: 'No Error', 942 2: 'Other', 943 3: 'Timeout', 944 4: 'AssertionError', 945 5: 'AttributeError', 946 6: 'decimal' 947 7: 'EOFError', 8: 'FileNotFoundError', 948 9: 'ImportError' 949 10: 'IndentationError', 950 11: 'IndexError', 951 12: 'KeyError', 952 13: 'MathDomainError', 953 14: 'MemoryError', 954 15: 'ModuleNotFoundError', 955 16: 'NameError', 956 17: 'OSError' 957 18: 'OverflowError', 19: 're.error', 958 20: 'RecursionError', 959 21: 'RuntimeError ' 960 22: 'StopIteration ' 961 'SyntaxError', 23: 962 'TabError' 24: 963 'TypeError ' 25: 964 'UnboundLocalError', 26: 965 'ValueError', 27: 966 28: 'ZeroDivisionError', 967 29: 'numpy.AxisError' 968 </Error list> 969 Please explain the logic of the implementation in the 970 \hookrightarrow "Implementation" section, especially how empty strings or 971 \hookrightarrow lists are handled. If the implementation is mostly correct

1011

 \hookrightarrow "No Error"; finally, the index of the identified runtime \hookrightarrow error that crashes the program in the "Conclusion" section, 975 \hookrightarrow being wrapped by <Conclusion > </Conclusion > .976 977 E.5 PROMPT TEMPLATE: ERROR DETECTION ON SWE-BENCH-LITE 978 979 A modification patch is proposed to resolve an issue with the 980 \hookrightarrow current github repo. This modification might contain 981 \hookrightarrow runtime errors that will crash the unit tests but cannot be 982 \hookrightarrow captured by static analysis tools. Your task is to check 983 \hookrightarrow whether those errors exist in the current modification. 984 985 First, you are provided with the problem statement, which 986 \hookrightarrow describes the issue and hints on how the modification patch 987 \rightarrow will be tested. The problem statement is as follows: 988 <Problem Statement> 989 {problem_statement} </Problem Statement> 990 991 Then, you will be provided with the original implementation of 992 \hookrightarrow python scripts containing modified functions: 993 <Original Implementation> 994 {original_implementation} 995 </Original Implementation> 996 997 Finaly, the modification patch is given below: 998 <Modificatoin Patch> 999 {modification_patch} 1000 </Modification Patch> 1001 Please identify potential runtime errors that can crash the 1002 \rightarrow program but cannot be captured by static analysis tools; 1003 \hookrightarrow and list them in the <Potential Errors></Potential Errors> 1004 \hookrightarrow section. Next, Prune errors that are unlikely to be 1005 \hookrightarrow relevant to the problem statement. Please list these errors \hookrightarrow in the "Remaining Errors" section, being wrapped by 1007 \hookrightarrow <Remaining Errors ></Remaining Errors >. 1008 1009 1010 E.6 PROMPT TEMPLATE: PATCH FIXING

 \hookrightarrow and should run without errors in most cases, please claim

Your task is to update the provided code files to prevent the 1012 \hookrightarrow previously detected runtime errors. You will be provided 1013 \hookrightarrow with relevant code chunks and identified errors. 1014 1015 Begin your response by providing a simple smoke test to test the 1016 \rightarrow updated code within <test></test> tags. The rest of your 1017 \hookrightarrow response should provide the updated code to prevent the 1018 \hookrightarrow runtime errors, matching the exact format of the provided 1019 \hookrightarrow <code> below, including the <code> and <file> tags, and the 1020 \hookrightarrow name and start-line attributes. If a code chunk does not 1021 \hookrightarrow need any modification, it can be omitted from your 1022 \hookrightarrow response. Each code chunk you update to solve the problem \hookrightarrow must be rewritten in full, including lines that are 1023 \hookrightarrow unchanged. The name and start-line XML attributes in your 1024 \hookrightarrow response should always match those in the code below 1025 \hookrightarrow exactly – do not change them. For example, if 100 lines of

1026 \hookrightarrow code are passed for a code chunk, but you only modify 5 1027 \rightarrow lines, you must still include the full code chunk in your 1028 \hookrightarrow response with the original start-line attribute. If you are 1029 \hookrightarrow able to solve the problem, provide 1030 \hookrightarrow <outcome>Complete </outcome> in your response, otherwise \hookrightarrow provide 1031 <outcome>Incomplete </outcome>, along with brief feedback and next 1032 \hookrightarrow steps within 1033 <assessment></assessment> tags. 1034 1035 Below is a simple example of a valid response: 1036 <example> 1037 <smoke_test> 1038 from path.to.file import combine_numbers 1039 combine_numbers (123, 456) 1040 </smoke_test> As requested, in the updated code below, I've rewritten the full 1041 \hookrightarrow chunks provided, even those parts that remain unchanged, 1042 \hookrightarrow such as the load_file function. 1043 <code> 1044 <file name="path/to/file1.py" start_line="5"> 1045 import numpy as np 1046 </file>1047 <file name="path/to/file1.py" start_line="23"> 1048 def load_file(path): 1049 with open (path, "r") as f: 1050 content = f.read()1051 return content 1052 def combine_numbers(a, b): 1053 return {{ 1054 "sum": a + b, 1055 "difference": a - b, 1056 "product": a * b, 1057 "quotient": a / b, 1058 "geometric_mean": geometric_mean(a, b), 1059 }} </file>1061 </code> 1062 <outcome>Incomplete </outcome> 1063 <assessment> Although this patch adds a geometric mean calculation, it does 1064 \hookrightarrow not import the required function to the file. The next step 1065 \hookrightarrow is to import the `geometric_mean` function to 1066 → path/to/file1.py 1067 </assessment> 1068 </example> 1069 1070 Here are the detected runtime errors: 1071 <Rumetime errors> 1072 {remaining_issues} 1073 </Runtime errors> 1074 1075 Here are the code chunks: <code> 1076 {code} 1077 </code> 1078 1079

1080 F ADDITIONAL QUALITATIVE RESULTS

	LLM prediction
I'11	analyze the modification patch for potential runtime errors:
<runt [Attr</runt 	<pre>ime Errors> ibuteError]: The modification adds `user.email` to the hash value calculation.</pre>
Howe	er, not all user models in Django are guaranteed to have an email field. If a custo
user <td>model is used that doesn't include an email field, this could raise an AttributeErr time Errors></td>	model is used that doesn't include an email field, this could raise an AttributeErr time Errors>
-(00	Junion
Unsat	e
<td>clusion></td>	clusion>
The m	odification introduces a potential runtime error that static analysis tools might r
catch While	. The main concern is the addition of `user.email` to the hash value calculation.
Abstr	actBaseUser class. Custom user models might not have this field, which could lead t
an At	tributeError at runtime when trying to access the non-existent email attribute.
This	change assumes that all user models will have an email field, which isn't always the
exist	ence of the email attribute before using it, or to provide a fallback value if the
email	attribute is not present.
	Evaluation log
Traceb	ack (most recent call last):
File	<pre>"/opt/pyenv/versions/3.6.15/lib/python3.6/unittest/case.py", line 59, in rtFroeutor viold</pre>
File	<pre>'/opt/pyenv/versions/3.6.15/lib/python3.6/unittest/case.py", line 523, in subTest</pre>
File	<pre>"/opt/django_django/tests/auth_tests/test_tokens.py", line 57, in with different email tk1 = n0 make token(user)</pre>
File	<pre>"/opt/djangodjango/django/contrib/auth/tokens.py", line 28, in make_token return</pre>
selfu	<pre>make_token_with_timestamp(user, selfnum_seconds(selfnow())) "/ont/diaggodiaggo/contrib/auth/tokens_py" line 70 in</pre>
make	<pre>coken_with_timestamp selfmake_hash_value(user, timestamp),</pre>
File return	<pre>"/opt/djangodjango/django/contrib/auth/tokens.py", line 98, in _make_hash_value str(user.pk) + user.password + str(login timestamp) + user.email + str(timestamp)</pre>
A + + + + + +	the Farmer - Louis and Frield - chiest has an attribute lensil
	Generated patch
diff -	git a/django/contrib/auth/tokens.py b/django/contrib/auth/tokens.py
a,	'django/contrib/auth/tokens.py
+++ b,	django/contrib/auth/tokens.py
<u>uu</u> –o.	used:
	1. The password field will change upon a password reset (even if the
-	2. The last_login field will usually be updated very shortly after
+	2. The email field will change and invalidate the token.
+	a password reset.
	Failing these things acttings DACCHORD DECET TIMEOUT
	railing those things, settings.PASSWUKD_RESE!_IIMEUUI eventually
@@ -94	railing inose things, settings.PASSWUKD_RESEI_IIMEUUI eventually invalidates the token. 7 +95,7 @@ class PasswordResetTokenGenerator:
@@ -94	<pre>raiting those things, settings.PASSWUKD_RESEL_IIMEOUT eventually invalidates the token. 1,7 +95,7 @c class PasswordResetTokenGenerator: # Truncate microseconds so that tokens are consistent even if the # database descels a so that tokens are consistent even if the</pre>
@@ -94	<pre>raiting those things, settings.PASSWORD_RESEL_TIMEOUT eventually invalidates the token. 4,7 +95,7 @@ class PasswordResetTokenGenerator: # Truncate microseconds so that tokens are consistent even if the # database doesn't support microseconds. login timestamp = '' if user.last login is None else</pre>
@@ -94 user.	<pre>railing inose things, settings.PASSWORD_RESEI_IIMEDUI eventually invalidates the token. i,7 +95,7 @@ class PasswordResetTokenGenerator: # Truncate microseconds so that tokens are consistent even if the # database doesn't support microseconds. login_timestamp = '' if user.last_login is None else .ast_login.replace(microsecond=0, tzinfo=None)</pre>
@@ -94 user. - +	<pre>railing inose things, settings.PASSWUKD_RESEI_IIMEDUI eventually invalidates the token. 4,7 +95,7 @@ class PasswordResetTokenGenerator: # Truncate microseconds so that tokens are consistent even if the # database doesn't support microseconds. login_timestamp = '' if user.last_login is None else .ast_login.replace(microsecond=0, tzinfo=None) return str(user.pk) + user.password + str(login_timestamp) + str(timestamp) return str(user.pk) + user.password + str(login timestamp) + user.email +</pre>
@@ -94 user. - + str(t:	<pre>railing inose things, settings.PASSWORD_RESEL_IIMEDUI eventually invalidates the token. 4,7 +95,7 @@ class PasswordResetTokenGenerator: # Truncate microseconds so that tokens are consistent even if the # database doesn't support microseconds. login_timestamp = '' if user.last_login is None else .ast_login.replace(microsecond=0, tzinfo=None) return str(user.pk) + user.password + str(login_timestamp) + str(timestamp) return str(user.pk) + user.password + str(login_timestamp) + user.email + .mestamp)</pre>
@@ -94 user. - + str(t:	<pre>railing inose things, settings.PASSWORD_RESEL_IIMEDUI eventually invalidates the token. 4,7 +95,7 @@ class PasswordResetTokenGenerator: # Truncate microseconds so that tokens are consistent even if the # database doesn't support microseconds. login_timestamp = '' if user.last_login is None else .ast_login.replace(microsecond=0, tzinfo=None) return str(user.pk) + user.password + str(login_timestamp) + str(timestamp) return str(user.pk) + user.password + str(login_timestamp) + user.email + .mestamp) lef _num_seconds(self, dt):</pre>
@@ -94 user. - + str(t:	<pre>invalidates the token. invalidates the token. ivvalidates the t</pre>
@@ -94 user. - + str(t:	<pre>railing inose things, settings.PASSWORD_RESEL_IIMEDUI eventually invalidates the token. 1,7 +95,7 @@ class PasswordResetTokenGenerator: # Truncate microseconds so that tokens are consistent even if the # database doesn't support microseconds. login_timestamp = '' if user.last_login is None else .ast_login.replace(microsecond=0, tzinfo=None) return str(user.pk) + user.password + str(login_timestamp) + str(timestamp) return str(user.pk) + user.password + str(login_timestamp) + user.email + .mestamp) lef _num_seconds(self, dt): return int((dt - datetime(2001, 1, 1)).total_seconds())</pre>
@@ -94 user. - + str(t:	<pre>railing inose things, settings.PASSWORD_RESET_IIMEDUI eventually invalidates the token. 4,7 +95,7 @@ class PasswordResetTokenGenerator: # Truncate microseconds so that tokens are consistent even if the # database doesn't support microseconds. login_timestamp = '' if user.last_login is None else .ast_login.replace(microsecond=0, tzinfo=None) return str(user.pk) + user.password + str(login_timestamp) + str(timestamp) return str(user.pk) + user.password + str(login_timestamp) + user.email + .mestamp) lef _num_seconds(self, dt): return int((dt - datetime(2001, 1, 1)).total_seconds()) Figure 8: Generated patch and LLM response on Django-13551 using Lingma.</pre>
@@ -94 user. - + str(t: (<pre>railing inose things, settings.PASSWORD_RESEL_ITMEDUT eventually invalidates the token. 4,7 +95,7 @@ class PasswordResetTokenGenerator: # Truncate microseconds so that tokens are consistent even if the # database doesn't support microseconds. login_timestamp = '' if user.last_login is None else .ast_login.replace(microsecond=0, tzinfo=None) return str(user.pk) + user.password + str(login_timestamp) + str(timestamp) return str(user.pk) + user.password + str(login_timestamp) + user.email + .mestamp) lef _num_seconds(self, dt): return int((dt - datetime(2001, 1, 1)).total_seconds()) Figure 8: Generated patch and LLM response on Django-13551 using Lingma.</pre>
@@ -94 user. + str(t: (<pre>railing inose things, settings.PASSWORD_RESET_IIMEDUI eventually invalidates the token. 4,7 +95,7 @@ class PasswordResetTokenGenerator: # Truncate microseconds so that tokens are consistent even if the # database doesn't support microseconds. login_timestamp = '' if user.last_login is None else .ast_login.replace(microsecond=0, tzinfo=None) return str(user.pk) + user.password + str(login_timestamp) + str(timestamp) return str(user.pk) + user.password + str(login_timestamp) + user.email + .mestamp) lef _num_seconds(self, dt): return int((dt - datetime(2001, 1, 1)).total_seconds()) Figure 8: Generated patch and LLM response on Django-13551 using Lingma.</pre>

1131 F.1 ANOTHER SUCCESSFUL EXAMPLE 1132

1133 This example is on instance *Django-10924* using *Lingma*. As shown in Figure 10, the modification fails due to calling a non-callable object *field.path*. According to the modification patch in Figure 10,

this happens if a non-callable *path* variable is input into the init function. As this error is only triggered during runtime, static analysis tools fail to capture it. On the other hand, as shown in Figure 10, the LLM succeeds to identify this situation and foresee the runtime error.

¹¹³⁸ F.2 ANOTHER FAILED EXAMPLE

This example is on instance *Sympy-13471* using *Mentabot*, which actually passed the unit test. As
shown in Figure 11, the LLM identifies a corner case where the variable *num* is not a string, which
will fail the program. However, since this corner case is not triggered, the instance passed the unit
test, making the LLM prediction a false positive alarm.

1145 G PATCH FIXING

As depicted in Figure 12, REDO identifies two potential runtime errors, including an AttributeError triggered by invoking the prepare_value attribute from JSONField. The evaluation log after executing the modified implementation confirms that this AttributeError indeed caused a crash.
Subsequently, only basing on the results from REDO and without knowing the true error, Figure 13 shows that the corrected patch avoids invoking JSONField, which successfully mitigates the AttributeError.

We remark that although the AttributeError is resolved, the fixed patch may introduce new AssertionError instances or other runtime errors. The AssertionError lies beyond the detection capabilities of REDO, and the emergence of additional runtime errors suggests that an iterative process may be required to fully eliminate all such errors.

I've	analyzed the modification patch and the original implementation. Here's my
asses	sment:
<runt< td=""><td>ime Errors></td></runt<>	ime Errors>
No po	rtential runtime errors detected.
<td>time Errors></td>	time Errors>
<conc< td=""><td>lusion></td></conc<>	lusion>
Safe	
<td>iclusion></td>	iclusion>
Tho m	adification patch appears to be safe and deep't introduce any obvious suntime
error	s. The changes are focused on preserving the GROUP BY clause for subqueries, which
a log	ical extension of the existing functionality. The modifications are made in a way
that'	s consistent with the existing code structure and don't introduce any new variable
or co	implex logic that could lead to runtime errors.
The m	ain changes involve:
1. Ad 2. En	wing a cneck to preserve the GROUP BY clause when building filters with subqueries
2. LI	Suring the should be clause is preserved when spritting excludes.
These	e modifications are straightforward and don't alter the fundamental behavior of th
query	v building process. They simply ensure that GROUP BY information is maintained when
work1	ng with subqueries, which should improve the correctness of the generated SQL with which new risks
111110	
	Evaluation log
r.	
Fil	<pre>Dack (most recent call last): e "/ont/diango_diango/db/backends/utils_nu"_line_96_ inoxecute</pre>
r	eturn self.cursor.execute(sql, params)
Fil	e "/opt/djangodjango/django/db/backends/sqlite3/base.py", line 396, in execute
r	eturn Database.Cursor.execute(self, query, params)
quit	es.uperationalError: no such column: U0.name
	Generated patch
dif	<pre>fgit a/django/db/models/fields/initpy b/django/db/models/fields/initp ex 0c64748ccd7efc0967c8 100644</pre>
dif ind	<pre>fgit a/django/db/models/fields/initpy b/django/db/models/fields/initp ex 0c64748ccd7efc0967c8 100644 a/django/db/models/fields/initpy</pre>
dif ind 	<pre>fgit a/django/db/models/fields/initpy b/django/db/models/fields/initp ex 0c64748ccd7efc0967c8 100644 a/django/db/models/fields/initpy b/django/db/models/fields/initpy</pre>
dif ind +++ @@	fgit a/django/db/models/fields/initpy b/django/db/models/fields/initp ex 0c64748ccd7efc0967c8 100644 · a/django/db/models/fields/initpy · b/django/db/models/fields/initpy -1663,7 +1663,9 @@ class FilePathField(Field):
dif ind +++ @@	<pre>fgit a/django/db/models/fields/initpy b/django/db/models/fields/init; ex 0c64748ccd7efc0967c8 100644 a/django/db/models/fields/initpy b/django/db/models/fields/initpy -1663,7 +1663,9 @@ class FilePathField(Field): definit(self, verbose_name=None, name=None, nath='', match=None,</pre>
dif ind +++ @@	<pre>fgit a/django/db/models/fields/initpy b/django/db/models/fields/init; ex 0c64748ccd7efc0967c8 100644 a/django/db/models/fields/initpy b/django/db/models/fields/initpy -1663,7 +1663,9 @@ class FilePathField(Field): definit(self, verbose_name=None, name=None, path='', match=None, recursive=False, allow_files=True, allow_folders=False, **kwargs):</pre>
dif ind +++ @@	<pre>fgit a/django/db/models/fields/initpy b/django/db/models/fields/init; ex 0c64748ccd7efc0967c8 100644 a/django/db/models/fields/initpy b/django/db/models/fields/initpy -1663,7 +1663,9 @@ class FilePathField(Field): definit(self, verbose_name=None, name=None, path='', match=None, recursive=False, allow_files=True, allow_folders=False, **kwargs): self.path, self.match, self.recursive = path, match, recursive</pre>
dif ind +++ @@ - +	<pre>fgit a/django/db/models/fields/initpy b/django/db/models/fields/initp ex 0c64748ccd7efc0967c8 100644 a/django/db/models/fields/initpy b/django/db/models/fields/initpy -1663,7 +1663,9 @@ class FilePathField(Field): definit(self, verbose_name=None, name=None, path='', match=None, recursive=False, allow_files=True, allow_folders=False, **kwargs): self.path, self.match, self.recursive = path, match, recursive self.path_callable = path if callable(path) else None colf path = path if callable(path) else None</pre>
dif ind +++ @@ - + + +	<pre>fgit a/django/db/models/fields/initpy b/django/db/models/fields/initp ex 0c64748ccd7efc0967c8 100644 a/django/db/models/fields/initpy b/django/db/models/fields/initpy -1663,7 +1663,9 @@ class FilePathField(Field): definit(self, verbose_name=None, name=None, path='', match=None,</pre>
dif ind +++ @@ - + + + +	<pre>fgit a/django/db/models/fields/initpy b/django/db/models/fields/initp ex 0c64748ccd7efc0967c8 100644 a/django/db/models/fields/initpy b/django/db/models/fields/initpy b/django/db/models/fields/initpy -1663,7 +1663,9 @@ class FilePathField(Field): definit(self, verbose_name=None, name=None, path='', match=None,</pre>
dif ind +++ @@ - + + + +	<pre>fgit a/django/db/models/fields/initpy b/django/db/models/fields/initp ex 0c64748ccd7efc0967c8 100644 a/django/db/models/fields/initpy b/django/db/models/fields/initpy b/django/db/models/fields/initpy -1663,7 +1663,9 @@ class FilePathField(Field): definit(self, verbose_name=None, name=None, path='', match=None,</pre>
dif ind +++ @@ - + + + +	<pre>fgit a/django/db/models/fields/initpy b/django/db/models/fields/initp ex 0c64748ccd7efc0967c8 100644 a/django/db/models/fields/initpy b/django/db/models/fields/initpy -1663,7 +1663,9 @@ class FilePathField(Field): definit(self, verbose_name=None, name=None, path='', match=None, recursive=False, allow_files=True, allow_folders=False, **kwargs): self.path, self.match, self.recursive = path, match, recursive self.path_callable = path if callable(path) else None self.path = path() if callable(path) else Path self.match, self.recursive = match, recursive self.allow_files, self.allow_folders = allow_files, allow_folders kwargs.setdefault('max_length', 100) super()init(verbose_name, name, **kwargs) =1687,7 +1680,9 @@ class FilePathField(Field):</pre>
dif ind +++ @@ + + + + +	<pre>fgit a/django/db/models/fields/initpy b/django/db/models/fields/initp ex 0c64748ccd7efc0967c8 100644 a/django/db/models/fields/initpy b/django/db/models/fields/initpy -1663,7 +1663,9 @@ class FilePathField(Field): definit(self, verbose_name=None, name=None, path='', match=None, recursive=False, allow_files=True, allow_folders=False, **kwargs): self.path, self.match, self.recursive = path, match, recursive self.path_callable = path if callable(path) else None self.path = path() if callable(path) else None self.match, self.recursive = match, recursive self.allow_files, self.allow_folders = allow_files, allow_folders kwargs.setdefault('max_length', 100) super()init(verbose_name, name, **kwargs) -1687,7 +1689,9 @@ class FilePathField(Field):</pre>
dif ind +++ @@ + + + +	<pre>fgit a/django/db/models/fields/initpy b/django/db/models/fields/initp ex 0c64748ccd7efc0967c8 100644 a/django/db/models/fields/initpy b/django/db/models/fields/initpy -1663,7 +1663,9 @@ class FilePathField(Field): definit(self, verbose_name=None, name=None, path='', match=None, recursive=False, allow_files=True, allow_folders=False, **kwargs): self.path, self.match, self.recursive = path, match, recursive self.path_callable = path if callable(path) else None self.path = path() if callable(path) else Path self.match, self.recursive = match, recursive self.allow_files, self.allow_folders = allow_files, allow_folders kwargs.setdefault('max_length', 100) super()init(verbose_name, name, **kwargs) -1687,7 +1689,9 @@ class FilePathField(Field): def deconstruct(self):</pre>
dif ind +++ @@ + + + +	<pre>fgit a/django/db/models/fields/initpy b/django/db/models/fields/initp ex 0c64748ccd7efc0967c8 100644 a/django/db/models/fields/initpy b/django/db/models/fields/initpy -1663,7 +1663,9 @@ class FilePathField(Field): definit(self, verbose_name=None, name=None, path='', match=None, recursive=False, allow_files=True, allow_folders=False, **kwargs): self.path, self.match, self.recursive = path, match, recursive self.path_callable = path if callable(path) else None self.path = path() if callable(path) else Path self.match, self.recursive = match, recursive self.allow_files, self.allow_folders = allow_files, allow_folders kwargs.setdefault('max_length', 100) super()init(verbose_name, name, **kwargs) -1687,7 +1689,9 @@ class FilePathField(Field): def deconstruct(self): name, path, args, kwargs = super().deconstruct() if erel6 rest default()</pre>
dif ind +++ @@ - + + + + + -	<pre>fgit a/django/db/models/fields/initpy b/django/db/models/fields/initp ex 0c64748ccd7efc0967c8 100644 a/django/db/models/fields/initpy b/django/db/models/fields/initpy -1663,7 +1663,9 @@ class FilePathField(Field): definit(self, verbose_name=None, name=None, path='', match=None, recursive=False, allow_files=True, allow_folders=False, **kwargs): self.path, self.match, self.recursive = path, match, recursive self.path_callable = path if callable(path) else None self.path = path() if callable(path) else Path self.match, self.recursive = match, recursive self.allow_files, self.allow_folders = allow_files, allow_folders kwargs.setdefault('max_length', 100) super()init(verbose_name, name, **kwargs) -1687,7 +1689,9 @@ class FilePathField(Field): def deconstruct(self): name, path, args, kwargs = super().deconstruct() if self.path := '': if self.path := '':</pre>
dif ind +++ + + + +	<pre>fgit a/django/db/models/fields/initpy b/django/db/models/fields/initp ex 0c64748ccd7efc0967c8 100644 a/django/db/models/fields/initpy b/django/db/models/fields/initpy -1663,7 +1663,9 @@ class FilePathField(Field): definit(self, verbose_name=None, name=None, path='', match=None, recursive=False, allow_files=True, allow_folders=False, **kwargs): self.path, self.match, self.recursive = path, match, recursive self.path_callable = path if callable(path) else None self.path = path() if callable(path) else Path self.match, self.recursive = match, recursive self.allow_files, self.allow_folders = allow_files, allow_folders kwargs.setdefault('max_length', 100) super()init(verbose_name, name, **kwargs) -1687,7 +1689,9 @@ class FilePathField(Field): def deconstruct(self): name, path, args, kwargs = super().deconstruct() if self.path != '': if self.path != self.path callable</pre>
diff ind +++ + + + + + + + +	<pre>fgit a/django/db/models/fields/initpy b/django/db/models/fields/initp ex 0c64748ccd7efc0967c8 100644 a/django/db/models/fields/initpy b/django/db/models/fields/initpy -1663,7 +1663,9 @@ class FilePathField(Field): definit(self, verbose_name=None, name=None, path='', match=None, recursive=False, allow_files=True, allow_folders=False, **kwargs): self.path, self.match, self.recursive = path, match, recursive self.path_callable = path if callable(path) else None self.path = path() if callable(path) else Path self.match, self.recursive = match, recursive self.allow_files, self.allow_folders = allow_files, allow_folders kwargs.setdefault('max_length', 100) super()init(verbose_name, name, **kwargs) -1687,7 +1689,9 @@ class FilePathField(Field): def deconstruct(self): name, path, args, kwargs = super().deconstruct() if self.path != '': if self.path != '': kwargs['path'] = self.path_callable elif self.path != '':</pre>
diff ind +++ (@@ - + + + + + + +	<pre>fgit a/django/db/models/fields/initpy b/django/db/models/fields/initp ex @c64748ccd7efc0967c8 100644 a/django/db/models/fields/initpy b/django/db/models/fields/initpy -1663,7 +1663,9 @@ class FilePathField(Field): definit(self, verbose_name=None, name=None, path='', match=None,</pre>
diff ind +++ + + + + + + + +	<pre>fgit a/django/db/models/fields/initpy b/django/db/models/fields/initp ex @c64748ccd7efc0967c8 100644 a/django/db/models/fields/initpy b/django/db/models/fields/initpy -1663,7 +1663,9 @@ class FilePathField(Field): definit(self, verbose_name=None, name=None, path='', match=None,</pre>
diff ind +++ + + + + (@@	<pre>fgit a/django/db/models/fields/initpy b/django/db/models/fields/initp ex @c64748ccd7efc0967c8 100644 a/django/db/models/fields/initpy b/django/db/models/fields/initpy -1663,7 +1663,9 @@ class FilePathField(Field): definit(self, verbose_name=None, name=None, path='', match=None,</pre>
diff ind ++ + + + + + (00) - + + + + + (00) (0)	<pre>fgit a/django/db/models/fields/initpy b/django/db/models/fields/initp ex @c64748ccd7efc@967c8 100644 a/django/db/models/fields/initpy b/django/db/models/fields/initpy -1663,7 +1663,9 @@ class FilePathField(Field): definit(self, verbose_name=None, name=None, path='', match=None,</pre>
diff ind +++ + + + + + (@@	<pre>fgit a/django/db/models/fields/initpy b/django/db/models/fields/initp ex @c64748ccd7efc@967c8 100644 a/django/db/models/fields/initpy b/django/db/models/fields/initpy -1663,7 +1663,9 @@ class FilePathField(Field): definit(self, verbose_name=None, name=None, path='', match=None,</pre>
diff ind +++ + + + + + (@@ 	<pre>fgit a/django/db/models/fields/initpy b/django/db/models/fields/initp ex @c64748ccd7efc0967c8 100644 a/django/db/models/fields/initpy b/django/db/models/fields/initpy -1663,7 +1663,9 @@ class FilePathField(Field): definit(self, verbose_name=None, name=None, path='', match=None,</pre>
diff ind ++ + + + + + (@@ 	<pre>fgit a/django/db/models/fields/initpy b/django/db/models/fields/initp ex @c64748ccd7efc0967c8 100644 a/django/db/models/fields/initpy b/django/db/models/fields/initpy -1663,7 +1663,9 @@ class FilePathField(Field): definit(self, verbose_name=None, name=None, path='', match=None,</pre>
diff ind ++ + + + + + + + (@@ - - + + + + + - -	<pre>fgit a/django/db/models/fields/initpy b/django/db/models/fields/initpy ex @c64748ccd7efc0967c8 100644 a/django/db/models/fields/initpy b/django/db/models/fields/initpy -1663,7 +1663,9 @@ class FilePathField(Field): definit(self, verbose_name=None, name=None, path='', match=None,</pre>
diff ind ++ + + + + + + + + + + + + + +	<pre>fgit a/django/db/models/fields/initpy b/django/db/models/fields/initpy ex 0c64748ccd7efc0967c8 100644 a/django/db/models/fields/initpy -1663,7 +1663,9 @@ class FilePathField(Field): definit(self, verbose_name=None, name=None, path='', match=None,</pre>
diff ind ++ + + + + + + (00) - + + + + + + + + - + + + + + +	<pre>fgit a/django/db/models/fields/initpy b/django/db/models/fields/initpy ex 0c64748ccd7efc0967c8 100644 a/django/db/models/fields/initpy b/django/db/models/fields/initpy -1663,7 +1663,9 @@ class FilePathField(Field): definit(self, verbose_name=None, name=None, path='', match=None,</pre>
diff ind +++ + + + + + + + + + (00) - + + + + + + +	<pre>fgit a/django/db/models/fields/initpy b/django/db/models/fields/initpy ex 0c64748ccd7efc0967c8 100644 a/django/db/models/fields/initpy -1663,7 +1663,9 @@ class FilePathField(Field): definit(self, verbose_name=None, name=None, path='', match=None,</pre>

	LLM prediction
I've a	analyzed the modification patch and the original implementation. Here's my
asses	smerrt.
<runt:< td=""><td>ime Errors></td></runt:<>	ime Errors>
If nur	rrorj: The modification assumes that hum[1] is always a string, but it might hol n[1] is already an integer (long). calling .rstrip('L') on it will raise a TypeE
<td>time Errors></td>	time Errors>
<conc<sup>°</conc<sup>	lusion>llnsafe
conc	
The mo	pdification patch introduces a potential runtime error. It assumes that num[1] i
impler	nentation, num[1] could already be an integer (long) if the input comes from a n
pickle	ed source. In such cases, calling .rstrip('L') on an integer would raise a TypeE
To mal	ke this modification safe, it should first check if num[1] is a string before
attem	oting to strip the 'L' suffix. For example:
```py	thon
if is:	<pre>instance(num[1], string_types):</pre>
num[1]	<pre></pre>
· · · ·	- cong(num[1]) 10/
This	change would make the modification safe for all input types
11115	change would make the moullication sale for all input types.
	Evaluation log
	======================================
test_ca	llable_path (model_fields.test_filepathfield.FilePathFieldIests) Traceback (most recent call last): Fil
"/opt/d	jangodjango/tests/model_fields/test_filepathfield.py", line 21, in test_callab
self.as	sertEqual(field.path(), path) TypeError: 'str' object is not callable
	Generated patch
dift	-git a/django/db/models/fields/initpy b/django/db/models/fields/initpy
index 0	)c64748ccd7efc0967c8 100644
index @	lc64748ccd7efc0967c8 100644 ljango/db/models/fields/initpy
index ( a/c +++ b/c	0c64748ccd7efc0967c8 100644 Ijango/db/models/fields/initpy Ijango/db/models/fields/initpy 3.7 +1663 9 @@ class FilePathField(Field):
index ( a/c +++ b/c @@ -166	bc64748ccd7efc0967c8 100644 ijango/db/models/fields/initpy ljango/db/models/fields/initpy 3,7 +1663,9 @@ class FilePathField(Field):
index ( a/c +++ b/c @@ -166	<pre>bc64748ccd7efc0967c8 100644 ijango/db/models/fields/initpy ljango/db/models/fields/initpy 3,7 +1663,9 @@ class FilePathField(Field): efinit(self, verbose_name=None, name=None, path='', match=None,</pre>
index @ a/c +++ b/c @@ -166 	<pre>bc64748ccd7efc0967c8 100644 ljango/db/models/fields/initpy ljango/db/models/fields/initpy i3,7 +1663,9 @@ class FilePathField(Field): efinit(self, verbose_name=None, name=None, path='', match=None,</pre>
index ( a/c +++ b/c @@ -166 	<pre>bc64748ccd7efc0967c8 100644 ljango/db/models/fields/initpy ljango/db/models/fields/initpy 03,7 +1663,9 @@ class FilePathField(Field): efinit(self, verbose_name=None, name=None, path='', match=None,</pre>
index ( a/c +++ b/c @@ -166 	<pre>bc64748ccd7efc0967c8 100644 ljango/db/models/fields/initpy ljango/db/models/fields/initpy i3,7 +1663,9 @@ class FilePathField(Field):  ffinit(self, verbose_name=None, name=None, path='', match=None,</pre>
a/c +++ b/c @@ -166 	<pre>bc64748ccd7efc0967c8 100644 tjango/db/models/fields/initpy tjango/db/models/fields/initpy ijango/db/models/fields/initpy i3,7 +1663,9 @@ class FilePathField(Field):  ffinit(self, verbose_name=None, name=None, path='', match=None,</pre>
index ( a/c +++ b/c @@ -166 	<pre>bc64748ccd7efc0967c8 100644 ljango/db/models/fields/initpy ljango/db/models/fields/initpy i3,7 +1663,9 @@ class FilePathField(Field): effinit(self, verbose_name=None, name=None, path='', match=None,</pre>
index ( a/c +++ b/c @@ -166 de - + + + +	<pre>bc64748ccd7efc0967c8 100644 tjango/db/models/fields/initpy tjango/db/models/fields/initpy tjango/db/models/fields/initpy 33,7 +1663,9 @@ class FilePathField(Field): effinit(self, verbose_name=None, name=None, path='', match=None,</pre>
index ( a/c +++ b/c @@ -166 ++ + + +	<pre>bc64748ccd7efc0967c8 100644 ijango/db/models/fields/initpy ijango/db/models/fields/initpy 33,7 +1663,9 @@ class FilePathField(Field): efinit(self, verbose_name=None, name=None, path='', match=None,</pre>
a/c +++ b/c @@ -166  + + + + @@ -168 de	<pre>bc64748ccd7efc0967c8 100644 ijango/db/models/fields/initpy ijango/db/models/fields/initpy 3,7 +1663,9 @@ class FilePathField(Field): efinit(self, verbose_name=None, name=None, path='', match=None,</pre>
	<pre>bc64748ccd7efc0967c8 100644 ijango/db/models/fields/initpy ijango/db/models/fields/initpy i3,7 +1663,9 @@ class FilePathField(Field): efinit(self, verbose_name=None, name=None, path='', match=None,</pre>
	<pre>bc64748ccd7efc0967c8 100644 ijango/db/models/fields/initpy ijango/db/models/fields/initpy i3,7 +1663,9 @@ class FilePathField(Field):  ffinit(self, verbose_name=None, name=None, path='', match=None,</pre>
	<pre>bc64748ccd7efc0967c8 100644 jjango/db/models/fields/initpy jjango/db/models/fields/initpy 33,7 +1663,9 @@ class FilePathField(Field): efinit(self, verbose_name=None, name=None, path='', match=None,</pre>
	<pre>bc64748ccd7efc0967c8 100644 jjango/db/models/fields/initpy jjango/db/models/fields/initpy 33,7 +1663,9 @@ class FilePathField(Field): efinit(self, verbose_name=None, name=None, path='', match=None,</pre>
	<pre>bc64748ccd7efc0967c8 100644 jjango/db/models/fields/initpy jjango/db/models/fields/initpy 33,7 +1663,9 @@ class FilePathField(Field): efinit(self, verbose_name=None, name=None, path='', match=None,</pre>
	<pre>bc64748ccd7efc0967c8 100644 jjango/db/models/fields/initpy jjango/db/models/fields/initpy 33,7 +1663,9 @@ class FilePathField(Field): efinit(self, verbose_name=None, name=None, path='', match=None,</pre>
ee	<pre>bc64748ccd7efc0967c8 100644 jjango/db/models/fields/initpy Jjango/db/models/fields/initpy Jjango/db/models/fields/initpy Jjango/db/models/fields/initpy Jjango/db/models/fields/initpy Jjango/db/models/fields/_initpy Jjango/db/models/fields/_initpy Jjango/db/models/fields/_initpy Jjango/db/models/fields/_initpy Self.path, self, verbose_name=None, name=None, path='', match=None, recursive=False, allow_files=True, allow_folders=False, **kwargs): self.path, self.match, self.recursive = path, match, recursive self.path_callable = path if callable(path) else None self.path = path() if callable(path) else path self.match, self.recursive = match, recursive self.allow_files, self.allow_folders = allow_files, allow_folders kwargs.setdefault('max_length', 100) super()init(verbose_name, name, **kwargs) JJ,7 +1689,9 @@ class FilePathField(Field): Ef deconstruct(self): name, path, args, kwargs = super().deconstruct() if self.path != '': if self.path != '': kwargs['path'] = self.path_callable elif self.path != '': kwargs['path'] = self.path if self.match is not None: kwargs['match'] = self.match N8,8 +1712,9 @@ class FilePathField(Field): return str(value)</pre>
index ( a/c +++ b/c @@ -166 - + + + + + + @@ -166 - - + + + + (@ -166 - (@ -166 () - () () () () () () () () () () () () ()	<pre>bc64748ccd7efc0967c8 100644 jjango/db/models/fields/initpy jjango/db/models/fields/initpy 33,7 +1663,9 @@ class FilePathField(Field): efinit(self, verbose_name=None, name=None, path='', match=None,</pre>
index ( a/c +++ b/c @@ -166 - + + + + + + (@@ -166 - + + + + + (@@ -166 ( - + + + + + + ( ( - - - - ( - - + + + +	<pre>bc64748ccd7efc0967c8 100644 jjango/db/models/fields/initpy jjango/db/models/fields/initpy 33,7 +1663,9 @@ class FilePathField(Field): efinit(self, verbose_name=None, name=None, path='', match=None,</pre>
	<pre>bc64748ccd7efc0967c8 100644 jjango/db/models/fields/initpy jjango/db/models/fields/initpy 33,7 +1663,9 @@ class FilePathField(Field):  ffinit(self, verbose_name=None, name=None, path='', match=None,</pre>
	<pre>bc64748ccd7efc0967c8 100644 ljango/db/models/fields/initpy ljango/db/models/fields/initpy bigago/db/models/fields/initpy bigagodit</pre>
	<pre>bc64748ccd.7efc0967c8 100644 ljango/db/models/fields/initpy ljango/db/models/fields/initpy bigago/db/models/fields/initpy bigao/db/models/fields/initpy bigao/db/models/fields/initpy bigacode/</pre>



_	
	LLM prediction
	I've analyzed the modification patch and the original implementation. Here's my assessment:
	<runtime errors=""></runtime>
	[TypeError]: The modification assumes that num[1] is always a string, but it might not
	<pre>//Runtime Errors&gt;</pre>
	<conclucion>Uncofez/Conclucion&gt;</conclucion>
	The modification patch introduces a potential runtime error. It assumes that num[1] is always a string and tries to call the retrip('1') method on it However in the origin
	implementation, num[1] could already be an integer (long) if the input comes from a nor pickled source. In such cases, calling .rstrip('L') on an integer would raise a TypeErr
	To make this modification safe, it should first check if num[1] is a string before
	attempting to strip the 'L' suffix. For example:
	```python
	<pre>if isinstance(num[1], string_types):</pre>
	num[1] = num[1].rstrip('L') # Strip 'L' suttix if present num[1] = long(num[1], 16)
	This change would make the modification safe for all input types.
	Generated patch
	diffait a/sympy/core/numbers.py b/sympy/core/numbers.py
	index 7bc96e68eb4c8301c6b2 100644
	a/sympy/core/numbers.py +++ b/sympy/core/numbers.py
	<pre>@@ -1042,6 +1042,7 @@ defnew(cls, num, dps=None, prec=None, precision=None):</pre>
	<pre># it's a hexadecimal (coming from a pickled object) # assume that it is in standard form</pre>
	num = list(num)
	<pre>+ num[1] = num[1].rstrip('L') # Strip 'L' suffix if present num[1] = long(num[1]. 16)</pre>
	mof = tuple(num)

	LLM prediction
I've ana assessme	yzed the modification patch and the original implementation. Here's my it:
<runtime [Attribu method d</runtime 	Errors> :eError]: The patch assumes that JSONField has a 'prepare_value' method. If thi wesn't exist on the JSONField class, it could raise an AttributeError.
[ImportE this imp/ ImportEr <td>rror]: The patch adds an import statement for JSONField from django.db.models. ort is not available in the version of Django being used, it could raise an or. : Errors></td>	rror]: The patch adds an import statement for JSONField from django.db.models. ort is not available in the version of Django being used, it could raise an or. : Errors>
<conclus< td=""><td>.on>Unsafe</td></conclus<>	.on>Unsafe
he modi tatic a	ication patch introduces potential runtime errors that cannot be captured by alysis tools. The main concerns are:
1. The pa method i	ntch assumes the existence of a 'prepare_value' method on JSONField. If this not part of the JSONField API, it could lead to an AttributeError at runtime.
2. The a Django,∣	lded import statement for JSONField might not be compatible with all versions o potentially causing an ImportError.
These is specific thorough API incl	ues wouldn't be caught by static analysis tools because they depend on the implementation and version of Django being used. Therefore, it's important to y test this change across different Django versions and ensure that the JSONFi ides the 'prepare_value' method before implementing this patch.
	Evaluation log
======================================	Evaluation log st_json_display_for_field (admin_utils.tests.UtilsTests) [<object at="" bbe90="" object="">] (value={('a', 'b'): 'c'})</object>
ERROR: to 0x7f98600 Traceback File "/ testPart	Evaluation log est_json_display_for_field (admin_utils.tests.UtilsTests) [<object at="" object="" bbe90="">] (value={('a', 'b'): 'c'}) (most recent call last): opt/pyenv/versions/3.6.15/lib/python3.6/unittest/case.py", line 59, in xecutor</object>
ERROR: to 0x7f98600 File "/ testPartf yielo File "/ yielo File "/	Evaluation log <pre>st_json_display_for_field (admin_utils.tests.UtilsTests) [<object at="" lbbe90="" object="">] (value={('a', 'b'): 'c'}) </object></pre> (most recent call last): opt/pyenv/versions/3.6.15/lib/python3.6/unittest/case.py", line 59, in xecutor opt/pyenv/versions/3.6.15/lib/python3.6/unittest/case.py", line 523, in subTes opt/diango_diango/tests/admin_utils/tests.py". line 192. in
ERROR: to 0x7f98600 File ", testParth yielo File ", yielo File ", test_json disp File ", retu	Evaluation log st_json_display_for_field (admin_utils.tests.UtilsTests) [<object at="" bbe90="" object="">] (value={('a', 'b'): 'c'}) (most recent call last): opt/pyenv/versions/3.6.15/lib/python3.6/unittest/case.py", line 59, in xecutor opt/pyenv/versions/3.6.15/lib/python3.6/unittest/case.py", line 523, in subTes opt/django_django/tests/admin_utils/tests.py", line 192, in _display_for_field ay_for_field(value, models.JSONField(), self.empty_value), opt/django_django/contrib/admin/utils.py", line 404, in display_for_fi</object>
ERROR: to 0x7f98600 File ", testPartl yielo File ", yielo File ", test_json disp File ", retui Attributo	Evaluation log sst_json_display_for_field (admin_utils.tests.UtilsTests) [<object at="" lbbe90="" object="">] (value={('a', 'b'): 'c'}) </object>
ERROR: t Øx7f98600 File ", testPartl yiel File ", yield File ", test_json disp File ", retu Attributo	Evaluation log est_json_display_for_field (admin_utils.tests.UtilsTests) [<object at="" lbbe90="" object="">] (value={('a', 'b'): 'c'}) c (most recent call last): opt/pyenv/versions/3.6.15/lib/python3.6/unittest/case.py", line 59, in ixecutor opt/pyenv/versions/3.6.15/lib/python3.6/unittest/case.py", line 523, in subTes opt/django_django/tests/admin_utils/tests.py", line 192, in i_display_for_field ay_for_field(value, models.JSONField(), self.empty_value), opt/django_django/contrib/admin/utils.py", line 404, in display_for_fi opt/django_django/test has no attribute 'prepare_value' sts in 0.779s</object>
ERROR: ti 0x7f9860 File ", testPartl yield File ", test_jsol disp File ", test_jsol disp File ", retu Attribute FAILED (d Destroyin ('file:me Destroyin ('file:me	Evaluation log est_json_display_for_field (admin_utils.tests.UtilsTests) [<object at="" lbbe90="" object="">] (value={('a', 'b'): 'c'}) ::(most recent call last): opt/pyenv/versions/3.6.15/lib/python3.6/unittest/case.py", line 59, in xecutor opt/pyenv/versions/3.6.15/lib/python3.6/unittest/case.py", line 523, in subTest opt/djangodjango/tests/admin_utils/tests.py", line 192, indisplay_for_field ay_for_field ay_for_field(value, models.JSONField(), self.empty_value), opt/django_django/contrib/admin/utils.py", line 404, in display_for_fi field.prepare_value(value) Error: 'JSONField' object has no attribute 'prepare_value' </object>
ERROR: to 0x7f9860 File ", testParti yield File ", test_jso disp File ", retu Attribute Canada State FAILED (d Destroyin ('file:me Destroyin ('file:me [django_	Evaluation log est_json_display_for_field (admin_utils.tests.UtilsTests) [<object at="" lbbe90="" object="">] (value={('a', 'b'): 'c'}) est_json_display_for_field (admin_utils.tests.UtilsTests) [<object at="" lbbe90="" object="">] (value={('a', 'b'): 'c'}) est_most recent call last): opt/pyenv/versions/3.6.15/lib/python3.6/unittest/case.py", line 59, in xecutor opt/pyenv/versions/3.6.15/lib/python3.6/unittest/case.py", line 523, in subTest opt/djangodjango/tests/admin_utils/tests.py", line 192, in</object></object>

	Originally generated patch
dif	fgit a/django/contrib/admin/utils.py b/django/contrib/admin/utils.py
inc	ex e4ddc8f017640614ce0b 100644
+++	· b/django/contrib/admin/utils.py
@@	-377,6 +377,8 @@ def help_text_for_field(name, model):
	return help_text
+†r +	om django.db.models import JSONField # Ensure this import is added
de	f display_for_field(value, field, empty_value_display): from django.contrib.admin.templatetags.admin_list import _boolean_icon
60	-398,6 +400,8 @@ def display_for_field(value, field, empty_value_display
	return formats.number_format(value)
	elif isinstance(field, models.FileField) and value: return format html(' {} ', value.url, value)
+	elif isinstance(field, JSONField):
+	return field.prepare_value(value)
	return display_for_value(value, empty_value_display)
	Fixed patch
	a/diango/contrib/admin/utils_nv
+++	· b/django/contrib/admin/utils.py
60	-377,6 +377,9 @@
	return netp_text
+fr	om django.db.models import JSONField
+	
de	f display_for_field(value, field, empty_value_display):
	<pre>trom django.contrib.admin.templatetags.admin_list import _boolean_icon</pre>
@@	-387,6 +390,7 @@
	elif isinstance(field, models.BooleanField):
	elif value is None:
+	
	return empty_value_display
	return formats.localize(timezone.template localtime(value))
@@	-398,6 +402,8 @@
	return formats.number_format(value)
	return format html(' {} ', value.url, value)
+	elif isinstance(field, JSONField):
+	<pre>return json.dumps(value, ensure_ascii=False, indent=2) </pre>
	return display for value(value, empty value display)
@@	-408,6 +414,7 @@
	return boolean icon(value)
	elif value is None:
+	
	return empty_value_aisplay
	return str(value)
	return str(value)
	Figure 13: Original generated patch and fixed patch
	Figure 13: Original generated patch and fixed patch