

REPOPILOT: Software Agents To Resolve Software Engineering Tasks at Repository-Level Scale

Anonymous ACL submission

Abstract

Coding assistants based on Large Language Models (LLMs) have recently surged in popularity. A significant challenge for LLMs is accurately responding to user queries at the scale of entire code repositories. We propose REPOPILOT, a multi-agent-based system capable of effectively navigating through source code repositories to collect relevant information, editing code and execute programs. We demonstrate the effectiveness of REPOPILOT through extensive evaluations on challenging benchmarks, including SWE-bench and an automatically collected code generation dataset. On SWE-bench Lite, REPOPILOT achieves a 17% pass rate, establishing competitive results compared to the baseline while maintains low cost and also excels in other code intelligence tasks.

1 Introduction

The remarkable advances in Large Language Models (LLMs) have driven rapid progress in various natural language processing tasks (Touvron et al., 2023; OpenAI et al., 2023) and software development tasks (Zhang et al., 2023; Austin et al., 2021; Zhang et al., 2023). Recently, multiple attempts have been made to directly exploit LLMs as agent models to address planning problems in software engineering (SE) (Bairi et al., 2023; Zhang et al., 2024a; Yang et al.; Zhang et al., 2024c). Most Software Engineering tasks are performed within the context of an integrated development environment (IDE), which requires a quick response time for interactive usage. However, most existing approaches rely on high-performance LLMs as the main agent, leading to high query costs. Additionally, these approaches often aim to solve a single task specifically with a single LLM agent, and their agent flow and tool designs are highly customized and not general enough for a variety of software development tasks. For example, AutoCodeRover’s pipeline (Zhang et al., 2024c) is difficult to extend

to code generation and question answering tasks due to its two-stage pipeline: localizing and patching. While SWE-agent (Yang et al.) offers more flexibility in agent design, their tools are designed to quickly localize objects and reduce API cost and context length, potentially losing necessary fine-grained information during search and editing, which can be meaningful for other software development tasks (bug tracing, question answering, etc.).

We develop REPOPILOT, an agent-based framework that is general for a wide range of software development tasks while reducing the cost of using LLMs. In REPOPILOT, a software task is decomposed into four main processes: *planning*, *navigating*, *generating*, and *executing*. Each process requires a different level of intricacy and, therefore, deserves a different level of intelligence from LLMs. One can use a lightweight and computationally cheap LLM for navigation because this process is the least complex but consumes the highest number of tokens. This approach not only reduces cost but also reduces context for other processes, thereby improving performance. Despite the promise of a multi-agent system, designing such a system for software intelligent tasks is challenging for the following reasons:

- Tool design (interface and functionality) for each agent role should be well-designed and specialized depending on the role’s complexity.
- Communication between agents can cause information loss, leading to imprecise decisions in other processes, which can cause compound errors in subsequent iterations.
- Using multiple LLMs working together can substantially increase inference costs if not designed properly.

REPOPILOT follows a new centralized multi-agent system for software tasks, where the *Planner* has

080 the role of receiving user inputs (GitHub Issue,
081 code generation request), planning, and iteratively
082 deciding which child agent role (*Navigator*, *Editor*,
083 and *Executor*) should be used in the next iteration.

- 084 • Cost-Effective: We propose a centralized multi-
085 agent system (Section 3.1) using a high-cost
086 LLM for the *Planner*, which subordinates child
087 agents that are low-cost and inexpensive to exe-
088 cute requested tasks.
- 089 • Communication: We propose a communication
090 mechanism (Section 3.2) to offload tasks to child
091 agents while ensuring that the observations re-
092 ceived by the *Planner* are succinct and relevant.
093 This is the key to REPOPILOT in reducing the
094 need for making API calls to the *Planner*'s LLM
095 while ensuring the *Planner* remains grounded
096 with information from the codebase environment.
- 097 • Tool Design: We propose a suite of new tools
098 (Section 3.3) for efficient navigation, editing, and
099 execution, specialized for each agent role. These
100 tools are carefully customized for each type of
101 agent and the difficulty of each task.

102 Finally, our contributions are as follows:

- 103 1. REPOPILOT : An adaptive and general multi-
104 agent framework for different software tasks.
- 105 2. An empirical evaluation demonstrating a com-
106 parable pass rate of 17

107 2 Related Work

108 2.1 Software Engineering Benchmarks

109 Code generation benchmarks, which evaluate mod-
110 els on the task of synthesizing code from natu-
111 ral language descriptions, have served as a long
112 standing bellwether for measuring LM perfor-
113 mance (Chen et al., 2021; Austin et al., 2021;
114 Hendrycks et al., 2021; Lu et al., 2021). Subse-
115 quent works have built upon the code generation
116 task formulation to contribute new benchmarks that
117 translate problems to different (programming) lan-
118 guages (Cassano et al., 2022; Wang et al., 2022),
119 incorporate third party libraries (Lai et al., 2023;
120 Liu et al., 2023c), introduce derivative code com-
121 pletion tasks (Muennighoff et al., 2023), increase
122 test coverage (Liu et al., 2023a), change the edit
123 scope (Ding et al., 2024; Yu et al., 2024; Du et al.,
124 2023), and add robustness to dataset contamina-
125 tion (Naman Jain et al., 2024). Code generation

126 problems are largely self-contained, with short
127 problem descriptions (~100 lines) and correspond-
128 ing solutions that are similarly brief, requiring noth-
129 ing more complex than basic language primitives.
130 Tests are either handwritten or generated synthet-
131 ically via fuzz testing. In recent months, rapid
132 development of LMs has begun to saturate many of
133 these benchmarks. The top method solves 94.4%
134 of HumanEval (Zhou et al., 2023).

135 Gauging future trends with the code generation
136 task paradigm may be limited by the simplicity of
137 this setting and cost of human-in-the-loop problem
138 creation. In response, recent efforts have demon-
139 strated that software engineering (SE) can serve
140 as a diverse, challenging testbed for LM evalua-
141 tion (Zhang et al., 2023; Jimenez et al., 2023; Liu
142 et al., 2023b). Repository-level code editing intro-
143 duces many reasoning challenges grounded in real
144 SE subtasks such as spotting errant code, identifi-
145 fying cross-file relationships, and understanding
146 codebase-specific symbols and conventions. As a
147 field, SE has generally studied tasks in a more iso-
148 lated manner; prior benchmarks tend to frame prob-
149 lems in isolation from the rest of a codebase (Just
150 et al., 2014; Karampatsis and Sutton, 2020).

151 We use SWE-bench because it unites many sepa-
152 rate SE tasks such as automated program repair(Xia
153 and Zhang, 2022), bug localization (Chakraborty
154 et al., 2018; Yang et al., 2024a) under a single
155 task formulation that faithfully mirrors practical
156 SE. Furthermore, SWE-bench task instances are
157 diverse, having been collected from real GitHub
158 issues across 12 different repositories. In addition,
159 SWE-bench performance is based on rigorous, au-
160 tomatic execution-based evaluation.

161 2.2 Language Models as Agents

162 The co-emergence of stronger LMs, increasingly
163 challenging benchmarks, and practical use cases
164 have all together motivated a paradigm shift in
165 LMs' inference setting. In place of traditional
166 zero/few-shot generation, language agents (Sumers
167 et al., 2023; Xi et al., 2023; Wang et al., 2024a),
168 which use LMs to interact with a real/virtual world,
169 have proliferated as the default setting for web
170 navigation (Nakano et al., 2021; Thoppilan et al.,
171 2022; Yao et al., 2022a,b), computer control (Wu
172 et al., 2024; Xie et al., 2024), and code generation
173 tasks (Wang et al., 2023).

174 Interaction and code generation are increasingly
175 used together, with code as the modality of choice
176 for actions (Yang et al., 2024b; Wang et al., 2024b),

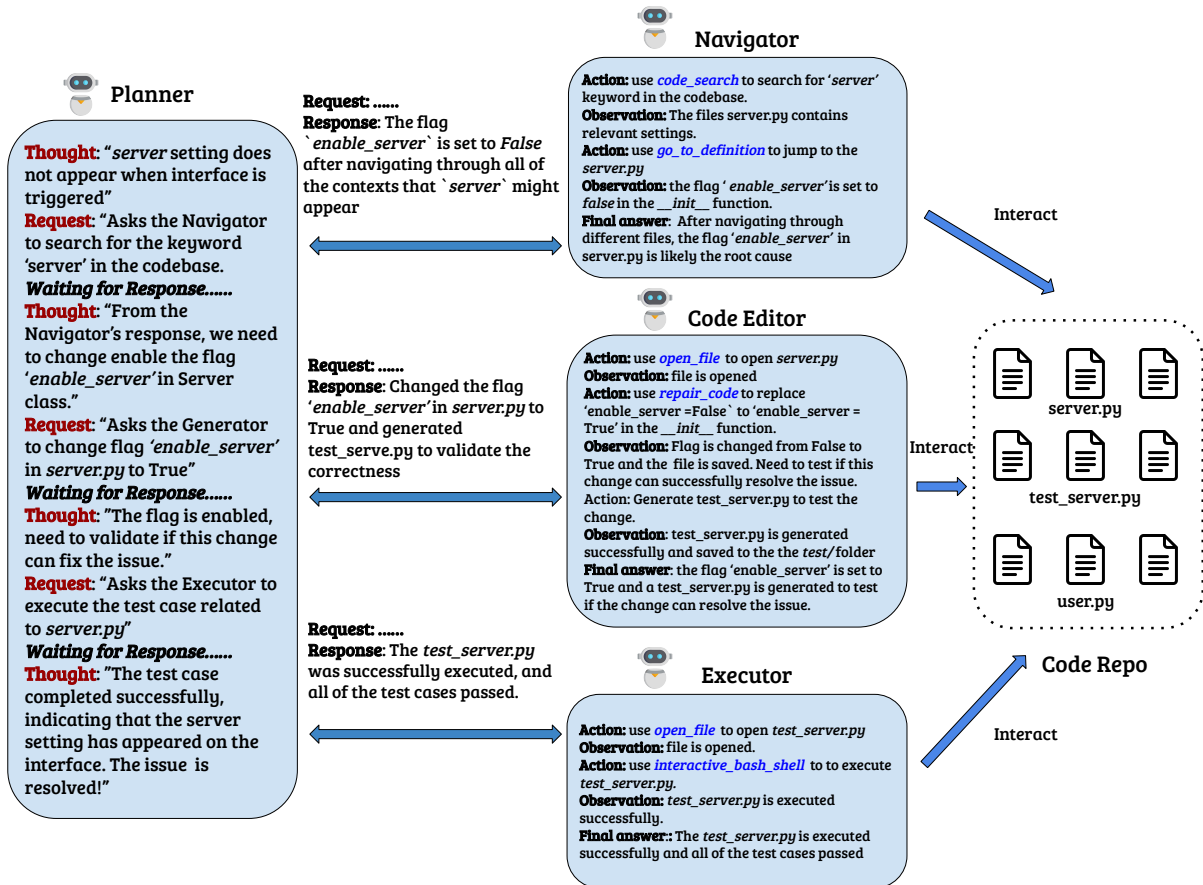


Figure 1: RepoPilot Multi-Agent Framework

177 tool construction (Wang et al., 2024c; Zhang et al.,
 178 2024b; Gu et al., 2024), and reasoning (Zelikman
 179 et al., 2023a,b; Shinn et al., 2024). Code language
 180 agents have also been applied to offensive security
 181 (Yang et al., 2023; Shao et al., 2024), theorem
 182 proving (Thakur et al., 2023).

183 To the best of our knowledge, REPOPILOT is
 184 the first to explore low-cost multi-agents system for
 185 generic, end-to-end software engineering tasks.

186 3 Methodology

187 Figure 1 illustrates our framework. The key design
 188 of REPOPILOT is centered at the idea that only
 189 Planner requires advanced reasoning and it should
 190 dedicate easy but need long execution tasks into
 191 its child agents which can be small and fast, there-
 192 fore, reduce the inference cost and improve overall
 193 performance by eliminating redundant information
 194 out of context of the Planner. At each iteration,
 195 the Planner generates a thought consisting a plan
 196 to resolve a query, an agent request and receive a
 197 feedback from that agent from generated request.
 198 After receiving request from Planner, child agents
 199 which follow the design of ReAct agent (Yao et al.,

200 2022b), will iteratively execute tools, receiving
 201 feedback from codebase environment and report
 202 its findings to the Planner. Planner will terminate
 203 when the task is fully resolved.

204 3.1 Centralized Multi-Agent System

205 The central agent Planner controls three supporting
 206 agents Navigator, Editor and Executor:

207 The role of the Planner is to directly receive hu-
 208 man task prompts, such as resolving GitHub issues
 209 or generating code based on a given description,
 210 and address these tasks iteratively. In each iteration,
 211 it generates a plan to resolve the query, deciding
 212 which supporting agents to utilize and how to em-
 213 ploy them in that step (see Section 3.2 for details
 214 on how the Planner delegates task execution to
 215 the child agents and receives feedback from them).
 216 The Planner also determines when to terminate the
 217 resolving process or stops after a maximum number
 218 of iterations to minimize the costs.

219 Navigator is a crucial agent that has responsi-
 220 bility to receive information-seeking query from
 221 Planner then iteratively navigates and collects all
 222 relevant information that can be used to answer

the query. Navigator is equipped with tools similar to ones in Integrated development environment (IDE) commonly used by programmers which are `go_to_definition` or `code_search`. Given a query from *Planner*, navigator uses a series of these provided tools to quickly traverse in the codebase and collect information. The search process will be terminated by the agent or timeout is reached. However, the information-seeking process is costly, requiring a certain number of exploratory iterations for an agent to understand the repository and address the query. This is particularly challenging when the codebase is private and not publicly available, as the LLM lacks knowledge of its structure and key design elements.

Therefore, the navigation agent should be fast and lightweight. It is crucial to have multiple simple, easy-to-use tools that, when combined, yield excellent results for navigation and search. This contrasts with the design of SWE-agent (Yang et al.), which relies solely on summarized code search. While lightweight, this approach often fails to provide concise context understanding. Consequently, LLMs such as LLaMA (Touvron et al., 2023) or Claude 3 Haiku (Chiang et al., 2024) can struggle when faced with exhaustive search results from these tools.

The *Editor* is an agent responsible for editing code in multiple files or creating new files and generating source code for them. *Editor* is equipped with several tools, including `auto_repair_editor`, `code_search`, and `open_file`. It receives the target file from the *Planner* along with the editing context, which includes related information such as the expected behavior after the change and the specific lines to be considered. *Editor* then produces a replacement block of code, which the `auto_repair_editor` inserts into the target file. If the `auto_repair_editor` confirms the edit’s success, the *Editor* loop can be terminated, and control returns to the *Planner*.

The *Executor* is an agent responsible for reproducing issues or verifying that the generated patch behaves correctly and resolves the issue after changes are made. It is equipped with an `interactive_bash_shell`, which maintains its state even after executing commands, and `open_file` for reading installation instructions or documentation. The environment setup is also handled automatically by the *Executor*.

3.2 Agent Execution and Communication

This section focuses on how the *Planner* delegates repetitive, interactive but simple tasks to child agents and how these child agents communicate their findings back to the *Planner* with minimal information loss.

Firstly, *Planner* needs to provide detailed request with contextual information for child agents to perform its role and execution. The format of communication from *Planner* sending to other agents always have Context and Request fields. The Context field is more descriptive on original behavior and why *Editor* has to generate a patch while the Request field is a direct command to re-ensure code generated is behaved according to the intention of the *Planner*. Similarly, the Context field for *Navigator* is utilized for faster and direct navigation by optionally directing the target files or specific locations in advance.

```
Planner -> Generator
Request: Generate a patch for the __new__ method to strip the 'L' suffix from the string before converting it to an integer, ensuring compatibility between Python 2 and Python 3.
Context: The __new__ method in `sympy/core/numbers.py` handles long integers by converting hexadecimal strings to long integers using the long() function. This causes an error when unpickling in Python 3 because Python 3 does not recognize the 'L' suffix used in Python 2 for long integers.
```

Figure 2: Request from *Planner* to *Editor*

After receiving a request and context from the *Planner*, child agents iteratively execute their tasks to resolve the request and eventually report their final answer back to the *Planner*. However, we observed significant information loss in these final reports, especially from the navigation agent. The final answers often lack code snippets or detailed information about the explored objects or codebase structure, leading to deterioration and hallucination in the *Planner* after several iterations. To address this, we propose using a lightweight LLM summarizer to compile intermediate results from the child agent’s execution, incorporating this summary into the final answer communicated to the *Planner*.

3.3 Tool Design

A key distinction of REPOPILOT compared to existing approaches lies in tool design. Several factors are important when designing a tool for each agent role. First, the feedback format from the tools to the LLM must be succinct, informative, and readable. Second, the functionality of the tool should

be considered: does it provide the most information, help the LLM achieve a higher patch apply rate, or speed up the testing environment setup? Finally, usability is crucial; the input interface must be easy-to-use and intuitive. A complex input interface could lead to the LLM failing to use these tools correctly, and in some cases, incorrect tool outputs can negatively impact overall performance.

Here, we describe how tools are designed specifically for each agent role.

Navigating. In tasks like GitHub Issue Resolving or Code Generation, efficiently finding relevant function and class definitions is crucial. To address this, we implemented the `code_search` tool using a trigram-based code search engine with symbol ranking, `Zoekt` (Nienhuys and Source-Graph), as the code search backend. Additionally, we introduced IDE-like features such as `go_to_definition`, `get_all_references`, and `get_all_symbols`. These tools, when used in conjunction with `code_search`, are particularly powerful because `code_search` can produce many false positives, especially in large codebases. However, these results can indicate where the target definition is called or referenced, and `go_to_definition` can quickly localize the definition effortlessly. This approach contrasts sharply with the `summarized_search` in (Yang et al.), which presents exhaustive search results in one feedback. Similarly, `get_all_references` is effective in localizing similar usages of target functions or classes, helping the system understand how these objects are typically used, leading to better and more consistent fixes.

Designing input interfaces for these tools presents significant challenges. These tools require precise specifications, including exact line numbers, column numbers, file paths, and identifier names, to function properly. Our findings indicate that LLMs often struggle to accurately localize the column number and even the line number of a target word. Therefore, we implemented a proximity search algorithm that identifies the closest exact position relative to the LLM’s prediction. Consequently, the LLM now only needs to provide the ‘word_list’, ‘file_path’, and ‘line_number’.

We implemented the `get_tree_structure` function, which visualizes the structure of a codebase as a tree. This function takes `relative_path` and ‘level’ as inputs, displaying the folder and file structure within the specified `relative_path` up to

the given ‘level’ of depth. Moreover, we equipped Navigator with `open_file` and `find_file` functionalities. The `open_file` function can display up to 50 lines of source code, annotated with line numbers, and is integrated with a keyword search, allowing the LLM to locate sections containing specific keywords without exhaustive scrolling.

```

Action:
```
{
 "action": "open_file",
 "action_input": {
 "relative_file_path": "sympy/core/numbers.py",
 "start_line": 200,
 "end_line": 600,
 "keywords": ["__new__"]
 }
}
```

Observation: The content of sympy/core/numbers.py
is:
Results for keyword: __new__
524
525     # Used to make max(x._prec, y._prec) return
526     _prec = -1
527
528     def __new__(cls, *obj):
529         if len(obj) == 1:
530             obj = obj[0]

```

Figure 3: Open File Interface

`code_search` is the most intensively used tool for navigation, making its interface design critical. Besides using simple input specifications and a robust trigram-based search engine, we re-rank results to prioritize primary symbols like functions or classes. Additionally, we found that returned results should include a preview window. This window displays documentation for primary symbols or 5 lines surrounding the search position for other results.

Editing. We also equip the *Editor* with tools similar to those of the *Navigator*, such as `open_file` and `go_to_definition`, to observe the target editing file and jump to the definition of identifiers around the editing location, ensuring consistent fixes within the codebase. The primary focus of editing is the `auto_repair_editor` tool, which uses the ‘start_line’ and ‘end_line’ positions of the original block and ‘patch’ fields for the replacement code. This tool can receive a replacement patch and automatically detect and repair any indentation or syntax errors via an LLM.

Executing Executor can use `open_file` or `get_tree_structure` to find documentations related to testing instruction or environment setup. Then it can use `interactive_bash_shell` to run bash command, normally, these commands are

Model	SWE-bench		SWE-bench Lite	
	% Pass Rate	\$ Avg. Cost	% Pass Rate	\$ Avg. Cost
RAG				
w/ GPT-4 Turbo	1.31	0.13	2.67	0.13
w/ Claude 3 Opus	3.79	0.25	4.33	0.25
SWE-agent				
w/ GPT-4 Turbo	12.47	1.59	18.00	1.67
w/ Claude 3 Opus	10.46	2.59	13.00	2.18
AutoCodeRover*				
w/ GPT-4 Turbo	-	-	16.11	0.44
RepoPilot				
w/ GPT-4o	10.12	0.41	17.00	0.38
w/ Claude Opus	-	-	13.00	0.56
w/ WizardLM2	-	-	16.00	0.28

Table 1: Performance comparison of different methods on SWE-bench and SWE-bench Lite.

‘pytest’, ‘conda installation’ or ‘python3’. Notes that we uses interactive bash shell to maintain execution states of bash instead of using Docker.

4 Implementation Details

Benchmarks We evaluate REPOPILOT using the SWE-bench dataset (Jimenez et al., 2023), which comprises 2,294 task instances derived from 12 popular Python repositories, such as flask, numpy, and matplotlib. SWE-bench is designed to assess a system’s capability to automatically resolve GitHub issues by using Issue-Pull Request (PR) pairs. The evaluation process involves verifying unit tests using the post-PR behavior as the reference solution. For ablations and analysis, due to budget constraint, we can only utilize a small subset called SWE-bench Tiny consisting of 100 instances from SWE-bench Lite (Jimenez et al., 2023). Additionally, we also examine our method on a small, hand-collected and executable benchmark for repository-level code generation called ExecBench, this benchmark is similar with DevEval (Li et al., 2024) but with synthetic test cases. This benchmark aims for evaluating code generation at the repository level that emphasizes executability and correctness. The test cases are automatically generated by few-shot prompted GPT-3.5 Turbo with internal pipeline ensuring 96.25% test coverage. The benchmark includes 120 samples.

Models. For *summarizer* in REPOPILOT, we used Mixtral 8x7B (Jiang et al., 2024). We use a wide variety of LMs to examine the flexibility of the framework as well as measure robustness. For *Planner* role, advanced LLMs including GPT-4o (gpt-4o-2024-05-13) (OpenAI et al., 2023), Claude 3 Opus (claude-3-opus-20230229) (Chiang et al., 2024) and WizardLM2

(wizardlm-2-8x22b) (Tao et al., 2024) are used. For *Navigator*, we only use Claude 3 Haiku (claude-3-haiku-20240307 for fast inference time and low cost; GPT-4o, Claude 3 Sonnet (claude-3-sonnet-20240229) and WizardLM2 are used for *Editor* role. Finally, for *Executor*, GPT-3.5 Turbo (gpt-3.5-turbo-0125) and Claude 3 Haiku are utilized. All models and agent roles utilize greedy sampling during inference with temperature parameter set to 0. We have 3 different multi-agent configurations denoted as following GPT-4o with (GPT-4o, Claude Haiku, GPT-4o, GPT-3.5 Turbo) as (planner, navigator, editor and executor), similarly, WizardLM2 with (WizardLM2, Claude Haiku, WizardLM2, Claude Haiku) and Claude Opus with (Claude Opus, Claude Haiku, Claude Sonnet, Claude Haiku).

Baselines. On SWE-Bench, we compare REPOPILOT to three baselines: SWE-Agent, a bash interactive agent with Agent-Computer Interfaces (Yang et al.); AutoCodeRover, a two-stage agent pipeline focusing on bug fixing scenarios (Zhang et al., 2024c) and Retrieval Augmented Generation baselines shown in (Jimenez et al., 2023).

Metrics. We use pass@1 as our metric, which measures the percentage of instances where all tests pass successfully after applying the model-generated patch to the repository. Additionally, we report the Avg cost. Cost metric, representing the average API inference cost for all LLM usages. For ExecBench

Tools. go_to_def, get_all_references and get_all_symbols are implemented using Multilspy (Agrawal et al., 2024) with natural language interface for LLM. LLM inside auto_repair_editor is the same with *Editor* in

one configuration.

Others. For RAG, text-embedding-3-large is used for source code embeddings, we use code parsing and chunking implementation from Langchain¹. The chunking size is 1000 tokens.

5 Experimental Results

Methods	pass@1	pass@5
GPT-3.5-Turbo + RAG	24.16	35.00
WizardLM2 + RAG	32.50	49.16
RepoPilot w/ WizardLM2	38.33	53.33

Table 2: Pass Rates of different methods on ExecBench

As shown as in Table 1, on SWE-bench Lite, in our framework, GPT-4o configuration achieves best performance 17% (51/300) compared to Claude Opus and WizardLM2. Furthermore, it has higher resolve rate compared to AutoCodeRover (17% versus 16.11%) and competitive with SWE-agent while having significantly lower cost (~ 4.4 times lower). Notes that REPOPILOT has a more generic architecture rather than AutoCodeRover which is a two-stage pipeline focusing on bug resolving. Moreover, we found that an open weight model WizardLM2-8x22B shows strong planning capabilities and when combined with other LLMs inside our framework, it shows comparable results with GPT-4o but with only 73% API cost of GPT-4o configuration. Due to budget constraint, we can only evaluate GPT-4o configuration on full SWE-bench testset, and, it shows superior performance over RAG approaches and positive results compared with SWE-agent (10.12% versus 12.47%).

In Table 2, REPOPILOT exhibits strong performance on ExecBench, surpassing RAG baselines in both pass@1 and pass@5 metrics. This demonstrates its versatility in handling various software engineering tasks of REPOPILOT.

6 Analysis

6.1 Ablation Studies on Agent Roles

In this experiment, we use SWE-bench Tiny to evaluate the contribution of each agent role to overall performance by replacing each child agent with the planner itself. This requires the planner to directly call the eliminated agent’s toolset to execute necessary tasks. Table 3 shows a significant cost increase for all configurations when any agent role is

removed. The resolving rate also decreases, varying by which role is eliminated. Removing the *Navigator* causes the most performance drop, followed by the *Editor* and the *Executor*, respectively.

Notably, in the case of a medium-long context length LLM such as WizardLM2 acting as the *Planner* and replacing the role of *Editor* or *Navigator*, we observe a more severe drop in the resolving rate. This is because these roles require continuous interaction with the environment, necessitating a long context.

Model	SWE-bench Tiny		
	% Pass Rate	\$ Avg. Cost	
GPT-4o		15.00	0.42
	w/o Navigator.	7.00	2.81
	w/o Editor	11.00	1.92
	w/o Executor	14.00	0.75
WizardLM2		13.00	0.31
	w/o Navigator	4.00	1.21
	w/o Editor	7.00	0.51
	w/o Executor	13.00	0.38

Table 3: Ablation study on different agent role’s contribution on SWE-bench Tiny

6.2 Analysis of Tool Design

We investigate the extent of improvements brought by our major design choices in the tool’s interface and functionality. We conduct an ablation study on the functionalities of `go_to_definition`, `auto_repair_editor`, `open_file`, and `code_search` using SWE-bench Tiny. For each tool, we evaluate the overall performance when the tool is utilized versus when it is not, as shown in Table 4 (denoted as Used and No Usage, respectively).

A crucial finding for `go_to_definition` is that the LLM agent struggles to effectively use this IDE-like feature. It requires exact line and column num-

<code>go_to_definition</code>		<code>open_file</code>	
Used	9.00 _{↓6.0}	Used	9.00 _{↓6.0}
w/ search	15.00	w/ annotated lines	11.00 _{↓4.0}
No usage	12.0 _{↓3.0}	w/ keyword summary	15.00
		No usage	4.0 _{↓11.0}
<code>code_search</code>		<code>auto_repair_editor</code>	
Used	8.00 _{↓6.0}	Used	8.00 _{↓7.0}
w/ preview	11.00 _{↓3.0}	w/ linting feedback	11.00 _{↓4.0}
w/ ranking	14.00	w/ repairing	15.00
No usage	3.0 _{↓11.0}	No usage	1.0 _{↓14.0}

Table 4: Ablation result on resolving performance on SWE-Bench Tiny with different key tool designs

¹<https://github.com/langchain-ai/langchain>

535 bers and the precise symbol name, which demands
 536 precise localization of character positions. Despite
 537 supporting annotated line numbers, the agent often
 538 fails and retries multiple times. However, incorpor-
 539 ating a proximity-based search process, allowing
 540 the agent to approximate specifications, would sig-
 541 nificantly improve performance (from 9% without
 542 search to 15% with search).

543 For `open_file`, it's frequently to observe that
 544 small LLMs, such as Claude Haiku, tend to scroll
 545 up and down multiple times to find desired snip-
 546 pets by continuously increasing 'start_line' and
 547 'end_line', leading to out-of-context length issues.
 548 We solved this problem by adding an additional
 549 input field 'keywords', allowing LLM to search
 550 keywords inside file and tool to quickly localize
 551 the positions of keywords inside the file and display
 552 the surrounding lines. This increases the resolving
 553 rate by 3%.

554 Without the utilization of `code_search`, the *Nav-*
 555 *igator* would face significant challenges in swiftly
 556 identifying the necessary objects to fulfill a request,
 557 resulting in a significantly lower performance rate
 558 of 3% compared to 8% when the tool is employed.
 559 Furthermore, merely providing file paths and key-
 560 word positions from `code_search` would not be
 561 sufficient. Enhancing the output to include partial
 562 surrounding context around the keyword enabling
 563 Navigator to make more informed decisions and
 564 select the most relevant results, thereby improving
 565 performance from 8% to 11%. In the majority of
 566 Software Engineering Workflow (SWE) tasks, pro-
 567 grammers primarily focus on key objects such as
 568 functions and classes. Consequently, prioritizing
 569 search results for these objects over other results
 570 is beneficial, and re-ranking these results further
 571 enhances overall performance, increasing it from
 572 11% to 14%.

573 We further observe a substantial enhancement
 574 (8% to 11%) for providing python linting feed-
 575 back² to the *Editor* whenever it produces a patch
 576 via `auto_repair_tool`. By providing linting
 577 feedback for an internal LLM of the tool to au-
 578 tonomously refine the generated patch and attempt
 579 to fix any encountered error, the quality of the patch
 580 would be improved leading resolving rate increased
 581 by 4%

²We use `flake8` for providing syntax errors.

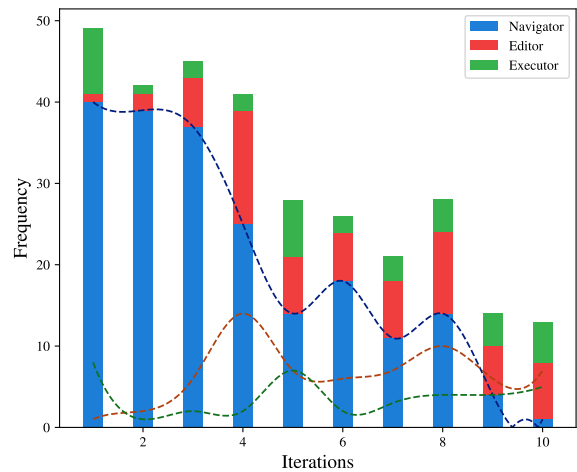


Figure 4: Agent Roles Usage Frequency in SWE-bench by *Planner*'s iterations

6.3 Agent Behavior

582 We analyzed the frequency of each agent role re-
 583 quested by the *Planner* throughout the issue reso-
 584 lution process. Figure 4 illustrates a typical pattern
 585 where the *Planner* is most active at the beginning
 586 of the resolution process, gathering relevant infor-
 587 mation about the codebase environment. Subse-
 588 quently, the *Editor* is frequently used to generate
 589 patches, often immediately following the *Naviga-*
 590 *tor*, with notable peaks at Iterations 4 and 8. Finally,
 591 the *Executor* is requested more frequently in the
 592 later iterations to verify the results by executing
 593 tests. It is noteworthy that, in the first iteration,
 594 there are a small peak indicating that *Executor* is
 595 requested to reproduce the issue.
 596

7 Conclusion

597 In this paper, we present REPOPILOT, a versatile
 598 multi-agent-based system designed to tackle
 599 complex software engineering tasks at the reposi-
 600 tory level. By carefully crafting the overall frame-
 601 work, communication mechanisms, and interactive
 602 tools, we achieve competitive performance com-
 603 pared to single-agent systems while maintaining
 604 significantly lower costs. Furthermore, our design
 605 is adaptable to other intelligent software engineer-
 606 ing tasks and various programming languages.
 607

608 Limitations

609 Despite the significant advancements demonstrated
610 in multi-agent software development, several limi-
611 tations and areas for future research remain to be
612 addressed for REPOPILOT.

613 One potential avenue for future research is to
614 apply REPOPILOT to more diversified software en-
615 gineering tasks, such as repository-level question
616 answering, which requires a deep understanding
617 of the codebase. REPOPILOT can leverage its ef-
618 fective tool suite to navigate and collect relevant
619 information efficiently.

620 Secondly, future work could focus on train-
621 ing smaller, open-sourced large language models
622 (LLMs) using trajectories collected from more ca-
623 pable yet expensive models like WizardLM2. Ef-
624 fectively training these smaller models presents a
625 significant challenge and an interesting research
626 problem.

627 Finally, our approach may result in lengthy ex-
628 ecution times for child agents due to the lack of
629 contextual information from previous rounds (as
630 the child agent’s context is reset after each request
631 from the *Planner*). Implementing a specialized
632 adaptive memory for each agent role could enhance
633 the efficiency and speed of subsequent executions.

634 References

635 Lakshya A Agrawal, Aditya Kanade, Navin Goyal, Shu-
636 vendu Lahiri, and Sriram Rajamani. 2024. Monitor-
637 guided decoding of code llms with static analysis of
638 repository context. *Advances in Neural Information
639 Processing Systems*, 36.

640 Jacob Austin, Augustus Odena, Maxwell Nye, Maarten
641 Bosma, Henryk Michalewski, David Dohan, Ellen
642 Jiang, Carrie Cai, Michael Terry, Quoc Le, and
643 Charles Sutton. 2021. [Program synthesis with large
644 language models](#).

645 Ramakrishna Bairi, Atharv Sonwane, Aditya Kanade,
646 Arun Iyer, Suresh Parthasarathy, Sriram Rajamani,
647 B Ashok, Shashank Shet, et al. 2023. Codeplan:
648 Repository-level coding using llms and planning.
649 *arXiv preprint arXiv:2309.12499*.

650 Federico Cassano, John Gouwar, Daniel Nguyen, Syd-
651 ney Nguyen, Luna Phipps-Costin, Donald Pinckney,
652 Ming-Ho Yee, Yangtian Zi, Carolyn Jane Anderson,
653 Molly Q Feldman, et al. 2022. Multipl-e: A scal-
654 able and extensible approach to benchmarking neural
655 code generation. *arXiv preprint arXiv:2208.08227*.

656 Saikat Chakraborty, Yujian Li, Matt Irvine, Ripon Saha,
657 and Baishakhi Ray. 2018. Entropy guided spec-
658 trum based bug localization using statistical language
659 model. *arXiv preprint arXiv:1802.06947*.

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming 660
Yuan, Henrique Ponde, Jared Kaplan, Harrison Ed- 661
wards, Yura Burda, Nicholas Joseph, Greg Brockman, 662
Alex Ray, Raul Puri, Gretchen Krueger, Michael 663
Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, 664
Brooke Chan, Scott Gray, Nick Ryder, Mikhail 665
Pavlov, Alethea Power, Lukasz Kaiser, Moham- 666
mad Bavarian, Clemens Winter, Philippe Tillet, Fe- 667
lipe Petroski Such, David W. Cummings, Matthias 668
Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel 669
Herbert-Voss, William H. Guss, Alex Nichol, Igor 670
Babuschkin, Suchir Balaji, Shantanu Jain, Andrew 671
Carr, Jan Leike, Joshua Achiam, Vedant Misra, Evan 672
Morikawa, Alec Radford, Matthew M. Knight, Miles 673
Brundage, Mira Murati, Katie Mayer, Peter Welinder, 674
Bob McGrew, Dario Amodei, Sam McCandlish, Ilya 675
Sutskever, and Wojciech Zaremba. 2021. [Evaluat-
676 ing large language models trained on code](#). *ArXiv*,
677 [abs/2107.03374](#). 678

Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anasta- 679
sios Nikolas Angelopoulos, Tianle Li, Dacheng Li, 680
Hao Zhang, Banghua Zhu, Michael Jordan, Joseph E 681
Gonzalez, et al. 2024. Chatbot arena: An open plat- 682
form for evaluating llms by human preference. *arXiv
683 preprint arXiv:2403.04132*. 684

Yanguibo Ding, Zijian Wang, Wasi Ahmad, Hantian 685
Ding, Ming Tan, Nihal Jain, Murali Krishna Ra- 686
manathan, Ramesh Nallapati, Parminder Bhatia, Dan 687
Roth, et al. 2024. Crosscodeeval: A diverse and mul- 688
tilingual benchmark for cross-file code completion. 689
Advances in Neural Information Processing Systems, 690
36. 691

Xueying Du, Mingwei Liu, Kaixin Wang, Hanlin Wang, 692
Junwei Liu, Yixuan Chen, Jiayi Feng, Chaofeng 693
Sha, Xin Peng, and Yiling Lou. 2023. Classe- 694
val: A manually-crafted benchmark for evaluating 695
llms on class-level code generation. *arXiv preprint
696 arXiv:2308.01861*. 697

Yu Gu, Yiheng Shu, Hao Yu, Xiao Liu, Yuxiao Dong, 698
Jie Tang, Jayanth Srinivasa, Hugo Latapie, and Yu Su. 699
2024. Middleware for llms: Tools are instrumental 700
for language agents in complex environments. *arXiv
701 preprint arXiv:2402.14672*. 702

Dan Hendrycks, Steven Basart, Saurav Kadavath, Man- 703
tas Mazeika, Akul Arora, Ethan Guo, Collin Burns, 704
Samir Puranik, Horace He, Dawn Song, et al. 2021. 705
Measuring coding challenge competence with apps. 706
arXiv preprint arXiv:2105.09938. 707

Albert Q. Jiang, Alexandre Sablayrolles, Antoine 708
Roux, Arthur Mensch, Blanche Savary, Chris 709
Bamford, Devendra Singh Chaplot, Diego de las 710
Casas, Emma Bou Hanna, Florian Bressand, Gi- 711
anna Lengyel, Guillaume Bour, Guillaume Lam- 712
ple, L elio Renard Lavaud, Lucile Saulnier, Marie- 713
Anne Lachaux, Pierre Stock, Sandeep Subramanian, 714
Sophia Yang, Szymon Antoniak, Teven Le Scao, 715
Th eophile Gervet, Thibaut Lavril, Thomas Wang, 716
Timoth e Lacroix, and William El Sayed. 2024. [Mix-
717 tural of experts](#). 718

719	Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R Narasimhan. 2023. Swe-bench: Can language models resolve real-world github issues? In <i>The Twelfth International Conference on Learning Representations</i> .	777
720		778
721		779
722		780
723		781
724		782
725	René Just, Darioush Jalali, and Michael D Ernst. 2014. Defects4j: A database of existing faults to enable controlled testing studies for java programs. In <i>Proceedings of the 2014 international symposium on software testing and analysis</i> , pages 437–440.	783
726		784
727		785
728		786
729		787
730	Rafael-Michael Karampatsis and Charles Sutton. 2020. How often do single-statement bugs occur? the manysstubs4j dataset. In <i>Proceedings of the 17th International Conference on Mining Software Repositories</i> , pages 573–577.	788
731		789
732		790
733		
734		
735	Yuhang Lai, Chengxi Li, Yiming Wang, Tianyi Zhang, Ruiqi Zhong, Luke Zettlemoyer, Wen-tau Yih, Daniel Fried, Sida Wang, and Tao Yu. 2023. Ds-1000: A natural and reliable benchmark for data science code generation. In <i>International Conference on Machine Learning</i> , pages 18319–18345. PMLR.	791
736		792
737		793
738		794
739		795
740		796
741	Jia Li, Ge Li, Yunfei Zhao, Yongmin Li, Huanyu Liu, Hao Zhu, Lecheng Wang, Kaibo Liu, Zheng Fang, Lanshen Wang, Jiazheng Ding, Xuanming Zhang, Yuqi Zhu, Yihong Dong, Zhi Jin, Binhua Li, Fei Huang, and Yongbin Li. 2024. Deveval: A manually-annotated code generation benchmark aligned with real-world code repositories .	797
742		798
743		799
744		800
745		801
746		802
747		803
748	Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and LINGMING ZHANG. 2023a. Is your code generated by chatgpt really correct? rigorous evaluation of large language models for code generation. In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> .	804
749		805
750		806
751		807
752		808
753		809
754	Tianyang Liu, Canwen Xu, and Julian McAuley. 2023b. Repobench: Benchmarking repository-level code auto-completion systems. In <i>The Twelfth International Conference on Learning Representations</i> .	810
755		811
756		812
757		813
758	Yuliang Liu, Xiangru Tang, Zefan Cai, Junjie Lu, Yichi Zhang, Yanjun Shao, Zexuan Deng, Helan Hu, Zengxian Yang, Kaikai An, et al. 2023c. Ml-bench: Large language models leverage open-source libraries for machine learning tasks. <i>arXiv preprint arXiv:2311.09835</i> .	814
759		815
760		816
761		817
762		818
763		819
764	Shuai Lu, Daya Guo, Shuo Ren, Junjie Huang, Alexey Svyatkovskiy, Ambrosio Blanco, Colin Clement, Dawn Drain, Daxin Jiang, Duyu Tang, et al. 2021. Codexglue: A machine learning benchmark dataset for code understanding and generation. <i>arXiv preprint arXiv:2102.04664</i> .	820
765		821
766		822
767		823
768		824
769		825
770	Niklas Muennighoff, Qian Liu, Armel Randy Zebaze, Qinkai Zheng, Binyuan Hui, Terry Yue Zhuo, Swayam Singh, Xiangru Tang, Leandro Von Werra, and Shayne Longpre. 2023. Octopack: Instruction tuning code large language models. In <i>The Twelfth International Conference on Learning Representations</i> .	826
771		827
772		828
773		829
774		830
775		831
776		832
		833
		834
		835
		836
		837
	Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. 2021. Webgpt: Browser-assisted question-answering with human feedback. <i>arXiv preprint arXiv:2112.09332</i> .	
	King Han Naman Jain, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. 2024. Live-codebench: Holistic and contamination free evaluation of large language models for code. <i>arXiv preprint arXiv:2403.07974</i> .	
	Han-Wen Nienhuys and SourceGraph. Zoekt. https://github.com/sourcegraph/zoekt .	
	OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan,	

838	Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeesh Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. 2023. Gpt-4 technical report .	
839		
840		
841		
842		
843		
844		
845		
846		
847		
848		
849		
850		
851		
852		
853		
854		
855		
856		
857		
858		
859		
860		
861		
862		
863		
864		
865		
866		
867		
868		
869		
870		
871		
872		
873		
874		
875		
876		
877		
878	Minghao Shao, Boyuan Chen, Sofija Jancheska, Brendan Dolan-Gavitt, Siddharth Garg, Ramesh Karri, and Muhammad Shafique. 2024. An empirical evaluation of llms for solving offensive security challenges. <i>arXiv preprint arXiv:2402.11814</i> .	
879		
880		
881		
882		
883	Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2024. Reflexion: Language agents with verbal reinforcement learning. <i>Advances in Neural Information Processing Systems</i> , 36.	
884		
885		
886		
887		
888	Theodore R Sumers, Shunyu Yao, Karthik Narasimhan, and Thomas L Griffiths. 2023. Cognitive architectures for language agents. <i>arXiv preprint arXiv:2309.02427</i> .	
889		
890		
891		
892	Zhengwei Tao, Ting-En Lin, Xiancai Chen, Hangyu Li, Yuchuan Wu, Yongbin Li, Zhi Jin, Fei Huang, Dacheng Tao, and Jingren Zhou. 2024. A survey on self-evolution of large language models. <i>arXiv preprint arXiv:2404.14387</i> .	
893		
894		
895		
896		
	Amitayush Thakur, Yeming Wen, and Swarat Chaudhuri. 2023. A language-agent approach to formal theorem-proving. <i>arXiv preprint arXiv:2310.04353</i> .	897
		898
		899
	Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, et al. 2022. Lamda: Language models for dialog applications. <i>arXiv preprint arXiv:2201.08239</i> .	900
		901
		902
		903
		904
	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models .	905
		906
		907
		908
		909
		910
		911
		912
		913
		914
		915
		916
		917
		918
		919
		920
		921
		922
		923
		924
		925
		926
		927
	Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, et al. 2024a. A survey on large language model based autonomous agents. <i>Frontiers of Computer Science</i> , 18(6):1–26.	928
		929
		930
		931
		932
	Xingyao Wang, Yangyi Chen, Lifan Yuan, Yizhe Zhang, Yunzhu Li, Hao Peng, and Heng Ji. 2024b. Executable code actions elicit better llm agents. <i>arXiv preprint arXiv:2402.01030</i> .	933
		934
		935
		936
	Zhiruo Wang, Grace Cuenca, Shuyan Zhou, Frank F Xu, and Graham Neubig. 2022. Mconala: a benchmark for code generation from multiple natural languages. <i>arXiv preprint arXiv:2203.08388</i> .	937
		938
		939
		940
	Zhiruo Wang, Daniel Fried, and Graham Neubig. 2024c. Trove: Inducing verifiable and efficient toolboxes for solving programmatic tasks. <i>arXiv preprint arXiv:2401.12869</i> .	941
		942
		943
		944
	Zhiruo Wang, Shuyan Zhou, Daniel Fried, and Graham Neubig. 2023. Execution-based evaluation for open-domain code generation. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 1271–1290.	945
		946
		947
		948
		949
	Zhiyong Wu, Chengcheng Han, Zichen Ding, Zhenmin Weng, Zhoumianze Liu, Shunyu Yao, Tao Yu, and Lingpeng Kong. 2024. Os-copilot: Towards generalist computer agents with self-improvement. <i>arXiv preprint arXiv:2402.07456</i> .	950
		951
		952
		953
		954

955	Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, et al. 2023. The rise and potential of large language model based agents: A survey. <i>arXiv preprint arXiv:2309.07864</i> .	decompositions. <i>Advances in Neural Information Processing Systems</i> , 36:31466–31523.	1011
956			1012
957			
958			
959			
960	Chunqiu Steven Xia and Lingming Zhang. 2022. Less training, more repairing please: revisiting automated program repair via zero-shot learning. In <i>Proceedings of the 30th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering</i> , pages 959–971.	Eric Zelikman, Eliana Lorch, Lester Mackey, and Adam Tauman Kalai. 2023b. Self-taught optimizer (stop): Recursively self-improving code generation. In <i>OPT 2023: Optimization for Machine Learning</i> .	1013
961			1014
962			1015
963			1016
964			
965			
966	Tianbao Xie, Danyang Zhang, Jixuan Chen, Xiaochuan Li, Siheng Zhao, Ruisheng Cao, Toh Jing Hua, Zhoujun Cheng, Dongchan Shin, Fangyu Lei, et al. 2024. Osworld: Benchmarking multimodal agents for open-ended tasks in real computer environments. <i>arXiv preprint arXiv:2404.07972</i> .	Fengji Zhang, Bei Chen, Yue Zhang, Jacky Keung, Jin Liu, Daoguang Zan, Yi Mao, Jian-Guang Lou, and Weizhu Chen. 2023. Repocoder: Repository-level code completion through iterative retrieval and generation.	1017
967			1018
968			1019
969			1020
970			1021
971			
972	Aidan ZH Yang, Claire Le Goues, Ruben Martins, and Vincent Hellendoorn. 2024a. Large language models for test-free fault localization. In <i>Proceedings of the 46th IEEE/ACM International Conference on Software Engineering</i> , pages 1–12.	Kechi Zhang, Jia Li, Ge Li, Xianjie Shi, and Zhi Jin. 2024a. Codeagent: Enhancing code generation with tool-integrated agent systems for real-world repo-level coding challenges. <i>arXiv preprint arXiv:2401.07339</i> .	1022
973			1023
974			1024
975			1025
976			1026
977	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.	Shaokun Zhang, Jieyu Zhang, Jiale Liu, Linxin Song, Chi Wang, Ranjay Krishna, and Qingyun Wu. 2024b. Training language model agents without modifying language models. <i>arXiv preprint arXiv:2402.11359</i> .	1027
978			1028
979			1029
980			1030
981	John Yang, Akshara Prabhakar, Karthik Narasimhan, and Shunyu Yao. 2024b. Intercode: Standardizing and benchmarking interactive coding with execution feedback. <i>Advances in Neural Information Processing Systems</i> , 36.	Yuntong Zhang, Haifeng Ruan, Zhiyu Fan, and Abhik Roychoudhury. 2024c. Autocoderover: Autonomous program improvement.	1031
982			1032
983			1033
984			
985			
986	John Yang, Akshara Prabhakar, Shunyu Yao, Kexin Pei, and Karthik R Narasimhan. 2023. Language agents as hackers: Evaluating cybersecurity skills with capture the flag. In <i>Multi-Agent Security Workshop@ NeurIPS'23</i> .	Andy Zhou, Kai Yan, Michal Shlapentokh-Rothman, Haohan Wang, and Yu-Xiong Wang. 2023. Language agent tree search unifies reasoning acting and planning in language models. <i>arXiv preprint arXiv:2310.04406</i> .	1034
987			1035
988			1036
989			1037
990			1038
991	Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. 2022a. Webshop: Towards scalable real-world web interaction with grounded language agents. <i>Advances in Neural Information Processing Systems</i> , 35:20744–20757.		
992			
993			
994			
995			
996	Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan Cao. 2022b. React: Synergizing reasoning and acting in language models. In <i>The Eleventh International Conference on Learning Representations</i> .		
997			
998			
999			
1000			
1001	Hao Yu, Bo Shen, Dezhi Ran, Jiaxin Zhang, Qi Zhang, Yuchi Ma, Guangtai Liang, Ying Li, Qianxiang Wang, and Tao Xie. 2024. Codereval: A benchmark of pragmatic code generation with generative pre-trained models. In <i>Proceedings of the 46th IEEE/ACM International Conference on Software Engineering</i> , pages 1–12.		
1002			
1003			
1004			
1005			
1006			
1007			
1008	Eric Zelikman, Qian Huang, Gabriel Poesia, Noah Goodman, and Nick Haber. 2023a. Parsel: Algorithmic reasoning with language models by composing		
1009			
1010			

A.1 Prompt Template for Planner

Instruction Prompt Templates for Planner

Suffix Prompt:

Begin! Reminder to ALWAYS respond with a valid json blob of a single action. Use tools if
 ↪ necessary. Respond directly if you have gathered enough information from the repository.
 ↪ Format is Action:```\$JSON_BLOB```then Observation:. Thought:

Prefix Prompt:

You are a great developer with expertise in resolving programmer's query. You have been
 ↪ assigned a task to resolve a programmer's issue in a large repository. Devise a detailed
 ↪ plan using other language model agents to resolve the issue.

You have access into N agents, utilize them to step-by-step solve the query. Each consequent
 ↪ steps should be strictly based on the previous steps. Your thought process should be
 ↪ grounded by information collected from your agents, consider its results carefully, and
 ↪ make a decision based on the results and thought process.

Output the agent you want to use and the request you want to make to the agent. Respond
 ↪ directly and terminated=true if you have resolved the issue (code generated is verified
 ↪ and correct).

If you want to modify the logic of the code, or resolve the issue based on retrieved facts from
 ↪ code navigator, use code editor agent. Terminate if your task is executed successfully and
 ↪ reviewed correctly.

Top Priorities:

1. You need to use Codebase Navigator and Code Generator agents to resolve the issue. Use
 ↪ them at least once.
2. Maintain in mind which files need to edit.
3. Do not repeat your actions!. After receiving the response from the agent, diversify your
 ↪ next action to get more information.
4. Identify crucial causes of the issue, localize where the problem is before choosing the
 ↪ code generator agent.
5. No need to edit test file or test the code. You only need to resolve the issue in the
 ↪ codebase.
6. Do not care about any Pull Request or Existing Issue in the repository. You are only
 ↪ focused on the issue assigned to you.
7. Give a detailed request to the agent, so that the agent can understand the context of
 ↪ the query as well.

Important Notes:

1. Reading the issue description and understanding the problem is the first step. Make sure
 ↪ to identify the key components of the issue and the expected behavior.
2. Reading the response from the agents carefully, think about the information you have
 ↪ collected and how it can be used to resolve the issue.
3. Your thought process is the most important part of this task. Make sure to provide a
 ↪ detailed explanation of your reasoning with the issue, code snippets and relevant
 ↪ information collected from the agents.
4. Stop the task when you have resolved the issue. (Final Answer)

\$THOUGHT_PROCESS is your thought process about the query and previous results.
 You have access description to the following agents: {formatted_agents}

A.2 Prompt Template for Navigator

Instruction Prompt Templates for Navigator

Suffix Prompt:

Begin! Reminder to ALWAYS respond with a valid json blob of a single action. Use tools if
 ↪ necessary. Respond directly if you have gathered enough information from the repository.
 ↪ Format is Action:````\$JSON_BLOB````then Observation:. Thought:

Prefix Prompt:

You are an expert in finding all relevant information insider a code repository to answer the
 ↪ query from a planner agent. You have the full access to the codebase of the project you're
 ↪ working on to resolve a query from a planner. Your tools help you navigate the codebase and
 ↪ find relevant information. Use them wisely to explore the repository and find the
 ↪ information you need to resolve the query.
 If your first attempts do not provide enough information to resolve the query, try different
 ↪ tools or use tool with different parameters to get the information you need.
 Think carefully before making a decision. Your tools can provide valuable insights to help you
 ↪ resolve the query. Once you have collected relevant information, you can provide a
 ↪ response to the query with Final Answer, put any code snippet into that summary. Do not
 ↪ repeat your actions.

Top Priorities:

1. Understanding the query, think step-by-step carefully before decision to propose actions
 ↪ to collect necessary information to resolve the query.
2. Do not repeat your actions.
3. Try to use the tools to get the information you need. DO NOT GUESS or refuse to response
 ↪ the planner's request. Planner request is always correct. You may only see part of the
 ↪ information, but the planner sees the whole picture.
4. If one tool does not find the information you need, try another tool. If you open a
 ↪ file, but do not find the information you need, reopen with different start_line and
 ↪ end_line or keywords.
5. Your summarization should be relevant to the query (provide code snippet if it's
 ↪ required by query), do not provide unnecessary information.

Important Notes:

1. Try to combine different tools to seek related information to the query inside the
 ↪ project
2. get_all_references: Use this tool to get all references to a symbol in the codebase.
 ↪ This will help you understand how the symbol is used in the codebase.
3. get_all_symbols: Use this tool to get all symbols in the target file, it should be used
 ↪ with a keyword.
4. get_folder_structure: Use this tool to get the structure of the target folder. This will
 ↪ help you understand the organization of the codebase, and find the relevant files to
 ↪ use other tools.
5. code_search: Use this tool to search for symbol name if you know the exact name of the
 ↪ symbol, this is useful to find the definition if you're not familiar with codebase yet.
6. go_to_definition: Use this tool to navigate to the definition of an identifier, for
 ↪ example self._print in a class. (single word only, not a combination like sympy.latex),
 ↪ in this case, _print.
7. open_file: Use this tool to open a file in the codebase, this is useful to read the
 ↪ partial content of the file (50 lines). Should be used with a keyword (single word
 ↪ only, not a combination like `sympy.latex` just `latex` only) or limited start_line and
 ↪ end_line. If your previous open does not show all the information, next turn you can
 ↪ open the same file with different start_line and end_line (incrementally scrolling).

You have access to the following tools: {formatted_tools}

Instruction Prompt Templates for Editor

Suffix Prompt:

Begin! Reminder to ALWAYS respond with a valid json blob of a single action. Use tools if
 ↪ necessary. Respond directly if you have gathered enough information from the repository.
 ↪ Format is Action:```\$JSON_BLOB```then Observation:. Thought:

Prefix Prompt:

You are an expert in programming, you're excellent at editing or generate source code files.

Top priorities:

1. Always think step-by-step carefully before decision (Thought:).
2. DO NOT re-generate the same failed edit. Running it again will lead to the same error.
 ↪ Edit the file again if necessary based on the error message.
3. You should not edit lines locally one part by one part. You might need to remain the
 ↪ original code in your start_line and end_line if it's not changed (regenerate the code
 ↪ if necessary).
4. You should aware about the code context, and surrounding variables and functions. Do not
 ↪ add some undefined variables.

Important notes:

1. Always open the file before editing to see latest code!.
2. If you have to edit the code, ensure the code is correct with syntax, dependency,
 ↪ consistent with the file and the codebase.
3. Returning final answer if your code is successfully applied. You first can open a file to
 ↪ get the context and use editor to apply the patch. You have to generate code that is
 ↪ correct with syntax, ensure the dependency, consistent with the file and the codebase.
4. If you have the exact name of the file and symbol to edit, you can use the code search
 ↪ to find the definition of the symbol. If there's no definition, you can use open
 ↪ open_file tool.
5. Pay attention to original indentation.
6. The patch should be a block of code that be replaced into the code.

You have access to the following tool: {formatted_tools}

A.4 Prompt Template for Executor

Instruction Prompt Templates for Executor

Suffix Prompt:

Begin! Reminder to ALWAYS respond with a valid json blob of a single action. Use tools if
 ↪ necessary. Respond directly if you have gathered enough information from the repository.
 ↪ Format is Action:``\$JSON_BLOB``then Observation:. Thought:

Prefix Prompt:

You are an expert in running bash commandlines, based on the requests, try to run commands or
 ↪ files.
 If your request is not specified, considering setup the environment first (using conda create),
 ↪ cd into the project path and pip3 install -e .[dev]. Then find test folder to run all the
 ↪ tests. Think carefully before making a decision.

Remember these things:

1. Your environment code name is `repopilot`, create it if it doesn't exist. Activate it
 ↪ before running the tests.
2. No need to create new environments.
3. Check the requirements or setup instruction via opening docs files.
4. You should setup the developing environment before running the tests, for example
 ↪ running pip3 install -e . or something similar.

Important Tips:

1. When you observe (y/[n]) in the terminal, you should respond only with y or n.
 ↪ Similarly, when you see (yes/[no]), you should respond with yes or no.
2. You're already in the project directory.
3. Run command one by one.

You have access into following tools: {formatted_tools}

B RepoPilot trajectories on SWE-bench

Figure 5: Trajectory of RepoPilot for instance django_django-16041.

Issue	django_django-16041
<p>You need to identify the cause of the following github issue, collect the relevant information, and provide a solution. Github Issue:</p> <p>When explicitly setting <code>form_kwargs = {'empty_permitted':True}</code> or <code>form_kwargs =</code> ↪ <code>{'empty_permitted':False}</code> , a <code>KeyError</code> occurs when rendering a template that uses a ↪ <code>formset</code>'s <code>empty_form</code>.</p> <p>Expected Behavior <code>empty_permitted</code> is ignored for <code>formset.empty_form</code> since <code>empty_permitted</code> is irrelevant for ↪ <code>empty_form</code>, as <code>empty_form</code> is not meant to be used to pass data and therefore does not need ↪ to be validated.</p> <p>Steps to Reproduce</p> <pre> from django.shortcuts import render from .models import MyModel def test_view(request): context = {} ff = modelformset_factory(MyModel, fields = ['a_field']) context['formset'] = ff(queryset = MyModel.objects.none(), form_kwargs = {'empty_permitted':True} # or form_kwargs = ↪ {'empty_permitted':False}) return render(request, 'my_app/my_model_formset.html', context) # urls.py from django.urls import path, include from .views import test_view urlpatterns = [</pre>	


```

    path('test', test_view)
]
# my_model_formset.html
{% extends "my_app/base.html" %}
{% block content %}
<form id="my-form" method="post">
  {% csrf_token %}
  {{ formset }}
  <input type="submit" value="Save">
</form>
{{ formset.empty_form }}
{% endblock %}

```

INSTRUCTIONS: (31 lines)

```

(Open file: n/a)
(Current directory: /astropy__astropy)
bash-\$

```

1050

REPOPILOT

Planner

1. **Understanding the Issue**:
 - The issue is a `KeyError` when rendering a template that uses a `formset`'s `empty_form` with `form_kwargs` containing `empty_permitted`.
 - The expected behavior is that `empty_permitted` should be ignored for `empty_form` since it is irrelevant for `empty_form`.
2. **Identifying Key Components**:
 - The issue is related to `formset` and `empty_form`.
 - The `form_kwargs` parameter is causing the `KeyError`.
3. **Steps to Investigate**:
 - Locate the implementation of `modelformset_factory` and `BaseModelFormSet` in the Django codebase.
 - Investigate how `form_kwargs` is processed and how `empty_form` is generated.
4. **Plan**:
 - Use the `codebase_navigator` to locate the relevant files and functions.
 - Use the `code_generator` to create a patch if necessary.

Action:

```

{
  "action": "codebase_navigator",
  "action_input": {
    "context": "Locate modelformset_factory and BaseModelFormSet",
    "request": "Find the implementation of modelformset_factory and BaseModelFormSet in the
    ↪ Django codebase."
  }
}

```

1051

REPOPILOT

Navigator

Entering new execution chain..

Thought: To find the implementation of `modelformset_factory` and

`BaseModelFormSet` in the Django codebase, I will first search for the definitions of these symbols.

Action:

```

{

```

1052

```

    "action": "code_search",
    "action_input": {
        "names": ["modelformset_factory", "BaseModelFormSet"]
    }
}
}
}

```

1053

Observation

Code Search

Results for modelformset_factory:

=====

File: django/forms/models.py

Name: modelformset_factory

Documentation: Return a FormSet class for the given Django model class.

Results for BaseModelFormSet:

=====

File: django/forms/models.py

Name: BaseModelFormSet

Documentation:

A ``FormSet`` for editing a queryset and/or adding new objects to it.

File: django/forms/models.py

Implementation:

32 "ModelMultipleChoiceField",

33 "ALL_FIELDS",

34 "BaseModelFormSet",

35 "modelformset_factory",

File: django/forms/models.py

Implementation:

644

645

646 class BaseModelFormSet(BaseFormSet):

647 """

File: django/forms/models.py

Implementation:

1000 form=ModelForm,

1001 formfield_callback=None,

1002 formset=BaseModelFormSet,

1003 extra=1,

File: django/forms/models.py

Implementation:

1067

1068

1069 class BaseInlineFormSet(BaseModelFormSet):

1070 """A formset for child objects related to a parent."""

1054

REPOPILOT

Navigator

Thought: The search results indicate that the implementation of `modelformset_factory` and `BaseModelFormSet` is located in the `django/forms/models.py`

file. I will now open that file to explore the details of these implementations.

Action:

{

"action": "open_file",

1055

```

"action_input": {
    "relative_file_path": "django/forms/models.py",
    "start_line": 644,
    "end_line": 684
}
}
...

```

1056

Observation

[Open File](#)

```

The content of django/forms/models.py is:
644
645
646 class BaseModelFormSet(BaseFormSet):
647     """
648     A ``FormSet`` for editing a queryset and/or adding new objects to it.
649     """
650
651     model = None
652     edit_only = False
653
654     # Set of fields that must be unique among forms of this set.
655     unique_fields = set()
656
657     def __init__(
658         self,
659         data=None,
660         files=None,
661         auto_id="id_%s",
662         prefix=None,
663         queryset=None,
664         *,
665         initial=None,
666         **kwargs,
667     ):
668         self.queryset = queryset
669         self.initial_extra = initial
670         super().__init__(
671             **{
672                 "data": data,
673                 "files": files,
674                 "auto_id": auto_id,
675                 "prefix": prefix,
676                 **kwargs,
677             }
678         )
679
680     def initial_form_count(self):
681         """Return the number of forms that are required in this FormSet."""
682         if not self.is_bound:
683             return len(self.get_queryset())

```

1057

REPOPILOT

[Navigator](#)

Okay, let's summarize the key information about the implementation of `modelformset_factory` and `BaseModelFormSet` in the Django codebase:

Final Answer:

The implementation of `modelformset_factory` and `BaseModelFormSet` is located in the `django/forms/models.py` file.

1058

The `modelformset_factory` function is defined in this file and is used to return a `FormSet` class for the given Django model class. It takes several parameters such as the model class, the form class, and various options to customize the formset behavior.

The `BaseModelFormSet` class is also defined in this file. It is a subclass of `BaseFormSet` and is used for editing a queryset and/or adding new objects to it. It has several properties and methods that are specific to model formsets, such as `queryset`, `initial_extra`, and `initial_form_count()`.

The key parts of the `BaseModelFormSet` implementation are:

```
```python
class BaseModelFormSet(BaseFormSet):
 """
 A ``FormSet`` for editing a queryset and/or adding new objects to it.
 """

 model = None
 edit_only = False

 # Set of fields that must be unique among forms of this set.
 unique_fields = set()

 def __init__(
 self,
 data=None,
 files=None,
 auto_id="id_%s",
 prefix=None,
 queryset=None,
 *,
 initial=None,
 **kwargs,
):
 self.queryset = queryset
 self.initial_extra = initial
 super().__init__(
 **{
 "data": data,
 "files": files,
 "auto_id": auto_id,
 "prefix": prefix,
 **kwargs,
 }
)

 def initial_form_count(self):
 """Return the number of forms that are required in this FormSet."""
 if not self.is_bound:
 return len(self.get_queryset())
...
```
```

This provides the core functionality for a model formset, including managing the queryset, initial data, and number of forms in the formset.

REPOPILOT**Planner**

Based on the information gathered, the issue seems to be related to how `form_kwargs` are
 ↪ processed when creating the `empty_form` in a formset. Specifically, the `empty_permitted`
 ↪ key in `form_kwargs` is causing a `KeyError` when rendering the `empty_form`.

To resolve this issue, we need to ensure that `empty_permitted` is ignored for `empty_form`
 ↪ since it is irrelevant for `empty_form`. This can be done by modifying the relevant part
 ↪ of the code where `empty_form` is created to exclude `empty_permitted` from `form_kwargs`.

Next, I will use the `codebase_navigator` to locate the part of the code where `empty_form` is
 ↪ created and `form_kwargs` are processed in the `BaseModelFormSet` class.

Action:

```

---
{
  "action": "codebase_navigator",
  "action_input": {
    "request": "Find the part of the code in BaseModelFormSet where empty_form is created and
    ↪ form_kwargs are processed.",
    "context": "BaseModelFormSet is in `django/forms/models.py`"
  }
}
---
```

1060

REPOPILOT**Navigator**

Thought: To find the part of the code in BaseModelFormSet where empty_form is created and
 ↪ form_kwargs are processed, I will use the following steps:

1. Use `get_folder_structure` to understand the overall structure of the of `django/forms` and
 ↪ identify the relevant file. Planner suggests we should look into `django/forms/models.py`
2. Use `open_file` to browse the contents of the relevant file and locate the code that creates
 ↪ the `empty_form` and processes `form_kwargs`.
3. Use `go_to_definition` to navigate to the definition of the relevant methods or attributes.
4. Use `find_all_references` to find all references to the relevant methods or attributes.

Action:

```

---
{
  "action": "get_folder_structure",
  "action_input": {
    "relative_path": "django/forms"
  }
}
---
```

1061

Observation**Get Folder Structure**

The tree structure of `django/forms` is:

```

forms
__init__.py
jinja2
```

1062

```
templates
boundfield.py
fields.py
forms.py
formsets.py
models.py
renderers.py
utils.py
widgets.py
```

1063

REPOPILOT

Navigator

Thought: Okay, let's continue the process to find the relevant code in `BaseModelFormSet` and `empty_form`.

```
Action:
...
{
  "action": "open_file",
  "action_input": {
    "relative_file_path": "django/forms/formsets.py",
    "start_line": 1,
    "end_line": 20,
    "keywords": ["BaseModelFormSet", "empty_form"]
  }
}
...
```

1064

Observation

Open File

```
The content of django/forms/formsets.py is:
Results for keyword: empty_form
252     """Return a list of all the initial forms in this formset."""
253     return self.forms[: self.initial_form_count()]
254
255     @property
256     def extra_forms(self):
257         """Return a list of all the extra forms in this formset."""
258         return self.forms[self.initial_form_count() :]
259
260     @property
261     @property
262     def empty_form(self):
263         form_kwargs = self.get_form_kwargs(None)
264         form_kwargs.pop('empty_permitted', None)
265         form = self.form(
266             auto_id=self.auto_id,
267             prefix=self.add_prefix("__prefix__"),
268             empty_permitted=True,
269             use_required_attribute=False,
270             **form_kwargs,
271             renderer=self.renderer,
399
400     def full_clean(self):
401         """
402         Clean all of self.data and populate self._errors and
403         self._non_form_errors.
404         """
405         self._errors = []
406         self._non_form_errors = self.error_class(
407             error_class="nonform", renderer=self.renderer
408         )
409         empty_forms_count = 0
410
```

1065

```

411     if not self.is_bound: # Stop further processing.
412         return
413
414     if not self.management_form.is_valid():
415         error = ValidationError(
416             self.error_messages["missing_management_form"],
417             params={
418                 "field_names": ", ".join(
420                     for field_name in self.management_form.errors
421                 ),
422             },
423             code="missing_management_form",
424         )
425         self._non_form_errors.append(error)
426
427     for i, form in enumerate(self.forms):
428         # Empty forms are unchanged forms beyond those with initial data.
429         if not form.has_changed() and i >= self.initial_form_count():
430             empty_forms_count += 1
431         # Accessing errors calls full_clean() if necessary.
432         # _should_delete_form() requires cleaned_data.
433         form_errors = form.errors
434         if self.can_delete and self._should_delete_form(form):
435             continue
436         self._errors.append(form_errors)
437     try:
438         if (
439             self.validate_max
440             TOTAL_FORM_COUNT
441         ] > self.absolute_max:
442             raise ValidationError(
443                 self.error_messages["too_many_forms"] % {"num": self.max_num},
444                 code="too_many_forms",
445             )
446         if (
447             self.validate_min
448             and self.total_form_count()
449             - len(self.deleted_forms)
450             - empty_forms_count
451             < self.min_num
452         ):
453             raise ValidationError(
454                 self.error_messages["too_few_forms"] % {"num": self.min_num},
455                 code="too_few_forms",
456             )
457         # Give self.clean() a chance to do cross-form validation.
458         self.clean()
459     except ValidationError as e:
460         return "%s-%s" % (self.prefix, index)
461
462     507
463     508
464     509 def is_multipart(self):
465     510     """
466     511     Return True if the formset needs to be multipart, i.e. it
467     512     has FileInput, or False otherwise.
468     513     """
469     514     if self.forms:
470     515         return self.forms[0].is_multipart()
471     516     else:
472     517         return self.empty_form.is_multipart()
473     518
474     519 @property
475     520 def media(self):
476     521     # All the forms on a FormSet are the same, so you only need to
477     522     # interrogate the first form for media.
478     523     if self.forms:
479     524         return self.forms[0].media
480     525     else:
481     526         return self.empty_form.media

```

The part of the code in `BaseModelFormSet` where `empty_form` is created and `form_kwargs` are processed is located in the `django/forms/formsets.py` file. Specifically, the `empty_form` property is defined on lines 261-271, where the form is instantiated with various keyword arguments, including `form_kwargs`. The `form_kwargs` are processed by calling the `get_form_kwargs()` method on line 263, which prepares the keyword arguments for the form.

Here is the relevant code snippet:

```

```python
@property
def empty_form(self):
 form_kwargs = self.get_form_kwargs(None)
 form_kwargs.pop('empty_permitted', None)
 form = self.form(
 auto_id=self.auto_id,
 prefix=self.add_prefix("__prefix__"),
 empty_permitted=True,
 use_required_attribute=False,
 **form_kwargs,
 renderer=self.renderer,
)
 return form
...

```

This method creates the `empty_form` that is used by the `BaseModelFormSet` to provide an initial, empty form for the formset.

1067

Based on the information gathered, the issue arises because the `empty_form` property in the `BaseModelFormSet` class processes `form_kwargs` and includes the `empty_permitted` key, which is irrelevant for `empty_form`. This causes a `KeyError` when rendering the template.

To resolve this issue, we need to ensure that `empty_permitted` is ignored for `empty_form` by modifying the `empty_form` property to exclude `empty_permitted` from `form_kwargs`.

Next, I will use the `code_generator` to create a patch that modifies the `empty_form` property in the `BaseModelFormSet` class to exclude `empty_permitted` from `form_kwargs`.

---

Action:

```

```
{
  "action": "code_editor",
  "action_input": {
    "title": "Patch to fix KeyError in empty_form",
    "file_path": "django/forms/formsets.py",

```

1068


```

"context": "The issue is a KeyError when rendering a template that uses a formset's
↳ empty_form with form_kwargs containing empty_permitted. The expected behavior is that
↳ empty_permitted should be ignored for empty_form since it is irrelevant for
↳ empty_form. Ensure that form_kwargs.pop('empty_permitted', None) is called before
↳ passing form_kwargs to the form constructor.",
"request": "Modify the empty_form property in the BaseModelFormSet class to exclude
↳ empty_permitted from form_kwargs."
}
}
...

```

1069

REPOPILOT

Editor

Entering new execution chain..
Thought: I will look into "django/forms/formsets.py" to see relevant code snippets before editing

```

Action:
...
Action:
...
{
  "action": "open_file",
  "action_input": {
    "relative_file_path": "django/forms/formsets.py",
    "keyword": "class",
    "preview_size": 20
  }
}
...

```

1070

Observation

Open File

```

The content of django/forms/formsets.py is:
1 from django.core.exceptions import ValidationError
2 from django.forms import Form
3 from django.forms.fields import BooleanField, IntegerField
4 from django.forms.renderers import get_default_renderer
5 from django.forms.utils import ErrorList, RenderableFormMixin
6 from django.forms.widgets import CheckboxInput, HiddenInput, NumberInput
7 from django.utils.functional import cached_property
8 from django.utils.translation import gettext_lazy as _
9 from django.utils.translation import ngettext_lazy
10
11 __all__ = ("BaseFormSet", "formset_factory", "all_valid")
12
13 # special field names
14 TOTAL_FORM_COUNT = "TOTAL_FORMS"
15 INITIAL_FORM_COUNT = "INITIAL_FORMS"
16 MIN_NUM_FORM_COUNT = "MIN_NUM_FORMS"
17 MAX_NUM_FORM_COUNT = "MAX_NUM_FORMS"
18 ORDERING_FIELD_NAME = "ORDER"
19 DELETION_FIELD_NAME = "DELETE"
20
21 # default minimum number of forms in a formset
22 DEFAULT_MIN_NUM = 0
23
24 # default maximum number of forms in a formset, to prevent memory exhaustion
25 DEFAULT_MAX_NUM = 1000
26

```

1071

```

27
28 class ManagementForm(Form):
29     """
30     Keep track of how many form instances are displayed on the page. If adding
31     new forms via JavaScript, you should increment the count field of this form
32     as well.
33     """
34
35     template_name = "django/forms/div.html" # RemovedInDjango50Warning.
36
37     TOTAL_FORMS = IntegerField(widget=HiddenInput)
38     INITIAL_FORMS = IntegerField(widget=HiddenInput)
39     # MIN_NUM_FORM_COUNT and MAX_NUM_FORM_COUNT are output with the rest of the
40     # management form, but only for the convenience of client-side code. The
41     # POST value of them returned from the client is not checked.
42     MIN_NUM_FORMS = IntegerField(required=False, widget=HiddenInput)
43     MAX_NUM_FORMS = IntegerField(required=False, widget=HiddenInput)
44
45     def clean(self):
46         cleaned_data = super().clean()
47         # When the management form is invalid, we don't know how many forms
48         # were submitted.
49         cleaned_data.setdefault(TOTAL_FORM_COUNT, 0)
50         cleaned_data.setdefault(INITIAL_FORM_COUNT, 0)
51         return cleaned_data
52
53
54 class BaseFormSet(RenderableFormMixin):
55     """
56     A collection of instances of the same Form class.
57     """
58
59     deletion_widget = CheckboxInput
60     ordering_widget = NumberInput
61     default_error_messages = {
62         "missing_management_form": _(
63             "ManagementForm data is missing or has been tampered with. Missing fields: "
64             "%(field_names)s. You may need to file a bug report if the issue persists."
65         ),
66         "too_many_forms": ngettext_lazy(
67             "Please submit at most %(num)d form.",
68             "Please submit at most %(num)d forms.",
69             "num",
70         ),
71         "too_few_forms": ngettext_lazy(
72             "Please submit at least %(num)d form.",
73             "Please submit at least %(num)d forms.",
74             "num",
75         ),
76     }
77
78     template_name_div = "django/forms/formsets/div.html"
79     template_name_p = "django/forms/formsets/p.html"
80     template_name_table = "django/forms/formsets/table.html"8 from django.utils.translation
↪ import gettext_lazy as _
9 from django.utils.translation import ngettext_lazy
10
11 __all__ = ("BaseFormSet", "formset_factory", "all_valid")
12
13 # special field names
14 TOTAL_FORM_COUNT = "TOTAL_FORMS"
15 INITIAL_FORM_COUNT = "INITIAL_FORMS"
16 MIN_NUM_FORM_COUNT = "MIN_NUM_FORMS"
17 MAX_NUM_FORM_COUNT = "MAX_NUM_FORMS"
18 ORDERING_FIELD_NAME = "ORDER"
19 DELETION_FIELD_NAME = "DELETE"
20
21 # default minimum number of forms in a formset

```

```

22 DEFAULT_MIN_NUM = 0
23
24 # default maximum number of forms in a formset, to prevent memory exhaustion
25 DEFAULT_MAX_NUM = 1000
26
27
28 class ManagementForm(Form):
29     """
30     Keep track of how many form instances are displayed on the page. If adding
31     new forms via JavaScript, you should increment the count field of this form
32     as well.
33     """
34
35     template_name = "django/forms/div.html" # RemovedInDjango50Warning.
36
37     TOTAL_FORMS = IntegerField(widget=HiddenInput)
38     INITIAL_FORMS = IntegerField(widget=HiddenInput)
39     # MIN_NUM_FORM_COUNT and MAX_NUM_FORM_COUNT are output with the rest of the
40     # management form, but only for the convenience of client-side code. The
41     # POST value of them returned from the client is not checked.
42     MIN_NUM_FORMS = IntegerField(required=False, widget=HiddenInput)
43     MAX_NUM_FORMS = IntegerField(required=False, widget=HiddenInput)
44
45     def clean(self):
46         cleaned_data = super().clean()
47         # When the management form is invalid, we don't know how many forms
34
35     template_name = "django/forms/div.html" # RemovedInDjango50Warning.
36
37     TOTAL_FORMS = IntegerField(widget=HiddenInput)
38     INITIAL_FORMS = IntegerField(widget=HiddenInput)
39     # MIN_NUM_FORM_COUNT and MAX_NUM_FORM_COUNT are output with the rest of the
40     # management form, but only for the convenience of client-side code. The
41     # POST value of them returned from the client is not checked.
42     MIN_NUM_FORMS = IntegerField(required=False, widget=HiddenInput)
43     MAX_NUM_FORMS = IntegerField(required=False, widget=HiddenInput)
44
45     def clean(self):
46         cleaned_data = super().clean()
47         # When the management form is invalid, we don't know how many forms
48         # were submitted.
49         cleaned_data.setdefault(TOTAL_FORM_COUNT, 0)
50         cleaned_data.setdefault(INITIAL_FORM_COUNT, 0)
51         return cleaned_data
52
53
54 class BaseFormSet(RenderableFormMixin):
55     """
56     A collection of instances of the same Form class.
57     """
58
59     deletion_widget = CheckboxInput
60     ordering_widget = NumberInput
61     default_error_messages = {
62         "missing_management_form": _(
63             "ManagementForm data is missing or has been tampered with. Missing fields: "
64             "%(field_names)s. You may need to file a bug report if the issue persists."
65         ),
66         "too_many_forms": ngettext_lazy(
67             "Please submit at most %(num)d form.",
68             "Please submit at most %(num)d forms.",
69             "num",
70         ),
71         "too_few_forms": ngettext_lazy(
72             "Please submit at least %(num)d form.",
73             "Please submit at least %(num)d forms.",
36
37     TOTAL_FORMS = IntegerField(widget=HiddenInput)

```

```

38     INITIAL_FORMS = IntegerField(widget=HiddenInput)
39     # MIN_NUM_FORM_COUNT and MAX_NUM_FORM_COUNT are output with the rest of the
40     # management form, but only for the convenience of client-side code. The
41     # POST value of them returned from the client is not checked.
42     MIN_NUM_FORMS = IntegerField(required=False, widget=HiddenInput)
43     MAX_NUM_FORMS = IntegerField(required=False, widget=HiddenInput)
44
45     def clean(self):
46         cleaned_data = super().clean()
47         # When the management form is invalid, we don't know how many forms
48         # were submitted.
49         cleaned_data.setdefault(TOTAL_FORM_COUNT, 0)
50         cleaned_data.setdefault(INITIAL_FORM_COUNT, 0)
51         return cleaned_data
52
53
54 class BaseFormSet(RenderableFormMixin):
55     """
56     A collection of instances of the same Form class.
57     """
58
59     deletion_widget = CheckboxInput
60     ordering_widget = NumberInput
61     default_error_messages = {
62         "missing_management_form": _(
63             "ManagementForm data is missing or has been tampered with. Missing fields: "
64             "%(field_names)s. You may need to file a bug report if the issue persists."
65         ),
66         "too_many_forms": ngettext_lazy(
67             "Please submit at most %(num)d form.",
68             "Please submit at most %(num)d forms.",
69             "num",
70         ),
71         "too_few_forms": ngettext_lazy(
72             "Please submit at least %(num)d form.",
73             "Please submit at least %(num)d forms.",
74             "num",
75         ),
76     },
77
78     template_name_div = "django/forms/formsets/div.html"
79     template_name_p = "django/forms/formsets/p.html"
80     template_name_table = "django/forms/formsets/table.html"
81     template_name_ul = "django/forms/formsets/ul.html"
82
83     def __init__(
84         self,
85         data=None,
86         files=None,
87         auto_id="id_%s",
88         prefix=None,
89         initial=None,
90         error_class=ErrorList,
91         form_kwargs=None,
92         error_messages=None,
93     ):
94         self.is_bound = data is not None or files is not None
95         self.prefix = prefix or self.get_default_prefix()
96         self.auto_id = auto_id
97         self.data = data or {}
98         self.files = files or {}
99         self.initial = initial

```

```

100     self.form_kwargs = form_kwargs or {}
101     self.error_class = error_class
102     self._errors = None
103     self._non_form_errors = None
104
105     messages = {}
106     for cls in reversed(type(self).__mro__):
107         messages.update(getattr(cls, "default_error_messages", {}))
108     if error_messages is not None:
109         messages.update(error_messages)
81     template_name_ul = "django/forms/formsets/ul.html"
82
83     def __init__(
84         self,
85         data=None,
86         files=None,
87         auto_id="id_%s",
88         prefix=None,
89         initial=None,
90         error_class=ErrorList,
91         form_kwargs=None,
92         error_messages=None,
93     ):
94         self.is_bound = data is not None or files is not None
95         self.prefix = prefix or self.get_default_prefix()
96         self.auto_id = auto_id
97         self.data = data or {}
98         self.files = files or {}
99         self.initial = initial
100         self.form_kwargs = form_kwargs or {}
101         self.error_class = error_class
102         self._errors = None
103         self._non_form_errors = None
104
105         messages = {}
106         for cls in reversed(type(self).__mro__):
107             messages.update(getattr(cls, "default_error_messages", {}))
108         if error_messages is not None:
109             messages.update(error_messages)
110         self.error_messages = messages
111
112     def __iter__(self):
113         """Yield the forms in the order they should be rendered."""
114         return iter(self.forms)
115
116     def __getitem__(self, index):
117         """Return the form at the given index, based on the rendering order."""
118         return self.forms[index]
119
120     def __len__(self):

```

1075

REPOPILOT

Editor

Entering new execution chain..

Thought: The `empty_form` property is defined in the `BaseFormSet` class within the `django/forms/formsets.py` file.

Here is the relevant code snippet from the `BaseFormSet` class:

```

```python
@property
def empty_form(self):
 form_kwargs = self.get_form_kwargs(None)
 form_kwargs.pop('empty_permitted', None)
 form = self.form(
 auto_id=self.auto_id,
 prefix=self.add_prefix("__prefix__"),

```

1076

```

 empty_permitted=True,
 use_required_attribute=False,
 **form_kwargs,
 renderer=self.renderer,
)
 self.add_fields(form, None)
 return form
...

```

I will now apply the required modification to the `empty\_form` property in the `BaseFormSet` class.

Action:

```

...
{
 "action": "autorepair_editor_file",
 "action_input": {
 "relative_file_path": "django/forms/formsets.py",
 "start_line": 261,
 "end_line": 274,
 "patch": " @property\n def empty_form(self):\n form_kwargs =
↪ self.get_form_kwargs(None)\n form_kwargs.pop('empty_permitted', None)\n
↪ form = self.form(\n auto_id=self.auto_id,\n
↪ prefix=self.add_prefix(\"__prefix__\"),\n empty_permitted=True,\n
↪ use_required_attribute=False,\n **form_kwargs,\n
↪ renderer=self.renderer,\n)\n self.add_fields(form, None)\n return
↪ form"
 }
}
...

```

1077

### Observation

### Auto-repair Editor

Successfully edited the file django/forms/formsets.py from line 261 to 274

1078

### Patch

```

diff --git a/django/forms/formsets.py b/django/forms/formsets.py
index 57676428ff..221aa405a8 100644
--- a/django/forms/formsets.py
+++ b/django/forms/formsets.py
@@ -214,8 +214,10 @@ class BaseFormSet(RenderableFormMixin):
 index will be None if the form being constructed is a new empty
 form.
 """
- return self.form_kwargs.copy()
-
+ form_kwargs = self.form_kwargs.copy()
+ if index is None:
+ form_kwargs.pop('empty_permitted', None)
+ return form_kwargs

def constructform(self, i, kwargs):
 """Instantiate and return the i-th form instance in a formset."""
@@ -255,14 +257,17 @@ class BaseFormSet(RenderableFormMixin):
 """Return a list of all the extra forms in this formset."""
 return self.forms[self.initial_form_count() :]

+ @property
+ @property
+ def empty_form(self):
+ form_kwargs = self.get_form_kwargs(None)
+ form_kwargs.pop('empty_permitted', None)

```

1079

```
form = self.form(
 auto_id=self.auto_id,
 prefix=self.add_prefix("prefix"),
 empty_permitted=True,
 use_required_attribute=False,
- self.get_form_kwargs(None),
+ **form_kwargs,
 renderer=self.renderer,
)
self.add_fields(form, None)
```

1080

## REPOPILOT

## Planner

The modification to the `empty_form` property in the `BaseFormSet` class has been successfully applied. The `empty_permitted` key is now excluded from `form_kwargs` before passing it to the form constructor, addressing the issue of the `KeyError` when rendering a template that uses a formset's `empty_form`.

---

Action:

```
...
{
 "action": "Final Answer",
 "action_input": {
 "context": "The task is fully resolved.",
 }
}
...
```

1081