

# RepoGraph: Enhancing AI Software Engineering with Repository-level Code Graph

Siru Ouyang<sup>1\*</sup>, Wenhao Yu<sup>2</sup>, Kaixin Ma<sup>2</sup>, Zilin Xiao<sup>3\*</sup>, Zhihan Zhang<sup>4\*</sup>, Mengzhao Jia<sup>4\*</sup>, Jiawei Han<sup>1</sup>, Hongming Zhang<sup>2</sup>, Dong Yu<sup>2</sup>

<sup>1</sup> University of Illinois Urbana-Champaign <sup>2</sup> Tencent AI Seattle Lab

<sup>3</sup> Rice University <sup>4</sup> University of Notre Dame

siruo2@illinois.edu

## Abstract

Large Language Models (LLMs) excel in code generation yet struggle with modern AI software engineering tasks. Unlike traditional function-level or file-level coding tasks, AI software engineering requires not only basic coding proficiency but also advanced skills in managing and interacting with code repositories. However, existing methods often overlook the need for repository-level code understanding, which is crucial for accurately grasping the broader context and developing effective solutions. On this basis, we present REPOGRAPH, a plug-in module that manages a repository-level structure for modern AI software engineering solutions. REPOGRAPH offers the desired guidance and serves as a repository-wide navigation for AI software engineers. We evaluate REPOGRAPH on the SWE-bench by plugging it into four different methods of two lines of approaches, where REPOGRAPH substantially boosts the performance of all systems, leading to *a new state-of-the-art* among open-source frameworks. Our analyses also demonstrate the extensibility and flexibility of REPOGRAPH by testing on another repo-level coding benchmark, CrossCodeEval. Our code is available at <https://github.com/ozyyshr/RepoGraph>.

## 1 Introduction

Recent advancements in large language models (LLMs) have showcased their powerful capabilities across various natural language processing tasks (OpenAI, 2023; Anil et al., 2023; Dubey et al., 2024), and now, coding-specific LLMs are emerging to tackle complex software engineering challenges (Hou et al., 2023; Fan et al., 2023), such as Code-Llama (Rozière et al., 2023) and StarCoder (Li et al., 2023a). These coding-specific LLMs are capable of assisting users with various

software engineering tasks, even achieving human-level performance in many function-level coding tasks, such as program synthesis (Chen et al., 2021; Austin et al., 2021), code annotation (Yao et al., 2019), bug fixing (Tufano et al., 2019), and code translation (Rozière et al., 2020).

Real-world software engineering often extends beyond single function or self-contained code files. Applications are typically built as repositories containing multiple interdependent files, modules, and libraries (Bairi et al., 2024). These complex structures require a holistic understanding of the entire codebase to perform tasks such as code completion (Shrivastava et al., 2023; Ding et al., 2023), feature addition (Liang et al., 2024), or issue resolving (Jimenez et al., 2024). Recent benchmarks like SWE-Bench (Jimenez et al., 2024) have been proposed to evaluate LLMs on real-world GitHub issues. It requires LLMs to modify the repository to resolve the issue, either by fixing a bug or introducing a new feature. This task is particularly challenging because it requires navigating complex code bases, understanding intricate dependencies between code files, and ensuring that changes integrate seamlessly without introducing new issues, which highlights the difficulties in scaling from function-level to repository-level understanding, as expounded in Figure 1.

A key step in addressing repository-level tasks is to understand the structure of a repository and identify related code. To achieve this, retrieval-augmented generation (RAG) and its variants (Zhang et al., 2023; Phan et al., 2024; Wu et al., 2024) have been leveraged, in a procedural manner, to retrieve relevant code files across the repository first, providing context for LLMs for further edition. However, indexing at file-level can only identify semantically similar but not genuinely related code snippets. Instead of using RAG, recent approaches like Agentless (Xia et al., 2024) construct a skeletal format for each file, and di-

\* This work was done when Siru, Zilin, Zhihan and Mengzhao were interns at Tencent AI Lab, Seattle.

### (a) Function-level Coding Problem

**Input text:** Write a python function to find the first repeated character in a given string.



```
1. def first_repeated_char(str1):
2.   for index,c in enumerate(str1):
3.     if str1[:index+1].count(c) > 1: return c
4.   return "None"
```

### (b) Repository-level Coding Problem

**Issue:** modeling's "`separability_matrix`" does not compute separability correctly for `CompoundModels...`

#### Code repository:

astropy core.py fitting.py  
coordinates earth.py ...

(i) Understand intricate dependencies



#### Generated patch:

```
diff --git a/astropy/modeling/
separable.py
--- a/astropy/modeling/separable.py
+++ b/astropy/modeling/separable.py
@@ -242,7 +242,7 @@ def _cstack(left,
right):...
```



#### Unit test:

test_coord_matrix	✗
test_cdot	✓
test_arith_oper	✗

(iii) Resolve without introducing new issues

(ii) Navigate complex codebases

Figure 1: The illustration of (a) a function-level coding problem from HumanEval (Chen et al., 2021) and (b) a repository-level coding problem from SWE-Bench (Jimenez et al., 2024).

rectly prompt LLMs to identify relevant files and code lines. However, this method still treats code repositories as flat documents (Zhang et al., 2024), which suffers from limitations of repository structure such as the intricate inter-dependencies cross files. An alternative approach is to design agent frameworks (Yang et al., 2024; Wang et al., 2024b), which enables LLMs to interact with repositories using actions. While LLM agents can freely determine the next action based on current observations, without the grasp of global repository structures, they tend to focus narrowly on specific files, resulting in local optimums. Addressing these limitations requires going beyond semantic matching and developing techniques that enable a deeper understanding of the codebase structure. This will allow LLMs to leverage fine-grained context across multiple files and function calls, facilitating more informed, repository-wide decision-making for coding tasks.

Motivated by this, we propose REPOGRAPH, a *plug-in* module designed to help LLMs-based AI programmers leverage the code structure of *an entire repository*. REPOGRAPH is a graph structure and operates at the line level, offering a more fine-grained approach compared to previous file-level browsing methods. Each node in the graph represents a line of code, and edges represent the dependencies of code definitions and references. REPOGRAPH is constructed via code line parsing and encodes the structured representation of the entire repository. Sub-graph retrieval algorithms are then used to extract ego-graphs centered around specific keywords. These ego-graphs can be smoothly integrated with both procedural and agent frameworks, offering key clues that provide a more comprehensive context for LLMs to solve real-world software engineering problems.

To assess REPOGRAPH’s effectiveness and versatility as a plug-in module, we integrate it with four existing software engineering frameworks and evaluate its performance using SWE-bench, a recent benchmark for AI software engineering. Experiment results show that REPOGRAPH boosts the success rate of existing methods for both agent and procedural frameworks by achieving an average relative improvement of 32.8%. We also test REPOGRAPH on CrossCodeEval to verify its transferability to general coding tasks that require repository-level code understanding. Additionally, we systematically analyze different sub-graph retrieval algorithms and integration methods. Together with error analyses, we hope to shed light on future works targeting modern AI software engineering.

## 2 Related Works

### 2.1 LLM-based methods for AI software engineering

Recently, there has been a significant increase in research focused on AI-driven software engineering, which can be broadly categorized into two primary approaches: (i) LLM agent-based frameworks and (ii) SWE-featured procedural frameworks. While this field has advanced rapidly, with most methods being released as proprietary solutions for industry applications (Cognition, 2024), our related work section will concentrate specifically on open-source frameworks.

**LLM agent-based framework** equips large language models (LLMs) with a set of predefined tools, allowing agents to iteratively and autonomously perform actions, observe feedback, and plan future steps (Yang et al., 2024; Zhang et al., 2024; Wang et al., 2024b; Cognition, 2024). While the exact set of tools may vary across dif-

ferent agent frameworks, they typically include capabilities such as opening, writing, or creating files, searching for code lines, running tests, and executing shell commands. To solve a problem, agent-based approaches involve multiple actions, with each subsequent turn depending on the actions taken in previous ones and the feedback received from the environment. For example, SWE-agent (Yang et al., 2024) facilitates interactions with the execution environment by designing a special agent-computer interface (ACI). There are various actions, including “search and navigation”, “file viewer and editor”, and “context management”. Another work, AutoCodeRover (Zhang et al., 2024), further offers fine-grained searching methods for LLM agents in better contexts without an execution process. Specifically, it supports class and function-level code search. OpenDevin (Wang et al., 2024b), initiated after Devin (Cognition, 2024), is a community-driven platform that integrates widely used agent systems. The action space design of OpenDevin is highly flexible, requiring LLM agents to generate code on the fly.

**SWE-featured procedural frameworks** typically follow a two-step *Localize/Search-Edit* approach, as seen in existing literature (Zhang et al., 2023; Wu et al., 2024; Liang et al., 2024; Xia et al., 2024). The *localize* step focuses on identifying relevant code snippets, while the *edit* step involves completing or revising the code. Some works introduce additional steps to further enhance performance, such as the *Search-Expand-Edit* approach (Phan et al., 2024). Retrieval (Lewis et al., 2020) is a popular technique used for localization, allowing models to search for relevant code snippets from large repositories by treating issue descriptions as queries and code snippets as indexed data. Some approaches use a sliding window to ensure completeness (Zhang et al., 2023). Besides, Agentless (Xia et al., 2024) is a recently developed method that uses LLMs to directly identify relevant elements for editing within code repositories. It first recursively traverses the repository structure to generate a format that aligns files and folders vertically, with indents for sub-directories. This structure and the issue description are then input into the LLM, which performs a hierarchical search to identify the top-ranked suspicious files requiring further inspection or modification.

Table 1: Comparison between our approach REPOGRAPH and existing methods for representing the repository on various aspects. **\*RepoUnderstander (Ma et al., 2024) and CodexGraph (Liu et al., 2024) are concurrent works to ours.**

Model	Line-level	File-level	Repo-level
DraCo	✗	✓	✗
Aider	✓	✗	✗
RepoUnderstander*	✗	✓	✓
CodexGraph*	✗	✓	✓
REPOGRAPH	✓	✓	✓

## 2.2 Repository-level Coding Capability

The evaluation of coding capabilities in AI systems has traditionally focused on function-level or line-level assessments (Lu et al., 2021; Chen et al., 2021; Austin et al., 2021), where individual code snippets or isolated functions are the primary units of analysis. Unlike previous studies, SWE-bench (Jimenez et al., 2024) highlights the trend of repository-level coding, driven by recent advances of coding-specific LLMs (Guo et al., 2024; Li et al., 2023b). It reflects the growing user demand to understand and contribute to entire projects rather than isolated functions (Ouyang et al., 2023), as well as solving real-world problems in an end-to-end and automatic manner.

While the pre-trained code LLMs mentioned earlier incorporate repository-level information such as file dependencies, tasks at the repository level often involve more intricate call relationships within their context. Recent works like RepoCoder (Zhang et al., 2023) and RepoFuse (Liang et al., 2024) have started integrating Retrieval-Augmented Generation (RAG) modules to harness additional information from repositories. Building on this, subsequent research has focused on embedding repository-level context into their methodologies. For instance, DraCo (Cheng et al., 2024) introduces importing relationships between files, while Aider (Gauthier, 2024) employs PageRank (Page, 1999) to identify the most significant contextual elements. RepoUnderstander (Ma et al., 2024) and CodexGraph (Liu et al., 2024) model code files as a knowledge graph. Despite similarities in representation, methods vary in how they retrieve information from these structures and utilize it for downstream tasks. Table 1 summarizes the differences between these methods and REPOGRAPH. REPOGRAPH surpasses previous approaches by effectively integrating context at the line, file, and

repository levels.

### 3 RepoGraph

This section introduces REPOGRAPH, a novel plug-in module that can be seamlessly integrated into existing research workflows for both agent-based and procedural frameworks. The primary goal of REPOGRAPH is to provide a structured way to analyze and interact with complex codebases, enabling detailed tracing of code dependencies, execution flow, and structural relationships across the repository. In the following sections, we will provide a detailed description of REPOGRAPH’s construction, its underlying representation, and its utility across various scenarios. The overall architecture is depicted in Figure 2, highlighting its key components and operational flow.

#### 3.1 Construction

Given a repository-level coding task, the first step is to carefully examine the repository structure so that the necessary information can be collected. The input for REPOGRAPH construction is a repository, i.e., a collection of its folders and files, while the output is a structured graph, where each node is a code line, and each edge represents the dependencies in between. REPOGRAPH enables tracing back to the root cause of the current issue and gathering dependent code context to help solve the problem. The construction process of REPOGRAPH could be divided into three key steps.

**Step 1: Code line parsing.** We first traverse the entire repository using a top-down approach to identify all code files as candidates for next-step parsing. This is accomplished by filtering based on file extensions, retaining only those with relevant code file suffixes (e.g., `.py`) while excluding non-essential file types (e.g., `.git` or `requirements.txt`), which are noisy and irrelevant for coding tasks. For each code file, we utilize tree-sitter<sup>1</sup> to parse the code, leveraging its Abstract Syntax Tree (AST) framework. The AST provides a tree-based representation of the abstract syntactic structure of the source code, enabling the identification of key elements such as functions, classes, variables, types, and other definitions. While recognizing these definitions is crucial, tracing their usage and references throughout the code is equally important. Tree-

<sup>1</sup><https://pypi.org/project/tree-sitter-languages/>

sitter facilitates this by capturing the definitions and tracking where they are utilized or referenced within the codebase. For example, in figure 2, we not only identify definitions like `class Model` and its inherent methods but also references like `self._validate_input_units()`. After processing each line of code with a tree-sitter, we selectively retain lines that involve function calls and dependency relations, discarding extraneous information. Our focus is primarily on the *functions* and *classes*, as these represent the core structural components of the code. By concentrating on these elements and their interrelationships, REPOGRAPH optimizes the analysis process by excluding less significant details, such as individual variables, which tend to be redundant and less relevant for further processing.

**Step 2: Project-dependent relation filtering.** After the previous parsing step, we obtain code lines with calling and dependency relations. However, not all relations are useful for fixing issues. Specifically, many default and built-in function/class calls could distract from the project-related ones. Therefore, we additionally introduce a filtering process that excludes the repository-independent relations. Two types of such relations exist: (i) *global relation* refers to Python standard and built-in functions and classes. (ii) *local relation* are introduced by third-party libraries, which are specific to the current code file. For global relations, we maintain a comprehensive list of methods from standard and built-in libraries, excluding any identified relations from this list. The list is empirically constructed by gathering methods of the `builtins` library and default methods such as `list` and `tuple`. For example, in figure 2, `line inputs = len(input)` is excluded since `len` is a default method. For local relations, we parse import statements in the code to identify third-party methods that are included, and exclude them accordingly.

**Step 3: Graph organization.** At this stage, we construct REPOGRAPH using code lines as the fundamental building units. The graph can be represented as  $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$ , where  $\mathcal{V}$  represents the set of nodes, with each node corresponding to a line of code, and  $\mathcal{E}$  represents the set of edges, capturing the relationships between these code lines. Each node in  $\mathcal{V}$  contains attributes to represent its meta-information, such as `line_number`, `file_name`, `directory`, etc. Additionally, we classify each code line as either a “definition” (def) or a “reference” (ref) to a particular module. A “def” node

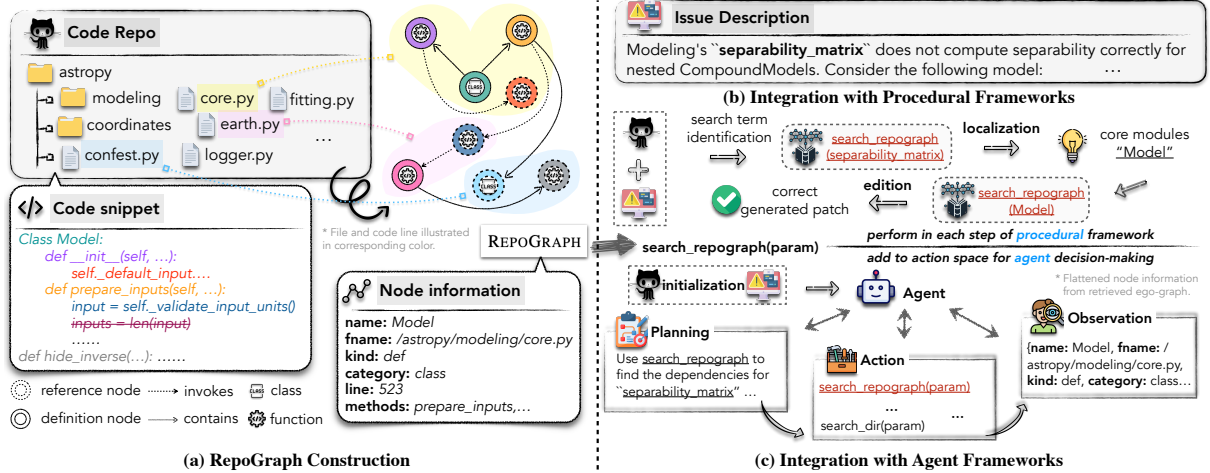


Figure 2: An in-depth illustration of (a) the construction, (b) the integration with procedural frameworks, and (c) the integration with agent frameworks of REPOGRAPH. Given a code repository, we first utilize AST to construct  $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$ , where  $\mathcal{G}$  consists of “reference” and “definition” node,  $\mathcal{E}$  includes “invoke” and “contain” relations. The constructed REPOGRAPH are then used in procedural frameworks by adding sub-retrieval results into each step, and agent frameworks by adding graph retrieval as an additional action “search\_reposgraph”.

corresponds to the line where a function, class, or variable is initially defined, while a “ref” node indicates a code line where this entity is referenced or invoked elsewhere in the code. Similar to soft links to “def” nodes, “ref” nodes also represent other variations of invoking methods. For example, in Figure 2, the class definition “class Model” would be a “def” node, while any subsequent usages of Model would be “ref” nodes. Each “def” node may have multiple “ref” nodes associated with it, as a single function or class can be referenced in various places throughout the code. We define two types of edges:  $\mathcal{E}_{invoke}$  and  $\mathcal{E}_{contain}$ . The triple  $(\mathcal{V}_1, \mathcal{E}_{contain}, \mathcal{V}_2)$  denotes that  $\mathcal{V}_1$  (e.g., a function definition) contains another module  $\mathcal{V}_2$  (e.g., an internal function or class). The edge  $\mathcal{E}_{contain}$  typically connects a “def” node to its internal components. In contrast,  $\mathcal{E}_{invoke}$  represents an invocation relationship, usually connecting a “def” node to a “ref” node, where the reference node includes a dependency on the definition.

### 3.2 Utility

The constructed REPOGRAPH serves as a structured representation of the current repository and facilitates better-related information collection and aggregation. For information collection based on REPOGRAPH, specifically, we use one search term each time for subgraph retrieval. Search terms are the key functions or classes that are determined by current states. For example, “separability\_matrix” is the initial search term in Figure 2(c). We retrieve

the  $k$ -hop ego-graphs (Hu et al., 2024) with the search term in the centric. The ego-graph is crucial for solving the problem because it focuses on the immediate relationships (Jin et al., 2024) around the search term, capturing the relevant dependencies and interactions within the repository, which is key to understanding the functional context. Additionally, the retrieved content explicitly contains information at both the method and line levels and implicitly expresses the grouping at the file level.

This process is abstracted via  $search\_reposgraph()$  as illustrated in the middle of Figure 2. The retrieved  $k$ -hop ego-graph will be flattened for further processing. We also tried other variants for integration later in Section 5.2 and their performance in Table 4. We narrate how REPOGRAPH could be plugged in with existing representative research lines in the following.

**Integration with procedural framework.** In a procedural framework, LLMs are usually prompted in “localization” and “edition” stages with the given repository context and issue description. In this case, we use  $search\_reposgraph()$  before both stages, leveraging our REPOGRAPH to assist in making more informed decisions at each step. For example, in Figure 2, we first include the subgraph of “separability\_matrix” for localization, and then use the localized result “Model” to search in the edition stage. To implement the strategy, we flatten the context of retrieved ego-graphs and append it as part of the prompt. As a result, the LLM generation is conditioned on both retrieved ego-graphs and the

context provided by baseline methods, helping the model preserve nuances.

**Integration with agent framework.** A significant difference in existing agent frameworks is the action space design, as expounded in Section 2. To leverage the power of REPOGRAPH, we put `search_repograph()` as an additional action in the action space. The agent decides when and where to use this action. The search term is also determined by the agent accordingly. The returned subgraph will be flattened and used as an observation for the next state.

## 4 Experiments

### 4.1 Setup

We evaluated REPOGRAPH as a plug-in component, i.e., integrated into existing baseline models of the two aforementioned research lines to assess its performance. We use the same baseline settings and configurations when incorporating REPOGRAPH to ensure a fair comparison.

**Dataset.** We test REPOGRAPH in SWE-bench-Lite<sup>2</sup>. Each problem in the dataset requires submitting a patch to solve the underlying issue described in the input issue description. The goal is to generate a patch that accurately revises the relevant portions of the codebase within the repository, ensuring that all test scripts included in the dataset are successfully executed.

**Baselines.** We integrate REPOGRAPH with representative methods from both aforementioned research lines. (i) For procedural frameworks, we evaluate the widely used traditional method, RAG (Lewis et al., 2020), as well as Agentless (Xia et al., 2024), an open-source state-of-the-art approach in this direction. For RAG, we follow its initial setting and use BM25 for file-level retrieval. After that, we append the context of REPOGRAPH after the code files as part of the prompt. Agentless first performs a hierarchical localization in terms of “file-class/function-edits” and then conducts repair based on localization. The context of REPOGRAPH is inserted in every step of Agentless. (ii) For agent frameworks, we consider SWE-agent (Yang et al., 2024) and AutoCodeRover (Zhang et al., 2024). For both frameworks, we add an additional action “search\_repograph” for the LLM agent as described in Section 3. All the choices in the two research

lines incorporate GPT-4 and GPT-4o based methods to ensure the generalization. Detailed implementations and prompts used can be found in Appendix A.

**Evaluation metrics.** We evaluate all methods across two key dimensions: Accuracy and Average Cost. (i) For accuracy, we report the *resolve rate* and *patch application rate*. The resolve rate represents the percentage of issues successfully resolved across all data points. An issue is considered resolved if the submitted patch passes all test scripts. To assess the patch application rate, we attempt to apply the generated patches to the repository using the patch program; successful applications contribute to this metric. (ii) To evaluate cost efficiency, we report two metrics: *average cost* and *average tokens*, which refer to the inference cost and the number of input/output tokens used when querying the LLMs, respectively.

**Configurations.** We use the same GPT version as in the baselines in the experiments. We used GPT-4o (2024-05-13) and GPT-4-Turbo (gpt-4-1106-preview) from OpenAI for evaluation and analyses in our experiments. All evaluation processes are performed in a containerized Docker environment<sup>3</sup>, ensuring stability and reproducibility, made possible through contributions from the open-source community. For plug-in with procedural frameworks, it usually takes around 2-3 hours to finish. For agent frameworks, the inference time is larger, up to around 10 hours.

### 4.2 Experiment results

Table 2 presents the main evaluation results of all baseline methods and the corresponding performance with REPOGRAPH (+REPOGRAPH) as a plug-in in the SWE-bench-Lite test set. We also report the number of correct samples for each method. The performance increase is marked by  $\uparrow_{num}$ . Based on the results, we have the following key observations:

**(i) REPOGRAPH brings consistent performance gain for all combinations of frameworks and LLM model bases.** Specifically, REPOGRAPH achieves an absolute improvement of +2.66 and +2.34 in terms of the resolve rate for RAG and Agentless, respectively, which is 99.63% and 8.56% of relative improvement. The notable improvement demonstrates the effectiveness of our REPOGRAPH in adapting to various scenarios by

<sup>2</sup>Current leaderboard could be found at <https://www.swebench.com/>

<sup>3</sup><https://github.com/aorwall/SWE-bench-docker>

<sup>4</sup><https://github.com/swe-bench/experiments/tree/main/evaluation/lite>

Table 2: Results of REPOGRAPH with open-source baselines in two research lines, including procedural and agent frameworks. Numbers of accuracy-related metrics are directly taken from the leaderboard, while the cost-related ones are computed from the corresponding trajectories<sup>4</sup>.

Methods	LLM	Accuracy			Avg. Cost	
		<i>resolve</i>	# samples	<i>patch apply</i>	\$ cost	# tokens
<i>Procedural frameworks</i>						
RAG	GPT-4	2.67	8	29.33	\$0.13	11,736
+REPOGRAPH	GPT-4	5.33 $\uparrow$ 2.66	16 $\uparrow$ 8	47.67 $\uparrow$ 18.34	\$0.17	15,439
Agentless	GPT-4o	27.33	82	97.33	\$0.34	42,376
+REPOGRAPH	GPT-4o	29.67 $\uparrow$ 2.34	89 $\uparrow$ 7	98.00 $\uparrow$ 0.67	\$0.39	47,323
<i>Agent frameworks</i>						
AutoCodeRover	GPT-4	19.00	57	83.00	\$0.45	38,663
+REPOGRAPH	GPT-4	21.33 $\uparrow$ 2.33	64 $\uparrow$ 7	86.67 $\uparrow$ 3.67	\$0.58	45,112
SWE-agent	GPT-4o	18.33	55	87.00	\$2.51	245,008
+REPOGRAPH	GPT-4o	20.33 $\uparrow$ 2.00	61 $\uparrow$ 6	90.33 $\uparrow$ 3.33	\$2.69	262,512

inducing relevant code context and performing precise code editions. Additionally, our best performance so far by plugging in Agentless, 29.67, achieves the *state-of-the-art* performance on the benchmark<sup>5</sup> among all open-source methods.

**(ii) Performance gain brought by REPOGRAPH is slightly larger on procedural frameworks than agent ones.** With procedural frameworks, REPOGRAPH correctly fixes more issues than agent ones. This could be due to two primary reasons. Firstly, we observed that mature procedural frameworks tend to achieve better baseline performance than agent-based frameworks on SWE-bench. The initial definitive nature of procedural frameworks, with their well-defined running flow and structure, allows them to leverage plugins more effectively. Another reason is that this deterministic behavior reduces the complexity that arises from dynamic decision-making, a key characteristic of agent-based systems, thereby enabling a smoother integration of performance improvements.

**(iii) Performance gain brought by REPOGRAPH does not rely on more costs.** We also compute and report each method’s average cost and token consumption. By introducing REPOGRAPH, we manage to reduce the costs associated with managing the entire repository while achieving comparable or even superior performance. As shown in Table 2, the performance improvements achieved by REPOGRAPH are generally proportional to, or slightly lower than, the corresponding increase in cost. This indicates REPOGRAPH’s performance gains are not mainly due to increased

token usage.

**(iv) Average costs are generally larger on agent frameworks with REPOGRAPH.** This phenomenon is especially obvious with SWE-agent, as it allows the agent to freely determine the next action based on the current observation. We also found that the integration with agent frameworks usually leads to larger cost increases, as exemplified by +0.13\$ and +0.18\$ with AutoCodeRover and SWE-agent, respectively. The reason lies in the large exploration space in agent frameworks. The agents might call the *search\_repograph()* action many times, which leads to the explosion of prompt contexts. We encourage users to be mindful of cost and to adopt a more granular approach to cost control when integrating REPOGRAPH into the agent framework in the future.

## 5 Analysis

This section presents a detailed analysis to demonstrate that the additional context provided by REPOGRAPH is beneficial for the task. We begin by analyzing localization accuracy in comparison to the gold-standard patch. Next, we explore various REPOGRAPH configurations, focusing on how the additional context can be effectively integrated into the existing system. Finally, we perform an in-depth error analysis, highlighting aspects where REPOGRAPH can be further improved. For more analyses including resolve rate in various aspects and action distributions of agent frameworks, please refer to Appendices C.

<sup>5</sup><https://www.swebench.com/>.

Table 3: Percentage of problems for accurate edition localizations with respect to file, function, and line levels. All the numbers are computed from the corresponding generated patches.

Methods	LLM	<i>file</i>	<i>function</i>	<i>line</i>
RAG	GPT-4	47.3	23.3	12.7
+REPOGRAPH	GPT-4	51.7 $\uparrow$ 4.4	25.3 $\uparrow$ 2.0	14.3 $\uparrow$ 1.6
Agentless	GPT-4o	68.7	51.0	34.3
+REPOGRAPH	GPT-4o	74.3 $\uparrow$ 5.6	54.0 $\uparrow$ 3.0	36.7 $\uparrow$ 2.4
<b>Agent frameworks</b>				
AutoCodeRover	GPT-4	62.3	42.3	29.0
+REPOGRAPH	GPT-4	69.0 $\uparrow$ 4.7	45.7 $\uparrow$ 3.4	31.7 $\uparrow$ 2.7
SWE-agent	GPT-4o	61.7	46.3	32.3
+REPOGRAPH	GPT-4o	67.3 $\uparrow$ 4.6	49.3 $\uparrow$ 3.0	35.0 $\uparrow$ 2.7

## 5.1 Localization Coverage

A crucial step in issue resolution is accurately identifying the correct locations within the code that require modification. Proper localization is essential, as it forms the foundation for generating an effective and accurate patch. Without this step, the quality of the fix may be compromised, leading to incomplete or incorrect solutions. In three granularity, we compute the percentage of problems where the edit locations match the ground truth patch. Namely, file-level, function-level, and line-level. We report that a patch contains the correct location if it edits a *superset* of all locations in the ground truth patch.

Table 3 presents the results of our analysis. We observed that integrating REPOGRAPH with all baseline methods significantly improves file-level accuracy, whereas the enhancement of accuracy in line-level is comparatively modest. This result aligns with our expectations, as file-level localization is the most coarse-grained, making it inherently easier to improve. In contrast, line-level localization, being the most fine-grained, poses a greater challenge due to its need for more precise identification of code segments. Additionally, we found that although line-level accuracy improvements are more pronounced for agent frameworks, their overall resolve rate is lower than that of procedural frameworks, as shown in Table 2. This discrepancy can be attributed to the fact that localization, while necessary for generating a final patch, is insufficient. The success of the final revision still heavily relies on the underlying capabilities of LLMs. Agent frameworks, designed to operate in a trial-and-error fashion, are particularly susceptible to error accumulation. As these frameworks iter-

atively refine their patches, small inaccuracies in earlier localization steps can propagate and magnify throughout the process, ultimately reducing the overall resolve rate. Procedural frameworks, on the other hand, follow a more structured and deterministic approach to localization and patch generation. They typically localize and fix issues in a single, more direct step, which can help mitigate the compounding of errors.

## 5.2 Investigation of REPOGRAPH Variants

In this section, we investigate the efficacy of various combinations of sub-graph retrieval and integration techniques as outlined in Section 3. We explore two sub-graph retrieval variants and two integration methods. Specifically, for sub-graph retrieval, we index  $k$ -hop ego-graphs where  $k$  is set to 1 and 2. We limit our exploration to  $k$  values up to 2 due to the extensive context required for integration and the potential introduction of noise or irrelevant nodes for  $k \geq 3$ . For the integration methods, we employ two distinct approaches: (i) directly flattening the textual sub-graph by explicitly detailing the relationships between the search term and its neighboring nodes, and (ii) leveraging an LLM first to summarize the sub-graph in terms of the core modules and salient dependencies, before proceeding with further processing. Detailed implementations of these variants and prompts used can be found in Appendix B.

We begin by presenting some statistics for REPOGRAPH and its various configurations. Performance evaluations are conducted on Agentless with REPOGRAPH integrated as a plug-in module. Table 4 reports the number of nodes and edges within the (sub)-graphs. Notably, the average number of nodes and edges for REPOGRAPH across the SWE-bench dataset is quite substantial, featuring over 1,000 nodes and 25,000 edges. This highlights the comprehensive nature of the constructed structure. For the different variants, when  $k = 1$ , the information within REPOGRAPH is concentrated around the search term, resulting in an average of 11.6 nodes and 37.1 edges. As  $k$  increases to 2, the retrieved ego-graph expands exponentially, reaching an average of 54.5 nodes and 89.9 edges. Moreover, directly flattening the retrieved ego-graph often significantly increases the token count, frequently reaching several thousand tokens. However, utilizing an LLM as an additional summarizer greatly reduces the token count, typically around a thousand.



Table 4: The number of nodes, edges, and tokens of REPOGRAPH and its variants. For different retrieval and integration variants, we report the average number on the test set. “summ.” refers to the summarized version by LLMs of the retrieved ego-graph.

Metrics	REPOGRAPH	1-hop + flatten	1-hop + summ.	2-hop + flatten	2-hop + summ.
# Nodes	1419.3	11.6	11.6	54.5	54.5
# Edges	26392.1	37.1	37.1	89.9	89.9
# tokens	-	2310.7	717.5	10505.3	1229.2
<i>resolve rate</i>	-	29.67	28.33	26.00	28.67

Table 5: Results on subset of CrossCodeEval with GPT-4o as the backbone LLM.

Methods	Code Match		Identifier Match	
	EM	ES	EM	F1
GPT-4o	10.8	59.4	16.7	48.2
+REPOGRAPH	28.5	68.3	36.1	61.9

Table 4 presents the resolve rates for four variants of our method. Notably, while the 2-hop variant incorporates additional information, directly flattening this information results in the poorest performance, with a resolve rate of 26.00%, even lower than the original baseline. In contrast, incorporating summarization via the LLM significantly alleviates context length constraints and enhances information organization, thereby improving performance. For the 1-hop variant, however, we observed that adding summarization actually degrades performance, reducing the resolve rate to 28.33%. We hypothesize that this occurs because the flat 1-hop ego-graph already contains comprehensive information that fits within the LLM’s context window; thus, summarization may introduce inevitable information loss.

### 5.3 Transferability Test

To demonstrate the representational power of REPOGRAPH for repositories and its transferability to tasks requiring an understanding of repository structures, we conducted experiments using the CrossCodeEval benchmark (Ding et al., 2023). CrossCodeEval is a code completion benchmark designed to emphasize numerous cross-file dependencies. The original dataset consists of 2,665 samples derived from real-world repositories. Due to resource constraints, we randomly selected 500 samples from CrossCodeEval, focusing on problems using Python as the evaluation programming language. For evaluation metrics, we follow the settings in the original paper and measure performance with code match and identifier match met-

rics, assessing accuracy with exact match (EM), edit similarity (ES), and F1 scores. In our experiments, the search terms are determined to be the function within which the current line is to be completed.

Table 5 demonstrates the results on CrossCodeEval using GPT-4o as the backbone language model. The original GPT-4o model struggles with repository-level tasks, evidenced by an EM score of only 10.8 for code matching and 16.7 for identifier matching. These results indicate significant limitations in handling code structure and variable usage in a broader repository context. However, with the integration of the REPOGRAPH method, there is a substantial improvement across all metrics. Code Match EM improves dramatically to 28.5, while identifier match EM more than doubles to 36.1. Similarly, the F1 score for identifier matching jumps from 48.2 to 61.9, and the ES for code increases from 59.4 to 68.3. These improvements suggest that incorporating repository-level knowledge, as facilitated by REPOGRAPH, greatly enhances the model’s ability to understand and generate more contextually accurate and semantically consistent code.

### 5.4 Error Analysis

We want to compare REPOGRAPH with the corresponding baselines to see the distribution of resolved cases and analyze the error reasons for unresolved cases. We plot a Venn diagram for representative methods in both procedural and agent frameworks in Figure 3, respectively. We manually examined all the unique error cases and defined three error categories. **(i) Incorrect localization** refers to the failure in accurately identifying code snippets, **(ii) contextual misalignment** happens when the generated patch fails to align with the broader context of the codebase, and **(iii) regressive fix** introduces new issues in resolving the original issues. More examples are in Appendix D.

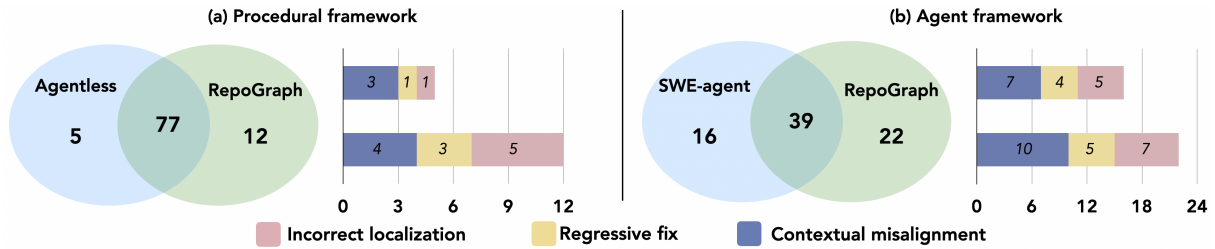


Figure 3: Venn diagram of REPOGRAPH and baseline methods on (a) procedural framework and (b) agent framework. We also plot the error distribution of failing cases against counterparts, e.g., detailed error distribution of 12 cases REPOGRAPH succeeds while Agentless fails.

We found that the improvement in agent frameworks is more complementary than procedural frameworks, with larger uniquely resolved cases of 22 compared with 12 in procedural frameworks. Together, they make even larger performance of 31.33% and 22.33%. The reason could also be attributed to the determinism of procedural frameworks. As a plug-in module, REPOGRAPH tends to make modifications on existing deterministic processes, resulting in larger overlaps in resolved issues compared with baselines. Agent frameworks, on the other hand, have quite different action distributions with REPOGRAPH as a plug-in (please refer to Appendix A.2). Therefore, the uniquely resolved cases are more compared with procedural frameworks. For error distributions, contextual misalignment is the most prevalent error type, followed by incorrect localization and regressive fixes for all methods, suggesting that while localization is often correct, the applied solutions may regress or fail to integrate contextually. The phenomenon also echoes with conclusion obtained in Section 5.1. This is intuitive as all the existing methods focus on providing a more comprehensive and desired context for LLMs to solve the task, which fundamentally depends on the power of backbone LLMs. We also found that when integrated with REPOGRAPH, the proportion of error types “incorrect localization” and “contextual misalignment” largely decreases, indicating that REPOGRAPH is specifically useful in aggregating the related contexts of the current to-be-fixed issue.

## References

Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M. Dai, Anja Hauth, Katie Millican, David Silver, Slav Petrov, Melvin Johnson, Ioannis Antonoglou, Julian Schrittwieser, Amelia Glaese, Jilin Chen, Emily Pitler, Timothy P. Lillcrap, Angeliki Lazaridou, Orhan Firat, James Molloy,

Michael Isard, Paul Ronald Barham, Tom Hennigan, Benjamin Lee, Fabio Viola, Malcolm Reynolds, Yuanzhong Xu, Ryan Doherty, Eli Collins, Clemens Meyer, Eliza Rutherford, Erica Moreira, Kareem Ayoub, Megha Goel, George Tucker, Enrique Piqueras, Maxim Krikun, Iain Barr, Nikolay Savinov, Ivo Danihelka, Becca Roelofs, Anaïs White, Anders Andreassen, Tamara von Glehn, Lakshman Yagati, Mehran Kazemi, Lucas Gonzalez, Misha Khalman, Jakub Sygnowski, and et al. 2023. [Gemini: A family of highly capable multimodal models](#). *CoRR*, abs/2312.11805.

Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. 2021. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*.

Ramakrishna Bairi, Atharv Sonwane, Aditya Kanade, Arun Iyer, Suresh Parthasarathy, Sriram Rajamani, B Ashok, and Shashank Shet. 2024. Codeplan: Repository-level coding using llms and planning. *Proceedings of the ACM on Software Engineering*, 1(FSE):675–698.

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebggen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. 2021. [Evaluating large language models trained on code](#).

Wei Cheng, Yuhan Wu, and Wei Hu. 2024. Dataflow-guided retrieval augmentation for repository-level code completion. In *ACL*.

Cognition. 2024. [Devin](#).

Yangruibo Ding, Zijian Wang, Wasi Uddin Ahmad, Hantian Ding, Ming Tan, Nihal Jain, Murali Krishna Ramanathan, Ramesh Nallapati, Parminder Bhatia, Dan Roth, and Bing Xiang. 2023. [Crosscodeeval: A diverse and multilingual benchmark for cross-file code completion](#). In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurélien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Rozière, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Al-lonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Grégoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel M. Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, and et al. 2024. [The llama 3 herd of models](#). *CoRR*, abs/2407.21783.

Angela Fan, Beliz Gokkaya, Mark Harman, Mitya Lyubarskiy, Shubho Sengupta, Shin Yoo, and Jie M Zhang. 2023. Large language models for software engineering: Survey and open problems. In *2023 IEEE/ACM International Conference on Software Engineering: Future of Software Engineering (ICSE-FoSE)*, pages 31–53. IEEE.

Paul Gauthier. 2024. Aider is ai pair programming in your terminal. <https://aider.chat/>.

Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, Y. Wu, Y.K. Li, Fuli Luo, Yingfei Xiong, and Wenfeng Liang. 2024. [Deepseek-coder: When the large language model meets programming – the rise of code intelligence](#).

Xinyi Hou, Yanjie Zhao, Yue Liu, Zhou Yang, Kailong Wang, Li Li, Xiapu Luo, David Lo, John Grundy,

and Haoyu Wang. 2023. Large language models for software engineering: A systematic literature review. *arXiv preprint arXiv:2308.10620*.

Yuntong Hu, Zhihan Lei, Zheng Zhang, Bo Pan, Chen Ling, and Liang Zhao. 2024. [GRAG: graph retrieval-augmented generation](#). *CoRR*, abs/2405.16506.

Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R Narasimhan. 2024. [SWE-bench: Can language models resolve real-world github issues?](#) In *The Twelfth International Conference on Learning Representations*.

Bowen Jin, Chulin Xie, Jiawei Zhang, Kashob Kumar Roy, Yu Zhang, Zheng Li, Ruirui Li, Xianfeng Tang, Suhang Wang, Yu Meng, and Jiawei Han. 2024. [Graph chain-of-thought: Augmenting large language models by reasoning on graphs](#). In *Findings of the Association for Computational Linguistics ACL 2024*, pages 163–184, Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.

Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474.

Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou, Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, Qian Liu, Evgenii Zheltonozhskii, Terry Yue Zhuo, Thomas Wang, Olivier Dehaene, Mishig Davaadorj, Joel Lamy-Poirier, João Monteiro, Oleh Shliashko, Nicolas Gontier, Nicholas Meade, Armel Zebaze, Ming-Ho Yee, Logesh Kumar Umapathi, Jian Zhu, Benjamin Lipkin, Muhtasham Oblokulov, Zhiruo Wang, Rudra Murthy V, Jason T. Stiller-man, Siva Sankalp Patel, Dmitry Abulkhanov, Marco Zocca, Manan Dey, Zhihan Zhang, Nour Fahmy, Urvasi Bhattacharyya, Wenhao Yu, Swayam Singh, Sasha Luccioni, Paulo Villegas, Maxim Kunakov, Fedor Zhdanov, Manuel Romero, Tony Lee, Nadav Timor, Jennifer Ding, Claire Schlesinger, Hailey Schoelkopf, Jan Ebert, Tri Dao, Mayank Mishra, Alex Gu, Jennifer Robinson, Carolyn Jane Anderson, Brendan Dolan-Gavitt, Danish Contractor, Siva Reddy, Daniel Fried, Dzmitry Bahdanau, Yacine Jernite, Carlos Muñoz Ferrandis, Sean Hughes, Thomas Wolf, Arjun Guha, Leandro von Werra, and Harm de Vries. 2023a. [Starcoder: may the source be with you!](#) *Trans. Mach. Learn. Res.*, 2023.

Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou, Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, et al. 2023b. [Starcoder: may the source be with you!](#) *arXiv preprint arXiv:2305.06161*.

Ming Liang, Xiaoheng Xie, Gehao Zhang, Xunjin Zheng, Peng Di, Hongwei Chen, Chengpeng Wang, Gang Fan, et al. 2024. [Repofuse: Repository-level](#)

- code completion with fused dual context. *arXiv preprint arXiv:2402.14323*.
- Xiangyan Liu, Bo Lan, Zhiyuan Hu, Yang Liu, Zhicheng Zhang, Wenmeng Zhou, Fei Wang, and Michael Shieh. 2024. Codexgraph: Bridging large language models and code repositories via code graph databases. *arXiv preprint arXiv:2408.03910*.
- Shuai Lu, Daya Guo, Shuo Ren, Junjie Huang, Alexey Svyatkovskiy, Ambrosio Blanco, Colin B. Clement, Dawn Drain, Daxin Jiang, Duyu Tang, Ge Li, Lidong Zhou, Linjun Shou, Long Zhou, Michele Tufano, Ming Gong, Ming Zhou, Nan Duan, Neel Sundaresan, Shao Kun Deng, Shengyu Fu, and Shujie Liu. 2021. Codexglue: A machine learning benchmark dataset for code understanding and generation. In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual*.
- Yingwei Ma, Qingping Yang, Rongyu Cao, Binhua Li, Fei Huang, and Yongbin Li. 2024. How to understand whole software repository? *arXiv preprint arXiv:2406.01422*.
- OpenAI. 2023. [GPT-4 technical report](#). *CoRR*, abs/2303.08774.
- Siru Ouyang, Shuohang Wang, Yang Liu, Ming Zhong, Yizhu Jiao, Dan Iter, Reid Pryzant, Chenguang Zhu, Heng Ji, and Jiawei Han. 2023. [The shifted and the overlooked: A task-oriented investigation of user-GPT interactions](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 2375–2393, Singapore. Association for Computational Linguistics.
- L Page. 1999. The page-rank citation ranking: Bringing order to the web.
- Huy N Phan, Hoang N Phan, Tien N Nguyen, and Nghi DQ Bui. 2024. Repohyper: Better context retrieval is all you need for repository-level code completion. *arXiv preprint arXiv:2403.06095*.
- Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, Ivan Evtimov, Joanna Bitton, Manish Bhatt, Cristian Canton-Ferrer, Aaron Grattafiori, Wenhan Xiong, Alexandre Défossez, Jade Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier, Thomas Scialom, and Gabriel Synnaeve. 2023. [Code llama: Open foundation models for code](#). *CoRR*, abs/2308.12950.
- Baptiste Rozière, Marie-Anne Lachaux, Lowik Chanussot, and Guillaume Lample. 2020. [Unsupervised translation of programming languages](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- Disha Shrivastava, Hugo Larochelle, and Daniel Tarlow. 2023. Repository-level prompt generation for large language models of code. In *International Conference on Machine Learning*, pages 31693–31715. PMLR.
- Michele Tufano, Cody Watson, Gabriele Bavota, Massimiliano Di Penta, Martin White, and Denys Poshyvanyk. 2019. An empirical study on learning bug-fixing patches in the wild via neural machine translation. *ACM Transactions on Software Engineering and Methodology (TOSEM)*, 28(4):1–29.
- Pengda Wang, Zilin Xiao, Hanjie Chen, and Frederick L. Oswald. 2024a. [Will the real linda please stand up...to large language models? examining the representativeness heuristic in LLMs](#). In *First Conference on Language Modeling*.
- Xingyao Wang, Boxuan Li, Yufan Song, Frank F Xu, Xiangru Tang, Mingchen Zhuge, Jiayi Pan, Yueqi Song, Bowen Li, Jaskirat Singh, et al. 2024b. Open Devin: An open platform for ai software developers as generalist agents. *arXiv preprint arXiv:2407.16741*.
- Di Wu, Wasi Uddin Ahmad, Dejiao Zhang, Murali Krishna Ramanathan, and Xiaofei Ma. 2024. Repoforger: Selective retrieval for repository-level code completion. *arXiv preprint arXiv:2403.10059*.
- Chunqiu Steven Xia, Yinlin Deng, Soren Dunn, and Lingming Zhang. 2024. Agentless: Demystifying llm-based software engineering agents. *arXiv preprint*.
- John Yang, Carlos E Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik Narasimhan, and Ofir Press. 2024. Swe-agent: Agent-computer interfaces enable automated software engineering. *arXiv preprint arXiv:2405.15793*.
- Ziyu Yao, Jayavardhan Reddy Peddamail, and Huan Sun. 2019. Coacor: Code annotation for code retrieval with reinforcement learning. In *The world wide web conference*, pages 2203–2214.
- Tongxin Yuan, Zhiwei He, Lingzhong Dong, Yiming Wang, Ruijie Zhao, Tian Xia, Lizhen Xu, Binglin Zhou, Fangqi Li, Zhuosheng Zhang, Rui Wang, and Gongshen Liu. 2024. [R-judge: Benchmarking safety risk awareness for llm agents](#). *Preprint*, arXiv:2401.10019.
- 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](#). In *The 2023 Conference on Empirical Methods in Natural Language Processing*.
- Yuntong Zhang, Haifeng Ruan, Zhiyu Fan, and Abhik Roychoudhury. 2024. [Autocoderover: Autonomous program improvement](#). *Preprint*, arXiv:2404.05427.

## A Detailed implementations for each baseline method

This section details the implementation of all four methods in procedural and agent research lines mentioned in Section 4.1.

### A.1 Procedural frameworks

Figure 4 and Figure 5 illustrate the instructions we used for procedural frameworks, including localization and edition, respectively.

We flattened the context of the retrieved ego-graph from REPOGRAPH into the template of the instructions. Specifically, in both “localization” and “edition” stages, REPOGRAPH is flattened in the part of Function/Class Dependencies so that the LLMs could better understand its context.

### A.2 Agent frameworks

LLM agents have a wide range of applications (Wang et al., 2024a; Yuan et al., 2024). We list all the instructions in every step of agent frameworks. The overall and system instruction is shown in Figure 6.

The system instruction defines the task setting and the template for each response. We add “search\_repopgraph” as a new action for the agent to use in the `command_docs`, with its signature listed in Figure 7.

We also plot the frequency for action invoked for both SWE-agent and SWE-agent with REPOGRAPH in Figure 8. We can see that with REPOGRAPH, the maximum turn of finishing the task is reduced from 38 to 35. Also, we computed the average turn to finish the task, which demonstrates a similar trend of 21.47 to 19.12, a significant improvement in efficiency while maintaining effectiveness. We also observe that the action “search\_repopgraph” is invoked mostly in the first 15 rounds of conversation with LLMs. It is more precise than the original action of “search\_dir”, “search\_file”, and “find\_file”.

## B Detailed implementations for each REPOGRAPH variant

In this section, we provide the implementations for variants of REPOGRAPH mentioned in Section 5.2.

In addition to directly flattening the retrieved ego-graph, we propose leveraging large language models (LLMs) to first summarize the context. The

full prompt, along with sample input and output, is provided in Figure 9.

## C Additional Results

### C.1 Resolve rate by repository

We plot the results of the resolve rate in terms of repository distribution in Figure 10. From the figure, it is clear that the resolution of issues varies significantly across different repositories. Notably, the Django and Sympy repositories have the most unresolved issues, with 75 and 61 unresolved issues, respectively. This may indicate a higher level of complexity in the issues or a larger backlog compared to the other repositories. On the other hand, Django has the highest number of resolved issues, with 39 cases. This highlights a strong effort to address issues, even though the unresolved count is still high. Sympy follows closely with 16 resolved issues, suggesting a similar trend. Other repositories like Scikit-learn, Sphinx, and Matplotlib have comparatively fewer issues overall, but their resolve rates are more balanced. For instance, Sphinx shows a ratio of 13 resolved to 3 unresolved issues, reflecting a more consistent effort in issue resolution.

### C.2 Resolve rate by time

We plot the results of the resolve rate in terms of the distribution of releasing time for repositories in Figure 11. Most of the issues were observed in recent years, starting from 2018, with a substantial increase in the total number of issues identified after 2018. In the early years (2012-2016), the number of issues remained relatively low, with both resolved and unresolved counts being minimal. Starting in 2017, there has been a noticeable increase in unresolved issues, with 13 unresolved and only 3 resolved issues. In 2019, the number of resolved and unresolved issues significantly increased, with 19 resolved out of 59 issues. This trend continued to rise until 2020, where 47 issues remained unresolved, and only 17 were resolved, marking the year with the highest number of unresolved issues in the dataset. By 2021 and 2022, the number of unresolved issues slightly decreased, while the resolve rate increased compared to 2020. This suggests an improvement in the system’s ability to address issues in these years. In 2023, although the total number of issues dropped to 30,

<b>Instruction</b>	Please review the following GitHub problem description and relevant files, and provide a set of locations that need to be edited to fix the issue. You will also be given a list of function/class dependencies to help you understand how functions/classes in relevant files fit into the rest of the codebase. The locations can be specified as class names, function or method names, or exact line numbers that require modification.		
<b>Template</b>	<pre> ### GitHub Problem Description ### {problem_statement}  ### Function/Class Dependencies ### {repo_graph}  ### Please provide the class name, function or method name, or the exact line numbers that need to be edited. </pre>	<pre> ### Related Files ### {file_contents} </pre>	
<b>Demonstrations</b>	<pre> ### Examples: ``` full_path1/file1.py          full_path2/file2.py          full_path3/file3.py line: 10                    function: MyClass2.my_method  function: my_function class: MyClass1             line: 12                    line: 24 line: 51                    line: 156 ``` Return just the location(s) </pre>		

Figure 4: Instructions used in the procedural framework to localize to detailed files and code lines of edition.

the proportion of resolved issues remained strong, with 8 resolved out of 22 issues.

## D Examples for error analyses

To help better understand the error category listed in Section 5.4, we provide one example for each category in Figure 12, Figure 13, and Figure 14.

## E Limitations and Future Work

(i) We only explored proprietary LLMs, i.e., GPT-4 series. Due to the poor performance of open-source models on this challenging task, we opted for proprietary models that have demonstrated superior results in code-related tasks. However, a comprehensive evaluation of open-source models such as Llama (Dubey et al., 2024) could be a valuable direction for future work, particularly as these models continue to improve.

(ii) Experiments were only conducted on the Lite set due to the high cost of running large-scale experiments with proprietary models. Exploring more efficient model deployment strategies and alternative cost-effective options for running experiments on larger datasets will be essential for broader applicability.

(iii) Although REPOGRAPH could be adapted to support other programming languages by adjusting

the parsing schemes in the implementation, we only explored Python in our main experiments. Future work could extend this approach to other widely used programming languages, such as JavaScript, Java, or C++, to evaluate the generalizability of our methodology across different programming paradigms.

## F Impact Statement

The impact of this paper lies in its substantial contribution to enhancing the capability of AI-driven software engineering, particularly with respect to repository-level code understanding. The introduction of REPOGRAPH not only significantly improves Large Language Models (LLMs) in navigating and comprehending entire codebases, but also showcases the potential of integrating repository-wide structures into AI workflows. By extending the scope from function-level tasks to holistic repository management, REPOGRAPH pushes the boundaries of AI’s utility in modern software engineering. This advancement opens new opportunities for using LLMs in complex engineering tasks such as automated debugging, repository maintenance, and large-scale refactoring. Furthermore, by highlighting the importance of repository-level context for accurate code generation and maintenance, the paper sets a new trajectory for future

<b>Instruction</b>	<p>We are currently solving the following issue within our repository. Here is the issue text:  {problem_statement}</p> <p>Below are some code segments, each from a relevant file. One or more of these files may contain bugs.  {content}</p> <p>To help you better understand the contexts of the code segments, we provide a set of dependencies of the code segments. The dependencies reflect how the functions/classes in the code segments are referenced in the codebase.  {dependencies}</p>
<b>Template</b>	<p>Please first localize the bug based on the issue statement, and then generate *SEARCH/REPLACE* edits to fix the issue.</p> <p>Every *SEARCH/REPLACE* edit must use this format:</p> <ol style="list-style-type: none"> <li>1. The file path</li> <li>2. The start of search block: &lt;&lt;&lt;&lt;&lt;&lt; SEARCH</li> <li>3. A contiguous chunk of lines to search for in the existing source code</li> <li>4. The dividing line: =====</li> <li>5. The lines to replace into the source code</li> <li>6. The end of the replace block: &gt;&gt;&gt;&gt;&gt;&gt; REPLACE</li> </ol>
<b>Demonstrations</b>	<p>Here is an example:</p> <pre> ### mathweb/flask/app.py &lt;&lt;&lt;&lt;&lt;&lt; SEARCH from flask import Flask ===== import math from flask import Flask &gt;&gt;&gt;&gt;&gt;&gt; REPLACE </pre> <p>Please note that the *SEARCH/REPLACE* edit REQUIRES PROPER INDENTATION. If you would like to add the line ' print(x)', you must fully write that out, with all those spaces before the code!  Wrap the *SEARCH/REPLACE* edit in blocks ```python...```.</p>

Figure 5: Instruction used for fixing an issue based on the identified locations in certain template.

research in AI and software engineering. It encourages deeper exploration of AI's ability to not only write code but also understand and manage large-scale software projects more efficiently. We foresee minimal risks or negative societal impacts from this work. All datasets and benchmarks used in the evaluation are publicly available, and we adhered to their respective licenses. Additionally, REPOGRAPH has been open-sourced, making it accessible to the research community, particularly to groups with limited access to extensive computing resources, thus fostering broader adoption and further development.

**SETTING:**

You are an autonomous programmer, and you're working directly in the command line with a special interface. The special interface consists of a file editor that shows you (WINDOW) lines of a file at a time. In addition to typical bash commands, you can also use the following commands to help you navigate and edit files.

**COMMANDS:**

[{command\\_docs}](#)

Please note that THE EDIT COMMAND REQUIRES PROPER INDENTATION.

If you'd like to add the line ' print(x)' you must fully write that out, with all those spaces before the code! Indentation is important and code that is not indented correctly will fail and require fixing before it can be run.

**RESPONSE FORMAT:**

Your shell prompt is formatted as follows:

(Open file: <path>) <cwd> \$

You need to format your output using two fields; discussion and command.

Your output should always include `_one_ discussion` and `_one_ command` field EXACTLY as in the following example:

**DISCUSSION**

First I'll start by using ls to see what files are in the current directory. Then maybe we can look at some relevant files to see what they look like.

ls -a


You should only include a \*SINGLE\* command in the command section and then wait for a response from the shell before continuing with more discussion and commands. Everything you include in the DISCUSSION section will be saved for future reference.

If you'd like to issue two commands at once, PLEASE DO NOT DO THAT! Please instead first submit just the first command, and then after receiving a response you'll be able to issue the second command.

You're free to use any other bash commands you want (e.g. find, grep, cat, ls, cd) in addition to the special commands listed above.

However, the environment does NOT support interactive session commands (e.g. python, vim), so please do not invoke them.

Figure 6: The signature of our new action "search\_repograph" for agent frameworks.

 **search\_repograph**

**docstring:** searches in the current repository with a specific function or class, and returns the def and ref relations for the search term.

**signature:** search\_repo <search\_term>

**arguments:**

- search\_term (string) [required]: function or class to look for in the repository.

Figure 7: The signature of our new action "search\_repograph" for agent frameworks.

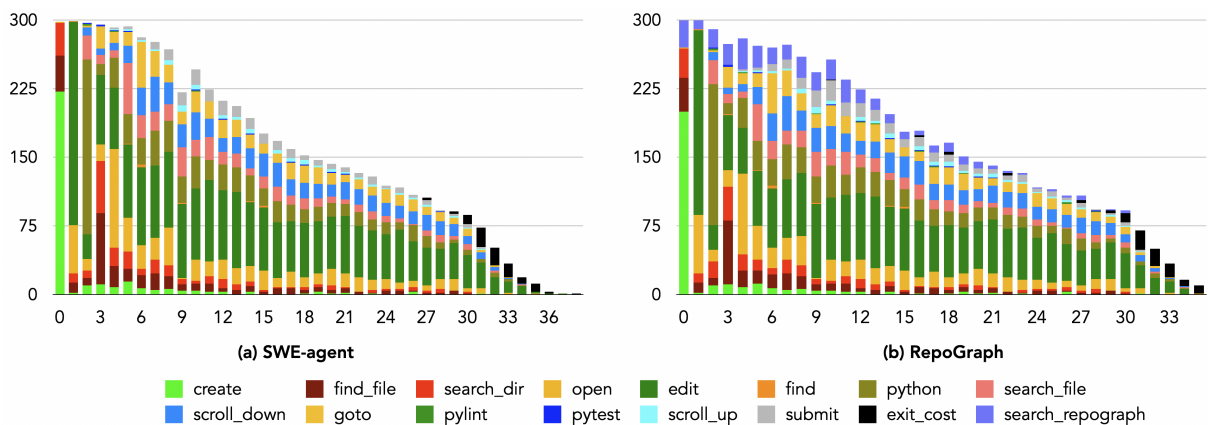


Figure 8: The frequency with which actions are invoked at each turn by (a) SWE-agent and (b) SWE-agent w/REPOGRAPH.



Instruction

You will be given a list of code lines and their meta data, please summarize them in terms of the most salient function/classes and the dependencies, including invocation and inheritance.

For each instance, please keep the following output format:

**\*\* Key Classes and Functions:\*\***

- [Class/Function name]:
  - [Description of the class/function]

**\*\*Core Dependencies:\*\***

- [Class/Function name]:
  - [Description of the dependency]

---

Input

**location:** astropy/modeling/core.py **lines** 1765 - 1885 **name:** render  
**contents:** \n\nclass Model(metaclass=\_ModelMeta):\n def render(self, out=None, coords=None):\n \n\n Evaluate a model at fixed positions, respecting the `bounding\_box`.\n\n The key difference relative to evaluating the model directly is that\n this method is limited to a bounding box if the Model.bounding\_box\n attribute is set.\n\n Parameters\n out : numpy.ndarray, optional\n ...

**location:** astropy/modeling/core.py **lines** 4537 - 4633 **name:** render\_model  
**contents:** \ndef render\_model(model, arr=None, coords=None):\n \n\n Evaluates a model on an input array. Evaluation is limited to\n a bounding box if the Model.bounding\_box attribute is set.\n\n Parameters\n -----\n model : Model\n Model to be evaluated.\n arr : numpy.ndarray, optional\n Array on which the model is evaluated.\n coords : array-like, optional\n Coordinate arrays mapping to arr, such that\n arr[coords] == arr.\n\n Returns:...

**location:** astropy/nddata/utills.py **lines** 546 - 654 **name:** \_\_init\_\_  
**contents:** \nclass Cutout2D:\n def \_\_init\_\_(\n self, data, position, size, wcs=None, mode=\"trim\", fill\_value=np.nan, copy=False\n ):\n if wcs is None:\n wcs = getattr(data, \"wcs\", None)\n\n if isinstance(position, SkyCoord):\n if wcs is None:\n raise ValueError(\"wcs must be input if position is a SkyCoord\")\n\n position = skycoord\_to\_pixel(position, wcs, mode=\"all\") # (x, y)\n\n if np.isscalar(size):\n size = np.repeat(size, 2) ...

**location:** astropy/modeling/core.py **lines** 3872 - 3986 **name:** render  
**contents:** \nclass CompoundModel(Model):\n def render(self, out=None, coords=None):\n \n\n Evaluate a model at fixed positions, respecting the bounding\_box.\n\n The key difference relative to evaluating the model directly is that\n this method is limited to a bounding box if the Model.bounding\_box\n attribute is set. ...

---

Output

**Key Classes and Functions:**

- Model (Class):
  - Location:** astropy/modeling/core.py (line 1765-1865)
  - Function:** This is the base class for models in the astropy library. It contains a `render` method, which evaluates the model at fixed positions while respecting the `bounding_box` attribute. This method ensures that the evaluation occurs over a defined region, with support for various inputs like coordinates and output arrays.
  - Invocation:** The `Model` class uses a custom metaclass `_ModelMeta`. The render method works with numpy arrays and uses helper methods like `extract_array` and `add_array` for array manipulation.
- Cutout2D (Class):
  - Location:** astropy/nddata/utills.py (line 546-654)
  - Function:** A utility class that creates a cutout from a larger 2D array (like an image) centered around a specific position. The `__init__` method handles coordinate transformations, shape adjustments, and manages World Coordinate System (WCS) integration for astronomical data.
  - Invocation:** .....
- `render_model` (Function): .....

**Core Dependencies:**

- `extract_array`:
  - Extracts a portion of a 2D or 3D array based on the provided coordinates and shape, allowing model evaluations to occur within specific sub-regions, used by `Model.render`, `CompoundModel.render`, and `render_model`.
- `skycoord_to_pixel`:
  - Astronomical utilities for converting between celestial coordinates and pixel positions, used in `Cutout2D`
- WCS:
  - .....

Figure 9: Instruction used for summarizing the flattened ego-graph.

Repository	Resolved	Unresolved
sphinx	3	13
django	39	75
pylint	0	1
scikit-learn	10	13
sympy	16	61
flask	0	3
matplotlib	6	17
xarray	0	5
seaborn	2	2
pytest	7	10
requests	4	2
astropy	0	5

Figure 10: Distribution of issues resolved by Agentless +REPOGRAPH plotted in terms of different repositories.

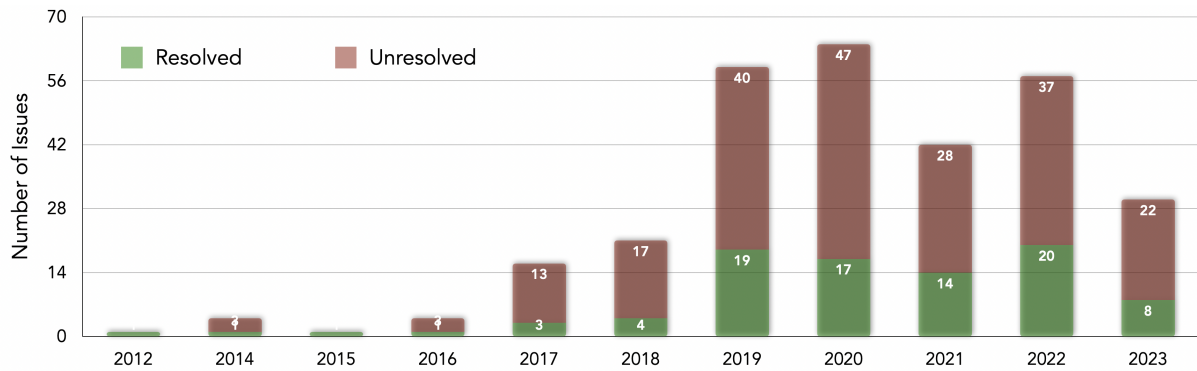


Figure 11: Distribution of issues resolved by Agentless +REPOGRAPH plotted in terms of release time.



**Issue**

Correct expected format in invalid DurationField error message

**Description**

If you enter a duration "14:00" into a duration field, it translates to "00:14:00" which is 14 minutes. The current error message for invalid DurationField says that this should be the format of durations: "[DD] [HH:[MM:]]ss[.uuuuuu]". But according to the actual behaviour, it should be: "[DD] [[HH:]MM:]ss[.uuuuuu]", because seconds are mandatory, minutes are optional, and hours are optional if minutes are provided. This seems to be a mistake in all Django versions that support the DurationField. Also the duration fields could have a default help\_text with the requested format, because the syntax is not self-explanatory.

```
diff --git a/django/forms/fields.py b/django/forms/fields.py
index 935044e..bfb3d28 100644
--- a/django/forms/fields.py
+++ b/django/forms/fields.py
@@ -468,9 +468,10 @@ class DateTimeField(BaseTemporalField):

class DurationField(Field):
    default_error_messages = {
-        'invalid': _('Enter a valid duration.'),
+        'invalid': _('Enter a valid duration in the format [DD]
+ [[HH:]MM:]ss[.uuuuuu].'),
        'overflow': _('The number of days must be between
{min_days} and {max_days}.')
    }
+    default_help_text = _('Enter a duration in the format [DD]
+ [[HH:]MM:]ss[.uuuuuu].')

    def prepare_value(self, value):
        if isinstance(value, datetime.timedelta):
```

(a) RepoGraph

```
diff --git a/django/db/models/fields/__init__.py
b/django/db/models/fields/__init__.py
index 77ba99e..de7d841 100644
--- a/django/db/models/fields/__init__.py
+++ b/django/db/models/fields/__init__.py
@@ -1587,7 +1587,7 @@ class DurationField(Field):
    empty_strings_allowed = False
    default_error_messages = {
        'invalid': _("%(value)s' value has an invalid format.
It must be in "
-        "[DD] [HH:[MM:]]ss[.uuuuuu] format.")
+        "[DD] [[HH:]MM:]ss[.uuuuuu] format.")
    }
    description = _("Duration")
```

(b) Groundtruth



**Tests**

test_dumping (model_fields.test_durationfield.TestSerialization)	✓
test_loading (model_fields.test_durationfield.TestSerialization)	✓
test_formfield (model_fields.test_durationfield.TestFormField)	✓
test_exact (model_fields.test_durationfield.TestQuerying)	✓
test_invalid_string (model_fields.test_durationfield.TestValidation)	✗
test_gt (model_fields.test_durationfield.TestQuerying)	✓
test_create_empty (model_fields.test_durationfield.TestSaveLoad)	✓
test_fractional_seconds (model_fields.test_durationfield.TestSaveLoad)	✓
test_simple_roundtrip (model_fields.test_durationfield.TestSaveLoad)	✓

Figure 12: An example of *incorrect localization*. The correct patch modifies the “DurationField” in `django/db/models/fields/__init__.py`. This is the correct place to handle the error message formatting for the “DurationField” used in Django models. REPOGRAPH, however, modifies `django/forms/fields.py`. This file handles Django form fields, not the model fields. While both model and form fields have overlapping behavior, in this case, the correction is required for the model field (DurationField in `__init__.py`), not the form field.

Trailing whitespace in DatasetGroupBy text representation  
 When displaying a DatasetGroupBy in an interactive Python session, the first line of output contains a trailing whitespace. The first example in the documentation demonstrate this:

```
>>> import xarray as xr, numpy as np
>>> ds = xr.Dataset(
...     {"foo": (("x", "y"), np.random.rand(4, 3))},
...     coords={"x": [10, 20, 30, 40], "letters": ("x", list("abba"))})
>>> ds.groupby("letters")
DatasetGroupBy, grouped over 'letters'
2 groups with labels 'a', 'b'.
```



**Issue**

There is a trailing whitespace in the first line of output which is "DatasetGroupBy, grouped over 'letters' ". This can be seen more clearly by converting the object to a string (note the whitespace before '\n'):

```
>>> str(ds.groupby("letters"))
"DatasetGroupBy, grouped over 'letters' \n2 groups with labels 'a', 'b'."
```

While this isn't a problem in itself, it causes an issue for us because we use flake8 in continuous integration to verify that our code is correctly formatted and we also have doctests that rely on DatasetGroupBy textual representation. Flake8 reports a violation on the trailing whitespaces in our docstrings. If we remove the trailing whitespaces, our doctests fail because the expected output doesn't match the actual output. So we have conflicting constraints coming from our tools which both seem reasonable. Trailing whitespaces are forbidden by flake8 because, among other reasons, they lead to noisy git diffs. Doctest want the expected output to be exactly the same as the actual output and considers a trailing whitespace to be a significant difference. We could configure flake8 to ignore this particular violation for the files in which we have these doctests, but this may cause other trailing whitespaces to creep in our code, which we don't want. Unfortunately it's not possible to just add '# NoQA' comments to get flake8 to ignore the violation only for specific lines because that creates a difference between expected and actual output from doctest point of view. Flake8 doesn't allow to disable checks for blocks of code either.

Is there a reason for having this trailing whitespace in DatasetGroupBy representation? Would it be OK to remove it? If so please let me know and I can make a pull request.

```
diff --git a/xarray/core/groupby.py b/xarray/core/groupby.py
index bf63803..7f9a25e 100644
--- a/xarray/core/groupby.py
+++ b/xarray/core/groupby.py
@@ -1011,5 +1011,10 @@ class DatasetGroupBy(GroupBy,
ImplementsDatasetReduce):
    return self.map(lambda ds: ds.assign(**kwargs))

+def __repr__(self):
+    # Original representation with trailing whitespace removed
+    groups = ", ".join(str(k) for k in self.groups)
+    return f"DatasetGroupBy, grouped over
{self._group.name}\n{len(self.groups)} groups with labels
'{groups}'."
+ops.inject_reduce_methods(DatasetGroupBy)
-ops.inject_binary_ops(DatasetGroupBy)
+ops.inject_binary_ops(DatasetGroupBy)
```

(a) RepoGraph

```
diff --git a/xarray/core/groupby.py b/xarray/core/groupby.py
index bf63803..9aac7fc 100644
--- a/xarray/core/groupby.py
+++ b/xarray/core/groupby.py
@@ -436,7 +436,7 @@ class GroupBy(SupportsArithmetic):
    return zip(self._unique_coord.values, self._iter_grouped())

    def __repr__(self):
-    return "{}, grouped over {} \n{} groups with labels
{}.format(
+    return "{}, grouped over {} \n{} groups with labels
{}."format(
        self.__class__.__name__,
        self._unique_coord.name,
        self._unique_coord.size,
```

(b) Groundtruth



**Tests**

xarray/tests/test_groupby.py::test_consolidate_slices	✓
xarray/tests/test_groupby.py::test_multi_index_groupby_sum	✓
xarray/tests/test_groupby.py::test_groupby_map_shrink_groups[obj0]	✓
xarray/tests/test_groupby.py::test_da_groupby_map_func_args	✓
xarray/tests/test_groupby.py::test_groupby_repr[obj0-x]	✗
.....	...
xarray/tests/test_groupby.py::test_groupby_reduce_dimension_error	✓
xarray/tests/test_groupby.py::test_groupby_repr_datetime[obj0]	✗
xarray/tests/test_groupby.py::test_groupby_getitem	✓

Figure 13: An example of *contextual misalignment*. The model-generated patch places the `__repr__` method in the wrong class `DatasetGroupBy`, whereas the correct patch modifies it in the `GroupBy` class. The `__repr__` method should be implemented in the `GroupBy` class because it deals with the general group-by functionality. “`DatasetGroupBy`” is a subclass and doesn’t require a new representation method if “`GroupBy`” already has one. Additionally, when fixing white space trailing, while functionally similar, REPOGRAPH’s generation doesn’t exactly match the style or intention of the correct patch.



**Issue**

expr.atoms() should return objects with no args instead of subclasses of Atom  
'expr.atoms()' with no arguments returns subclasses of 'Atom' in 'expr'. But the correct definition of a leaf node should be that it has no '.args'.

This should be easy to fix, but one needs to check that this doesn't affect the performance.

```
diff --git a/sympy/core/basic.py b/sympy/core/basic.py
index 97eb81a..4fa2f3c 100644
--- a/sympy/core/basic.py
+++ b/sympy/core/basic.py
@@ -500,15 +500,11 @@ class Basic(metaclass=ManagedProperties):
     {!*pi, 2*sin(y + l*pi)}
-    if types:
-        types = tuple(
-            [t if isinstance(t, type) else type(t) for t in types])
+    if not types:
+        result = {expr for expr in preorder_traversal(self) if not expr.args}
     else:
-        types = (Atom,)
-        result = set()
-        for expr in preorder_traversal(self):
-            if isinstance(expr, types):
-                result.add(expr)
+        types = tuple(t if isinstance(t, type) else type(t) for t in types)
+        result = {expr for expr in preorder_traversal(self) if isinstance(expr,
types)}
     return result
```

(a) RepoGraph

```
diff --git a/sympy/core/basic.py b/sympy/core/basic.py
index 97eb81a..f1f411f 100644
--- a/sympy/core/basic.py
+++ b/sympy/core/basic.py
@@ -507,7 +507,7 @@ class
Basic(metaclass=ManagedProperties):
     types = (Atom,)
     result = set()
     for expr in preorder_traversal(self):
-        if isinstance(expr, types):
+        if isinstance(expr, types) and not expr.args:
+            result.add(expr)
     return result
```

(b) Groundtruth



**Tests**

test_literal_evalf_is_number_is_zero_is_comparable	✓
test_subs_with_unicode_symbols	✓
test_free_symbols_empty	✓
test_preorder_traversal	✓
test_atomic	✗
.....	...
test_canonical_variables	✓
test_sizeof	✗
test_CommaOperator	✓

Figure 14: An example of *regressive fix*. The model-generated patch successfully resolves the issue which requires the return of objects with no args instead of subclasses of Atom by adding the key code line `if not expr.args`. However, the fix introduces other new issues such as the size of return objects, as exemplified in the unit tests.