# SWE-EXT: EXTENDING AND SCALING AUGMENTED DATA FOR REPOSITORY-LEVEL CODING TASKS

# **Anonymous authors**

Paper under double-blind review

#### **ABSTRACT**

Repository-level benchmarks such as SWE-Bench have highlighted the challenges of scaling language models to complex software engineering tasks. However, current training data remains narrow in scope, primarily focusing on monolingual issue resolving and feature implementation. In this work, we introduce SWE-Ext, a large-scale effort to extend and scale augmented data for repository-level coding tasks. SWE-Ext broadens existing data along two key dimensions: multilingual coverage (spanning 10 languages) and an auxiliary code completion task. We uncover distinct transfer mechanisms: data from other programming languages provides transferable signals that generally enhance localization and editing capabilities in single-language (Python) settings, while code completion data strengthens code editing capabilities, particularly for feature implementation tasks requiring substantial new code generation. These extensions yield consistent improvements on Python repository-level benchmarks like SWE-Bench and FEA-Bench. Our method offers a simple yet effective way to leverage more opensource data for advancing repository-level code models.

# 1 Introduction

Repository-level coding tasks have emerged as a critical frontier in code generation, shifting focus from isolated script-level challenges to comprehensive software engineering workflows in full repositories. Early benchmarks emphasized standalone function synthesis, such as HumanEval (Chen et al., 2021) for Python problems and MBPP (Austin et al., 2021) for basic algorithms. However, these evaluations often ignored broader repository contexts, including dependencies, multi-file interactions, and real-world specifications (Xu et al., 2022). Recent advancements have introduced more holistic benchmarks, with SWE-Bench (Jimenez et al., 2024) exemplifying repository-level issue resolution by evaluating models on GitHub pull requests (PRs) tied to verifiable test outcomes. Complementary efforts like FEA-Bench (Li et al., 2025) target feature implementation.

To enhance model capabilities on these demanding tasks, substantial work has explored post-training strategies. Broadly, these approaches fall into two categories: augmented data and verified data. Augmented data methods construct training instances directly from real-world GitHub PRs, leveraging ground-truth information without requiring execution environments (Xie et al., 2025; Wang et al., 2025). For instance, SWE-Fixer (Xie et al., 2025) curates PR-based data for instruction tuning, while MCTS-Refine (Wang et al., 2025) synthesizes reasoning chains via search algorithms. In contrast, verified data approaches build executable environments to collect agent trajectories, filtering for those that pass unit tests or achieve successful outcomes (Jain et al., 2025; Pan et al., 2024; Yang et al., 2025b). Techniques like R2E (Jain et al., 2025) and SWE-Gym (Pan et al., 2024) generate validated trajectories in sandbox for real PRs, and SWE-Smith (Yang et al., 2025b) extends synthetic tasks with verification. While verified data ensures high-quality supervision through execution feedback, it demands significant resources for environment setup and scaling. Augmented data, conversely, offers greater flexibility and avoids these complexities, enabling broader exploration of data domains.

Despite their advantages, existing augmented datasets remain limited in scope, predominantly focusing on Python repositories and monolingual issue resolution or feature implementation (Xie et al., 2025). This narrow focus underutilizes the vast diversity of open-source GitHub data, restricting model generalization across languages and task types. In this work, we introduce **SWE-Ext**, a

scalable pipeline to extend and augment repository-level coding data along two orthogonal dimensions: multilingual coverage and auxiliary task inclusion (e.g., code completion). Starting from high-quality PRs crawled from GitHub Archive <sup>1</sup> across ten programming languages (Python, Go, JavaScript, Ruby, PHP, Java, TypeScript, C#, C++, and C), we construct datasets of four complementary sub-tasks: (1) file localization to identify relevant files, (2) component localization to pinpoint functions or methods, (3) code editing to generate patches for code repositories, and (4) code completion derived from short-description PRs to provide code in-filling on fixed-positions. Our contributions are as follows:

- We develop a scalable data extension pipeline that broadens the utility of GitHub PRs for repository-level tasks, enabling extension across multilingual and completion-based domains while preserving baseline effectiveness.
- We demonstrate consistent performance gains on repository-level benchmarks like SWE-Bench and FEA-Bench, with multilingual and completion extensions yielding up to +1.4% and +2.5% improvements on 32B models, respectively, and up to +5.4% on 7B models.
- We provide preliminary validation of strong cross-task and cross-language transfer effects, showing that multilingual or completion data can boost single-language downstream performance and highlighting the value of diverse augmentation for generalization.

# 2 PROBLEM DEFINITION

**Notation.** Let  $\mathbb{R}$  denote a software repository represented as a finite collection of files  $\{f_i\}_{i=1}^N$  together with auxiliary structures (ASTs, file-level skeletons, import graph, etc.). Let q denote a problem statement (e.g., a PR description, an issue body, a failing test, or a natural-language specification). Let  $\mathbb{P}$  be the universe of addressable program positions (file identifiers  $\times$  component spans). A *localization* is a subset  $\ell \subseteq \mathbb{L}$ . The space of edits (patches) is denoted  $\mathbb{P}$ ; an edit (patch)  $\delta \in \mathbb{P}$  is an operator  $p: \mathbb{R} \to \mathbb{R}'$  that applies insert/replace/delete/move operations to positions in  $\mathbb{P}$ . Finally,  $L(\mathbb{R}',q)$  is a task loss measuring how well the edited repository  $\mathbb{R}'$  satisfies q. Note that  $P(\cdot)$  denotes a probability distribution.

**Repository-level coding task.** Given  $(\mathbb{R}, q)$ , the objective of a repository-level coding agent is to produce a patch  $\delta \in \mathbb{P}$  that minimizes the task loss:

$$\hat{\delta} = \arg\min_{\delta \in \mathbb{P}} L(\delta(\mathbb{R}), q), \tag{1}$$

or equivalently, under a probabilistic modeling view,

$$\hat{\delta} = \arg \max_{\delta \in \mathbb{P}} P(\delta \mid \mathbb{R}, q). \tag{2}$$

**Decomposition into localization and editing.** Because modern LLMs cannot feasibly consume the entirety of large repositories at once, we decompose the posterior over patches by marginalizing over possible localizations:

$$P(\delta \mid \mathbb{R}, q) = \sum_{\ell \in \mathbb{L}} P(\ell \mid \mathbb{R}, q) P(\delta \mid \mathbb{R}, q, \ell).$$
 (3)

This decomposition separates (i) a *localization* model that identifies candidate regions  $\ell$  of interest, and (ii) an *editing* model that produces a concrete patch conditional on the selected regions. In practice, localization and editing are performed in separate iterative phases: the system first iterates to identify a sufficiently accurate set of relevant files or components, and only then initiates an editing phase where patches can be generated, potentially with feedback.

**Code completion as a special case.** Code completion corresponds to the special case where the localization  $\ell^*$  is given *a priori* (a fixed position or contiguous span) and the patch is constrained to be a continuation or replacement at that position:

$$\hat{\delta} = \arg \max_{\delta \in \mathbb{P}} P(\delta \mid \mathbb{R}, q, \ell^*). \tag{4}$$

<sup>&</sup>lt;sup>1</sup>https://www.gharchive.org

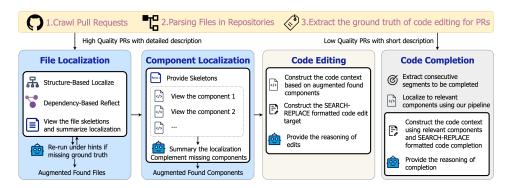


Figure 1: Data collection pipeline for 4 types of tasks in SWE-Ext. We crawl and process pull requests from GH Archive, and then process well described PRs for localization-editing task and process other PRs for completion task.

# 3 SWE-EXT: EXTEND THE REPOSITORY-LEVEL CODING DATA

#### 3.1 Data Collection

Source data and ground truth. We crawl pull requests (PRs) from GH Archive (recent 10 years) and retain only high-quality candidate samples. We keep PRs whose repository has both (i) star count > 100 and (ii) total PR count > 100, and we exclude PRs that modify only non-code files or that were not merged. For each retained PR we collect metadata (PR description, linked issue text if any, timestamps) and patch (diff). We parse repository files with tree-sitter parsers for the following languages: Python, Go, JavaScript, Ruby, PHP, Java, TypeScript, C#, C++, and C. From each file we extract a compact *file skeleton* (top-level classes, functions/methods). Using the patch we record the ground-truth modified files  $\mathbb{F}_{\rm gt}$  and the ground-truth modified components  $\mathbb{C}_{\rm gt}$ . To limit tractability, we retain only PRs with  $||\mathbb{F}_{\rm gt}|| \leq 5$  and  $||\mathbb{C}_{\rm gt}|| \leq 10$ . The problem statement q is the concatenation of the PR description and the linked issue body (when available). Formally each retained example is stored as a tuple

$$s = (\mathbb{R}, q, \mathbb{F}_{gt}, \mathbb{C}_{gt}, \delta_{gt}), \tag{5}$$

where  $\delta_{\rm gt}$  is the ground-truth patch.

**File localization.** To approximate  $P(\ell \mid \mathbb{R}, q)$  at file granularity, we construct file-localization data. We employ a three-stage pipeline driven by a expert model: (i) rank files according to the file tree the problem statement; (ii) expand via file-level dependency to form a candidate set; (iii) re-rank using file skeletons and the expert model to produce a final top-n set  $\mathbb{S}$ . If  $\mathbb{F}_{\text{gt}} \not\subseteq \mathbb{S}$ , we re-run the expert with rationalization hints so that S always covers the ground truth. The resulting dataset pairs  $(\mathbb{R}, q, \mathbb{S})$  provide supervision for learning to predict high-recall file sets that contain the true patch locations, thereby reducing the effective search space for downstream editing.

Component localization. Once file localization is complete, we build component-level (function/class level) localization data to approximate  $P(c \mid \mathbb{R}, q, \ell)$ . We adapt iterative selection methods from CoSIL (Jiang et al., 2025) to allow the expert model to sequentially inspect and select components from candidate files. The process terminates once up to K=10 components are selected. If some ground-truth components are missing, they are inserted into the iterative process and final list. Each example  $(\mathbb{R}, q, \mathbb{F}, \mathbb{S}_{comp})$  thus can serves as a high-recall label for a policy that selects components relevant to q. This supports models in learning to identify fine-grained patch locations and to maximize recall under constrained inspection budgets.

**Code editing.** We construct code-editing data directly from  $\delta_{\rm gt}$ . For each selected localization context  $\ell_{\rm ctx} = \mathbb{S}_{\rm comp}$ , we convert patches into *before/after* code fragments (search–replace pairs  $\delta_{sp}$ ). To enrich supervision, a expert model generates a rationale explaining the patch. Each example  $(\mathbb{R}, q, \ell_{\rm ctx}, \delta_{sp})$  provides supervision for  $P(\delta \mid \mathbb{R}, q, \ell_{\rm ctx})$ . This partition trains models to propose concrete patches given a localized context.

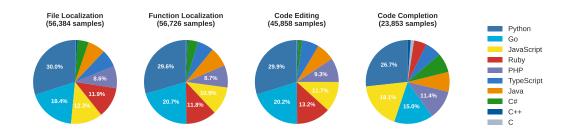


Figure 2: Language distribution of dialogues for 4 sub-tasks in SWE-Ext dataset.

**Code completion.** Finally, we derive a completion dataset from PRs whose problem statements are short (fewer than 50 words), which are possibly not complete instructions for code editing. In these cases, we extract newly added functions or contiguous inserted regions and frame them as completion tasks with a fixed localization  $\ell^*$ . Each example provides supervision for  $P(\delta \mid \mathbb{R}, q, \ell^*)$ . Because completions originate from real patches and are abundant, they supply local-generation signals that improve fluency and act as an effective curriculum for training editing models.

**Summary.** The construction pipeline yields four complementary datasets, each aligned with a distinct conditional distribution in the factorization of  $P(\delta \mid \mathbb{R}, q)$ . By supervising different stages separately, we enable scalable training for repository-level coding tasks while ensuring that the final system can integrate these capabilities into a coherent editing pipeline. For further details, please refer to Appendix D.

#### 3.2 Dataset Characteristics

The SWE-Ext corpus is a *training* dataset constructed from real-world pull requests and intended to supervise models for repository-level coding tasks. The dataset is organized into four complementary tasks: file localization, component (function/method) localization, code editing, and code completion. Table 1 reports core statistics. Below we summarize the most important characteristics and how they relate to model training.

**Origin and construction.** Table 1 presents the statistics of the SWE-Ext dataset. For the three localization and editing tasks (File-Localization, Component-Localization, and Code-Editing), we retain a canonical pool of approximately 56k to 46k PRs. Crucially, these tasks exhibit high data diversity, being sourced from  $4{,}319$  to  $4{,}746$  unique GitHub repositories, which strongly demonstrates the broad applicability of our data augmentation strategy. The Code-Completion task collects newly-added functions or contiguous insertion ranges from  $\sim 23.9k$  PRs, whose short descriptions make end-to-end supervision weak but whose insertions are realistic completion targets.

**Augmentation of localizations.** The **Localizations** column in Table 1 highlights our data augmentation strategy: The large candidate set size (e.g., 4.83 files vs. 1.66 ground truth files for File-Localization) confirms that we incorporate highly relevant, non-ground-truth files or components into the training data. This design is crucial because it makes the dataset more realistic, aligning with scenarios where agents' localization is often relevant but not perfectly precise in real-world software engineering tasks.

Scale, context length and edit complexity. The dataset provides long, realistic contexts: average input contexts (**Dialogue Len**) range from  $\sim 5.9$ k tokens (completion) up to  $\sim 9.8$ k tokens (file-localization), reflecting the long context challenge for real-word repository-level tasks. Average **Answer Len** varies by task ( $\approx 60-955$  tokens). Edit complexity differs across modalities: code-editing patches average  $\sim 32$  lines (explicit before/after pairs), while completion targets are larger on average ( $\sim 88$  lines) because they often correspond to newly-added contiguous code regions.

**Language coverage.** SWE-Ext covers ten languages: Python, Go, JavaScript, Ruby, PHP, Type-Script, Java, C#, C+++, and C. Language distributions (See Figure 2) show consistent dominance

Table 1: Statistics of the dialogue data for 4 tasks in the SWE-Ext dataset. Token counts are computed using <code>Qwen2.5-Coder</code>. **Dialogue Len** represents the total token count of the full dialogue after applying a chat template. **Answer Len** is the average token count of all "assistant" responses within the dialogue. The 'Answer Len' are the final macro-averaged values. The **Patch** metric shows the average number of lines changed in the ground truth patches. For **Localizations**, the values indicate the number of ground truth modified locations versus the total number of locations in the augmented data.

Dataset	Unique repos	Samples (#)	Dialogue Len (Avg tok.)	Answer Len (Avg tok.)	Patch (Avg lines)	Localizations (Avg GT/Total)
File-Localization	4746	56384	9 787.8	552.1	-	1.66/4.83
Component-Localization	4702	56726	8 164.3	59.8	-	2.55/6.32
Code-Editing	4319	45 858	6447.7	955.0	32.0	-
Code-Completion	1 642	23 853	5 929.7	430.1	88.0	-

of Python and Go across modalities (roughly 27--30% and 15--21% respectively), moderate representation for JavaScript, Ruby and PHP ( $\approx 9\text{--}13\%$ ), and smaller but present contributions from Java/TypeScript/C#. Low-resource languages (C, C++) appear at the long tail (< 1--1.1% in most partitions). This multi-language composition enables cross-lingual training and evaluation of model robustness while reflecting real-world repository populations.

# 4 EXPERIMENTS

**Training.** To empirically validate whether our augmented data, derived from real-world GitHub PRs, can enhance the proficiency of large language models on complex repository-level coding tasks, and to explore the synergistic effects of incorporating multilingual and code completion data, we employ a supervised finetuning approach tailored for a multi-turn dialogue format.

**Models.** Our methodology leverages GPT-4o (gpt-4o-2024-05-13) (Hurst et al., 2024) as a data augmentation expert to transform raw GitHub PR data into a multi-task dataset, as detailed in Section 3.1. For the training phase, we utilize <code>Qwen2.5-Coder-Instruct</code> (Hui et al., 2024) as our foundational model. For details of training and models, please refer to Appendix B.

**Agent System.** Our agent system adopts a staged approach. We first perform file and component localization following the same iterative process used for data generation. Subsequent stages, including line-level localization, patch generation, and verification, are executed in an Agentless manner (Xia et al., 2024). We have termed this pipeline **CosAgentless**, signifying the integration of the iterative and fine-grained localization process from CoSIL (Xia et al., 2024) into the standard Agentless inference pipeline.

**Data.** Our dataset is partitioned into a 96% training set and a 4% validation set. For a comprehensive ablation study and a fair comparison with prior works, we define three distinct training configurations. First, the **SWE-Ext-Baseline** model is trained on the standard Python-only data for the first three tasks (localization and editing), consistent with existing agent systems. Second, the **SWE-Ext-Multilingual** model is trained on the first three tasks of our multilingual data. Finally, the **SWE-Ext-Completion** model is trained on Python-only data but across all four tasks. This setup enables a clear analysis of the performance gains derived from extending our data along both multilingual and task-specific dimensions.

**Evaluation.** We evaluate our models on two distinct repository-level coding benchmarks: the well-established SWE-bench (Jimenez et al., 2024) for resolving issues and FEA-Bench (Li et al., 2025) for implementing features. Both the benchmarks only contain Python task instances and the primary metric is the task pass rate (%resolved). For a more granular analysis, we also report the Top-x hit rates, MRR (Mean Reciprocal Rank), and MPP (Mean Precision at Position)(Manning, 2008) for the file and component localization stages. In our evaluation, we use a modified Agentless framework to perform a single complete attempt per task. A task is considered a failure if the process stalls or errors out, resulting in an empty patch.

# 5 RESULTS

# 5.1 SOFTWARE ENGINEERING BENCHMARKS

Model	System	Expert Model	Exec	% Resolved
	Closed Weight Mo	odels		
<b>GPT-40</b> (Hurst et al., 2024)	Agentless	-	-	38.8%
Claude 3.5 Sonnet (Anthropic, 2024)	Agentless	_	-	50.8%
Claude 3.7 Sonnet (Anthropic, 2025a)	SWE-agent	_	-	58.2%
Claude 4 Sonnet (Anthropic, 2025b)	SWE-agent	_	-	72.7%
<b>Llama3-SWE-RL-70B</b> (Wei et al., 2025)	Agentless	-	-	41.0%
	Open Weight Mo	dels		
<b>DeepSeek-V3-671B</b> (Liu et al., 2024)	Agentless	-	-	42.0%
Kimi K2-1TB (Team et al., 2025)	Agentless	-	-	65.8%
<b>Lingma-SWE-GPT-72B</b> (Ma et al., 2024)	SWE-SynInfer	-	_	28.8%
<b>Qwen3-235B-A22B</b> (Yang et al., 2025a)	OpenHands	_	-	34.4%
<b>SWE-gym-32B</b> (Pan et al., 2024)	OpenHands	Hybrid	<b>√</b>	20.6%
<b>R2E-Gym-32B</b> (Jain et al., 2025)	OpenHands	Claude 3.5 Sonnet	<b>√</b>	34.4%
<b>SWE-smith-7B</b> (Yang et al., 2025b)	SWE-agent	Claude 3.7 Sonnet*	<b>✓</b>	15.2%
<b>SWE-smith-32B</b> (Yang et al., 2025b)	SWE-agent	Claude 3.7 Sonnet*	<b>/</b>	40.2%
<b>SWE-fixer-72B</b> (Xie et al., 2025)	SWE-Fixer	_	X	32.8%
<b>SoRFT-Qwen-32B</b> (Ma et al., 2025)	Agentless	Claude 3.5 Sonnet	X	30.8%
MCTS-Refine-32B (Wang et al., 2025)	Agentless	DeepSeek-v3	X	32.4%
	SWE-Ext Mode	rls		
SWE-Ext-Baseline-32B	CosAgentless	GPT-4o	X	31.2%
SWE-Ext-Multilingual-32B	CosAgentless	GPT-40	X	32.6%
SWE-Ext-Completion-32B	CosAgentless	GPT-4o	X	32.2%

Table 2: Resolve rates for existing solutions on **SWE-bench Verified**, collected from (Yang et al., 2025b) and Kimi-K2 (Team et al., 2025) technical reports. **Expert Model** indicates the large language models that generated content during the data construction process. **Exec** indicates whether execution-based feedback is used in the data construction process. All performance numbers are pass@1 (Single attempt using agent systems). \*Indicates the primary data for training is *mainly* generated by the specified expert model.

Model	System	% Resolved						
Zero-Short Inference								
GPT-40 (Hurst et al., 2024) 01 (Jaech et al., 2024) DeepSeek-V3-671B (Liu et al., 2024)	Agentless Agentless Agentless-Lite	9.0% 14.0% 11.0%						
SWE-Ext Models								
SWE-Ext-Baseline-32B SWE-Ext-Multilingual-32B SWE-Ext-Completion-32B	CosAgentless CosAgentless CosAgentless	10.0% 11.5% 12.5%						

Table 3: Resolve rates for existing solutions on **FEA-Bench Lite**, collected from Li et al. (2025). All performance numbers are pass@1 (Single attempt using agent systems).

Competitive performance under limited resources. According to Table 2, although the absolute performance of SWE-Ext models is lower than the strongest closed-weight systems (e.g., Claude 4 Sonnet at 72.7%), our models achieve competitive results among open-weight models even with the baseline data. It is important to note that most higher-performing models either (i) use significantly larger parameter counts, or (ii) are distilled from expert models with higher success rates. In contrast, SWE-Ext only leverages augmented training data without relying on such privileged resources, yet already achieves notable improvements. The results validate our pipeline of data collection.

Consistent improvements with extended data. Table 2 and Table 3 demonstrate that the models trained with SWE-Ext consistently outperform the baseline across both benchmarks. On SWE-Bench Verified, our SWE-Ext-Multilingual/Completion-32B achieve 32.6% and 32.2% resolve rates, respectively, compared to 31.2% of the baseline. On FEA-Bench Lite, similar trends are observed: the model trained with extended data outperform the baseline by up to +2.5%. These results confirm that the proposed data extensions are effective in enhancing model performance.

**Multilingual vs. Completion extension.** We observe distinct patterns between multilingual and completion data extension. Multilingual extension brings more consistent improvements in issue resolving, suggesting that cross-lingual signals help models capture generalizable reasoning strategies. In contrast, completion extension is particularly effective for feature implementation, where generating new functionality requires stronger code completion capabilities. These results highlight that different types of target complementary aspects of repository-level development, and jointly contribute to the overall performance gains.

**Ablation insights.** By comparing the multilingual extension and the completion-based extension, we find that both augmentation strategies yield consistent gains. This suggests that when the amount of original supervision is limited, expanding the dataset in orthogonal directions, either through multilingual variants or through completion-style tasks, can provide complementary signals and enhance model generalization. This result underscores the value of leveraging diverse forms of augmented data for repository-level coding tasks.

Overall, these findings highlight that SWE-Ext contributes a practical approach for improving the software engineering capabilities, providing consistent benefits across different benchmarks.

#### 5.2 LOCALIZATION ANALYSIS

System	Component-level Localization				File-level Localization			
	Hit@1	Hit@3	MAP	MRR	Hit@1	Hit@3	MAP	MRR
CoSIL	43.0	54.3	46.1	48.9	60.7	77.3	69.8	69.4
CosAgentless	47.4	61.8	47.4	54.8	69.0	86.2	76.1	77.7
CosAgentless	55.0	65.4	53.0	60.8	72.8	87.6	77.0	80.2
Cos Agentless	<b>57.2</b> 52.8	<b>68.2</b> 62.0	<b>55.3</b> 51.3	<b>63.1</b> 57.9	<b>75.8</b>	<b>90.2</b>	<b>79.4</b>	<b>82.8</b> 79.8
	CoSIL CosAgentless CosAgentless	Hit@1     CoSIL   43.0     CosAgentless   47.4     CosAgentless   55.0     CosAgentless   57.2	Hit@1         Hit@3           CoSIL         43.0         54.3           CosAgentless         47.4         61.8           CosAgentless         55.0         65.4           CosAgentless         57.2         68.2	Hit@1         Hit@3         MAP           CoSIL         43.0         54.3         46.1           CosAgentless         47.4         61.8         47.4           CosAgentless         55.0         65.4         53.0           CosAgentless         57.2         68.2         55.3	Hit@1         Hit@3         MAP         MRR           CoSIL         43.0         54.3         46.1         48.9           CosAgentless         47.4         61.8         47.4         54.8           CosAgentless         55.0         65.4         53.0         60.8           CosAgentless         57.2         68.2         55.3         63.1	Hit@1         Hit@3         MAP         MRR         Hit@1           CoSIL         43.0         54.3         46.1         48.9         60.7           CosAgentless         47.4         61.8         47.4         54.8         69.0           CosAgentless         55.0         65.4         53.0         60.8         72.8           CosAgentless         57.2         68.2         55.3         63.1         75.8	Hit@1         Hit@3         MAP         MRR         Hit@1         Hit@3           CoSIL         43.0         54.3         46.1         48.9         60.7         77.3           CosAgentless         47.4         61.8         47.4         54.8         69.0         86.2           CosAgentless         55.0         65.4         53.0         60.8         72.8         87.6           CosAgentless         57.2         68.2         55.3         63.1         75.8         90.2	Hit@1         Hit@3         MAP         MRR         Hit@1         Hit@3         MAP           CoSIL         43.0         54.3         46.1         48.9         60.7         77.3         69.8           CosAgentless         47.4         61.8         47.4         54.8         69.0         86.2         76.1           CosAgentless         55.0         65.4         53.0         60.8         72.8         87.6         77.0           CosAgentless         57.2         68.2         55.3         63.1         75.8         90.2         79.4

Table 4: Localization performance on the **SWE-bench Verified** test set, consolidating Component-level and File-level results. All scores are reported in percentages (%).

To better understand how augmented data improves final resolution rates, we analyze the intermediate step of *localization* on SWE-bench Verified. Localization accuracy determines whether the correct files or components are identified, and thus directly affects downstream success.

Effect of multilingual extension. Introducing multilingual data consistently improves localization: component Hit@1 increases from 55.0% to 57.2%, and file Hit@1 from 72.8% to 75.8% (Table 4). This shows that data from other programming languages provides transferable structural and semantic signals that strengthen the model's ability to locate relevant code regions. As a result, localization accuracy and overall resolution both improve, even when the target task is restricted to Python repositories.

Effect of completion extension. In contrast, extending with completion data leads to a different effect: component Hit@1 decreases from 55.0% to 52.8%, while file-level metrics remain nearly unchanged. This is expected, since completion data focuses on the code editing objective rather than localization. Formally, such data trains the model according to Eq. 4, which optimizes the synthesis of code edits  $\delta$  given a fixed ground-truth location  $\ell^*$ , without providing supervision for selecting  $\ell$ . Despite the degradation in localization, the strengthened editing capability yields higher overall resolution rates, demonstrating that the positive impact on code edit outweighs the negative impact on localization.

**Implication.** These results reveal distinct mechanisms: multilingual extension simultaneously improves localization and editing by enabling transfer across languages, whereas completion extension primarily enhances code editing ability, especially in scenarios that require substantial additions of new content, such as feature implementation.

#### 5.3 SCALING ANALYSIS

To validate the effectiveness and robustness of our data construction approach, we conduct scaling experiments across both model sizes and training data volumes. We evaluate our method on <code>Qwen2.5-Coder-7B-Instruct</code> and systematically sample 20%, 40%, and 100% of our augmented training data to examine scaling behaviors. The results are presented in Figure 3.

Consistent scaling patterns across model sizes. Our experiments reveal that scaling laws hold consistently across different model capacities. The 32B model substantially outperforms its 7B counterpart across all configurations, achieving an average improvement of +12.8%. Notably, the relative performance trends between different extension strategies remain remarkably stable: multilingual extension consistently provides the largest gains (+5.4% for 7B, +1.4% for 32B over baseline), while completion extension shows modest improvements. This consistency across model scales validates the robustness of our data construction pipeline and suggests that both multilingual and completion extensions scale predictably with model parameters.

Extension strategies enhance data scaling. Figure 3(a) shows that applying our extension strategies yields clear logarithmic scaling patterns as the training set size increases. The largest gains (+9.4%) appear when moving from 0% to 20% of the data, with diminishing returns thereafter. Compared with the SWE-Ext-Baseline, which simply discards non-Python samples or PRs with omitted description and achieves only 18.4% on SWE-bench Verified, our approach lifts the entire scaling curve. By converting previously unusable data into effective training signals through multilingual and completion extensions, we substantially improve both data efficiency and the scaling potential of repository-level code generation models.

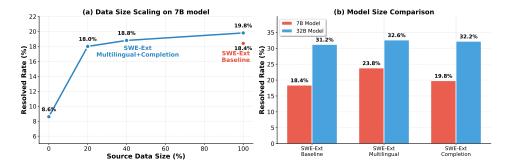


Figure 3: Scaling law analysis on SWE-bench. (a) Performance scaling with training data size follows an approximate logarithmic curve. (b) Model size comparison across three configurations demonstrates consistent improvements from 7B to 32B models.

# 6 RELATED WORK

# 6.1 Repository-level Coding Task

Early benchmarks for code generation focused on function-level tasks, such as HumanEval (Chen et al., 2021) and MBPP (Austin et al., 2021) for simple algorithms. These evaluations emphasized isolated synthesis but overlooked repository contexts. Subsequent benchmarks introduced class- and multi-file challenges, including ClassEval (Du et al., 2023), BigCodeBench (Zhuo et al., 2024), and LiveCodeBench (Jain et al., 2024). Evaluation of real problem in cpde repositories is advanced with SWE-bench (Jimenez et al., 2024) for GitHub issue resolution, alongside variants like DevEval (Li et al., 2024b) and EvoCodeBench (Li et al., 2024a) that align with real repositories. FEA-Bench (Li et al., 2025) targets feature additions via pull requests. Our work aims to enhance the capabilities

of large language models in repository-level coding tasks through post-training with augmented real-world data.

#### 6.2 SOFTWARE ENGINEERING AGENTS

Although basic code large language models (Hui et al., 2024) can achieve good performance on many code benchmarks through direct generation, for repository-level tasks involving numerous files and complex edits, the common practice is to incorporate agent frameworks (Yao et al., 2023). Agents for repository-level tasks often employ iterative processes for issue resolution. SWE-agent (Yang et al., 2024) uses agent-computer interfaces for navigation and editing, while AutoCodeRover (Zhang et al., 2024) integrates fault localization for repairs. CodePlan (Bairi et al., 2024) focuses on planning modifications. Localization-specific agents include CoSIL (Jiang et al., 2025), and LocAgent (Chen et al., 2025) for multi-hop reasoning via heterogeneous graphs. Agentless (Xia et al., 2024) simplifies to a three-phase process without complex tooling. Multi-agent platforms like OpenHands (Wang et al., 2024) facilitate collaboration. Our data construction method draws from the commonalities of these agents, employing a pipeline to build data for four subtasks, thereby enhancing model capabilities in repository-level coding.

# 

# 6.3 TRAINING DATA FOR REPOSITORY-LEVEL CODING

Traditional instruction tuning methods enhance code models' adaptability for general code generation, such as WizardCoder (Luo et al., 2023), WaveCoder (Yu et al., 2024), and Magicoder (Wei et al., 2024). However, constructing training data for repository-level tasks differs significantly, as it requires selectively building input-output and reasoning processes to handle complex interactions.

Mainstream approaches involve building executable environments to collect agent trajectories for rejection sampling and supervised fine-tuning. R2E (Jain et al., 2025) and SWE-Gym (Pan et al., 2024) create runtime environments to gather verified trajectories from limited task instances. SWE-Smith (Yang et al., 2025b) extends synthetic data generation to produce more verifiable tasks under constrained environments. Other methods focus on reinforcement learning-compatible data. SWE-RL (Wei et al., 2025) refines reasoning using software evolution data, while SoRFT (Ma et al., 2025) employs subtask-oriented fine-tuning with rejection sampling and PPO.

Due to the diversity requirements of code repositories, some works augment real-world data or synthetic data without constructing environments. SWE-Fixer (Xie et al., 2025) gathers data with chain-of-thought, and MCTS-Refine (Wang et al., 2025) builds reasoning chains via MCTS to form instruction data. In relation to these efforts, while training data construction for repository-level tasks predominantly focuses on Python, our SWE-Ext pipeline extends the data scope by systematically gathering multi-language data of different subtasks, demonstrating that even out-of-distribution data like code completion can further enhance model performance in repository-level scenarios.

# 

#### 7 CONCLUSION

In this work, we introduced SWE-Ext, a scalable pipeline for extending augmented data in repository-level coding tasks. By broadening coverage to multilingual PRs across ten languages and incorporating code completion as an auxiliary task, we expand the scope of GitHub-derived training data beyond traditional Python-centric datasets. Our experiments demonstrate consistent improvements on benchmarks like SWE-bench and FEA-Bench, with multilingual extensions enhancing overall abilities and completion data strengthening code editing. Moreover, we reveal robust cross-task and cross-language transfer, where diverse data sources benefit monolingual performance, underscoring the potential of data extension in different dimensions for model generalization.

**Limitations and future directions.** SWE-Ext primarily focuses on augmented data without integrating execution-based verification, which could further refine supervision quality. Additionally, due to the constraints on computational and API resources, we cannot use better expert models for data construction and carry out more experiments for more data combinations and on more base models. Future work could explore verifiable pipelines or extend to additional languages and task

types, such as refactoring or generation of unit tests.

# **ETHICS STATEMENT**

We are committed to the responsible and ethical development of artificial intelligence. Our source data is derived exclusively from publicly available, open-source code on GitHub, which inherently poses minimal privacy risks.

During the data augmentation process, we acknowledge the potential for the expert model to generate content that could be considered harmful. However, as these outputs are constrained to code-related reasoning and dialogues, the risk is considered manageable.

As is the case with any model trained on software engineering data, the models trained on our dataset may produce outputs containing security vulnerabilities or even malicious code. Executing these outputs carries a risk of compromising a system. Therefore, for all evaluation and deployment purposes, we strongly recommend that users operate within a securely isolated environment, such as a Docker container, to mitigate any potential harm.

#### REPRODUCIBILITY STATEMENT

We are committed to ensuring the reproducibility of our results. To facilitate the replication of our work, we provide a comprehensive description of our methodology and experimental setup. In Appendix D, we detail the prompts and the full pipeline used for data construction. Furthermore, Appendix B provides all specified hyperparameters and configurations of our training, which were conducted using a standard open-source training framework. For further reference and insight into the data format, representative data examples are included in the supplementary materials. This detailed documentation provides all the necessary components for researchers to reproduce our findings.

#### REFERENCES

- Anthropic. Introducing Claude 3.5 Sonnet, 2024. URL https://www.anthropic.com/news/claude-3-5-sonnet.
- Anthropic. Introducing Claude 3.7 Sonnet, 2025a. URL https://www.anthropic.com/news/claude-3-7-sonnet.
- Anthropic. Introducing Claude 4, 2025b. URL https://www.anthropic.com/news/claude-4.
- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language models. *ArXiv preprint*, abs/2108.07732, 2021. URL https://arxiv.org/abs/2108.07732.
- Ramakrishna Bairi, Atharv Sonwane, Aditya Kanade, Arun Iyer, Suresh Parthasarathy, Sriram Rajamani, B Ashok, and Shashank Shet. Codeplan: Repository-level coding using llms and planning. *Proceedings of the ACM on Software Engineering*, 1(FSE):675–698, 2024.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. *ArXiv preprint*, abs/2107.03374, 2021. URL https://arxiv.org/abs/2107.03374.
- Zhaoling Chen, Xiangru Tang, Gangda Deng, Fang Wu, Jialong Wu, Zhiwei Jiang, Viktor Prasanna, Arman Cohan, and Xingyao Wang. Locagent: Graph-guided llm agents for code localization. *arXiv* preprint arXiv:2503.09089, 2025.
- Xueying Du, Mingwei Liu, Kaixin Wang, Hanlin Wang, Junwei Liu, Yixuan Chen, Jiayi Feng, Chaofeng Sha, Xin Peng, and Yiling Lou. Classeval: A manually-crafted benchmark for evaluating llms on class-level code generation. *ArXiv preprint*, abs/2308.01861, 2023. URL https://arxiv.org/abs/2308.01861.

- Binyuan Hui, Jian Yang, Zeyu Cui, Jiaxi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Keming Lu, et al. Qwen2. 5-coder technical report. *ArXiv preprint*, abs/2409.12186, 2024. URL https://arxiv.org/abs/2409.12186.
  - Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *ArXiv preprint*, abs/2410.21276, 2024. URL https://arxiv.org/abs/2410.21276.
    - Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. Openai of system card. *ArXiv* preprint, abs/2412.16720, 2024. URL https://arxiv.org/abs/2412.16720.
    - Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free evaluation of large language models for code. *ArXiv preprint*, abs/2403.07974, 2024. URL https://arxiv.org/abs/2403.07974.
    - Naman Jain, Jaskirat Singh, Manish Shetty, Liang Zheng, Koushik Sen, and Ion Stoica. R2e-gym: Procedural environments and hybrid verifiers for scaling open-weights swe agents. *arXiv* preprint *arXiv*:2504.07164, 2025.
    - Zhonghao Jiang, Xiaoxue Ren, Meng Yan, Wei Jiang, Yong Li, and Zhongxin Liu. Cosil: Software issue localization via llm-driven code repository graph searching. *arXiv* preprint *arXiv*:2503.22424, 2025.
    - Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R Narasimhan. SWE-bench: Can language models resolve real-world github issues? In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=VTF8yNQM66.
    - Jia Li, Ge Li, Xuanming Zhang, Yihong Dong, and Zhi Jin. Evocodebench: An evolving code generation benchmark aligned with real-world code repositories. *ArXiv preprint*, abs/2404.00599, 2024a. URL https://arxiv.org/abs/2404.00599.
    - Jia Li, Ge Li, Yunfei Zhao, Yongmin Li, Huanyu Liu, Hao Zhu, Lecheng Wang, Kaibo Liu, Zheng Fang, Lanshen Wang, et al. Deveval: A manually-annotated code generation benchmark aligned with real-world code repositories. *ArXiv preprint*, abs/2405.19856, 2024b. URL https://arxiv.org/abs/2405.19856.
    - Wei Li, Xin Zhang, Zhongxin Guo, Shaoguang Mao, Wen Luo, Guangyue Peng, Yangyu Huang, Houfeng Wang, and Scarlett Li. Fea-bench: A benchmark for evaluating repository-level code generation for feature implementation. *arXiv preprint arXiv:2503.06680*, 2025.
    - Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *ArXiv preprint*, abs/2412.19437, 2024. URL https://arxiv.org/abs/2412.19437.
    - Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xiubo Geng, Wenxiang Hu, Chongyang Tao, Jing Ma, Qingwei Lin, and Daxin Jiang. Wizardcoder: Empowering code large language models with evol-instruct. *ArXiv preprint*, abs/2306.08568, 2023. URL https://arxiv.org/abs/2306.08568.
    - Yingwei Ma, Rongyu Cao, Yongchang Cao, Yue Zhang, Jue Chen, Yibo Liu, Yuchen Liu, Binhua Li, Fei Huang, and Yongbin Li. Lingma swe-gpt: An open development-process-centric language model for automated software improvement. *arXiv preprint arXiv:2411.00622*, 2024.
  - Zexiong Ma, Chao Peng, Pengfei Gao, Xiangxin Meng, Yanzhen Zou, and Bing Xie. Sorft: Issue resolving with subtask-oriented reinforced fine-tuning. *arXiv preprint arXiv:2502.20127*, 2025.
  - Christopher D Manning. *Introduction to information retrieval*. Syngress Publishing,, 2008.
  - Jiayi Pan, Xingyao Wang, Graham Neubig, Navdeep Jaitly, Heng Ji, Alane Suhr, and Yizhe Zhang. Training software engineering agents and verifiers with swe-gym. *arXiv* preprint *arXiv*:2412.21139, 2024.

- Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. Hybridflow: A flexible and efficient rlhf framework. *arXiv preprint arXiv:* 2409.19256, 2024.
- Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen, Jiahao Chen, Ningxin Chen, Ruijue Chen, Yanru Chen, Yuankun Chen, Yutian Chen, et al. Kimi k2: Open agentic intelligence. *arXiv* preprint arXiv:2507.20534, 2025.
- Xingyao Wang, Boxuan Li, Yufan Song, Frank F Xu, Xiangru Tang, Mingchen Zhuge, Jiayi Pan, Yueqi Song, Bowen Li, Jaskirat Singh, et al. Openhands: An open platform for ai software developers as generalist agents. *arXiv preprint arXiv:2407.16741*, 2024.
- Yibo Wang, Zhihao Peng, Ying Wang, Zhao Wei, Hai Yu, and Zhiliang Zhu. Mcts-refined cot: High-quality fine-tuning data for llm-based repository issue resolution. *arXiv* preprint *arXiv*:2506.12728, 2025.
- Yuxiang Wei, Zhe Wang, Jiawei Liu, Yifeng Ding, and Lingming Zhang. Magicoder: Empowering code generation with oss-instruct. In *Forty-first International Conference on Machine Learning*, 2024.
- Yuxiang Wei, Olivier Duchenne, Jade Copet, Quentin Carbonneaux, Lingming Zhang, Daniel Fried, Gabriel Synnaeve, Rishabh Singh, and Sida I Wang. Swe-rl: Advancing llm reasoning via reinforcement learning on open software evolution. *arXiv* preprint arXiv:2502.18449, 2025.
- Chunqiu Steven Xia, Yinlin Deng, Soren Dunn, and Lingming Zhang. Agentless: Demystifying Ilm-based software engineering agents. *ArXiv preprint*, abs/2407.01489, 2024. URL https://arxiv.org/abs/2407.01489.
- Chengxing Xie, Bowen Li, Chang Gao, He Du, Wai Lam, Difan Zou, and Kai Chen. Swe-fixer: Training open-source llms for effective and efficient github issue resolution. *arXiv* preprint arXiv:2501.05040, 2025.
- Frank F. Xu, Uri Alon, Graham Neubig, and Vincent Josua Hellendoorn. A systematic evaluation of large language models of code. In *Proceedings of the 6th ACM SIGPLAN International Symposium on Machine Programming*, MAPS 2022, pp. 1–10, New York, NY, USA, 2022. Association for Computing Machinery. ISBN 9781450392730. doi: 10.1145/3520312.3534862. URL https://doi.org/10.1145/3520312.3534862.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025a.
- 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. *ArXiv preprint*, abs/2405.15793, 2024. URL https://arxiv.org/abs/2405.15793.
- John Yang, Kilian Lieret, Carlos E Jimenez, Alexander Wettig, Kabir Khandpur, Yanzhe Zhang, Binyuan Hui, Ofir Press, Ludwig Schmidt, and Diyi Yang. Swe-smith: Scaling data for software engineering agents. *arXiv preprint arXiv:2504.21798*, 2025b.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. In *International Conference on Learning Representations (ICLR)*, 2023.
- Zhaojian Yu, Xin Zhang, Ning Shang, Yangyu Huang, Can Xu, Yishujie Zhao, Wenxiang Hu, and Qiufeng Yin. WaveCoder: Widespread and versatile enhancement for code large language models by instruction tuning. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 5140–5153, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.280. URL https://aclanthology.org/2024.acl-long.280/.

Yuntong Zhang, Haifeng Ruan, Zhiyu Fan, and Abhik Roychoudhury. Autocoderover: Autonomous program improvement. In *Proceedings of the 33rd ACM SIGSOFT International Symposium on Software Testing and Analysis*, ISSTA 2024, pp. 1592–1604, New York, NY, USA, 2024. Association for Computing Machinery. ISBN 9798400706127. doi: 10.1145/3650212.3680384. URL https://doi.org/10.1145/3650212.3680384.

Yanli Zhao, Andrew Gu, Rohan Varma, Liang Luo, Chien-Chin Huang, Min Xu, Less Wright, Hamid Shojanazeri, Myle Ott, Sam Shleifer, et al. Pytorch fsdp: experiences on scaling fully sharded data parallel. *arXiv preprint arXiv:2304.11277*, 2023.

Terry Yue Zhuo, Minh Chien Vu, Jenny Chim, Han Hu, Wenhao Yu, Ratnadira Widyasari, Imam Nur Bani Yusuf, Haolan Zhan, Junda He, Indraneil Paul, et al. Bigcodebench: Benchmarking code generation with diverse function calls and complex instructions. *ArXiv preprint*, abs/2406.15877, 2024. URL https://arxiv.org/abs/2406.15877.

#### A USE OF LLMS

Large Language Models (LLMs) were applied in this research in two specific ways. First, LLMs assisted in the development of the codebase, primarily through code completion functionalities. Second, LLMs were employed for language refinement to improve the clarity and readability of the manuscript. All research ideas, experimental designs, and methodological frameworks were independently conceived and implemented by the authors. The authors take full responsibility for the content and conclusions presented in this work.

# B TRAINING SETTINGS

From our preliminary experiments, we found that using <code>Qwen2.5-Coder</code> as the base model yields strong performance, consistent with prior works on constructing training datasets (Yang et al., 2025b; Jain et al., 2025; Wang et al., 2025). We adopt <code>Qwen2.5-Coder-32B-Instruct</code> and <code>Qwen2.5-Coder-7B-Instruct</code> as the backbone models for fine-tuning. For the training framework, we use the multi-turn dialogue SFT framework provided by <code>verl</code> (Sheng et al., 2024), which is implemented on top of FSDP (Zhao et al., 2023).

For the 32B model, we fine-tune on two GPU workstations, each equipped with 8 NVIDIA H200 GPUs. Training is performed in float 32 precision with CPU offloading enabled. We set the maximum sequence length to 32,768. The training configuration includes a sequence parallel size of 2, a learning rate of  $5 \times 10^{-6}$ , a training batch size of 128, and training for 3 epochs. The total training time is approximately 2, 7, and 3 days on the SWE-EXT-BASELINE, SWE-EXT-MULTILINGUAL, and SWE-EXT-COMPLETION datasets, respectively.

For the 7B model, we fine-tune on a single workstation with 8 NVIDIA A100 GPUs. The setup remains the same as the 32B model, except that the sequence parallel size is increased to 4. The training time is similar to that of the 32B model.

#### C EVALUATION SETTINGS

All execution-based evaluations were conducted on an cloud computing instance equipped with a 32-core AMD EPYC 7763 processor @ 2.45GHz and 256 GB RAM, running Ubuntu 22.04.5 LTS with Linux kernel 6.8.0-1027-azure. It is important to note that in our evaluation environment, we encountered compatibility issues with certain test cases where 5 out of 500 gold patches in SWE-bench Verified and 9 out of 200 gold patches in FEA-Bench Lite failed to execute successfully due to environment-specific dependencies or configuration conflicts. As a result, the reported performance metrics may represent a slight underestimate of the true capabilities of our approach.

# D PROMPTS OF COSAGENTLESS AND DETAILS OF AUGMENTED DATA

In this section, we detail the prompts used and the specifics of how they construct our four task datasets. All of our data is built upon real-world GitHub pull requests and augmented based on

their extracted ground truth results, which is why we refer to it as "augmented data". This data augmentation process consists of transforming the raw pull requests into a dialogue-based format, following the steps of our CosAgentless pipeline.

#### D.1 FILE LOCALIZATION

**Phase 1: File Localization.** The first phase of the task is to select a top-5 list of relevant files based on the problem statement (problem\_statement) and the provided file tree (structure). The prompt used for this task is as follows:

```
The problem description is as follows:
### GitHub Problem Description ###
{problem_statement}
###
### Candidate Files ###
{structure}
###
Let's locate the relevant file step by step using reasoning.
In order to locate accurately, you can pre-select {pre_select_num} files,
and finally confirm {top_n} file names.
Based on the available information, confirm and provide complete name of
the top-5 most likely relevant files that need to be edited for the
problem.
You should output your reasoning process first.
Since your final answer will be processed automatically, please give your
final answer of relevant files in the format as follows.
The returned files should be separated by new lines ordered by most to
least important and wrapped with '''.
file1.py
file2.py
file3.py
file4.py
file5.py
Replace the 'file1.py' with the actual file path.
For example,
sklearn/linear_model/__init__.py
sklearn/base.py
```

**Phase 2: Dependency Analysis and Reflection.** The second phase analyzes the dependencies of the files selected in the previous stage based on their import relationships import\_content. These relationships are derived from code segments identified by regular expressions. The model is then prompted to reflect its selection and choose a new list of up to 10 relevant files.

```
Please look through the following problem description and repository structure and provide a list of files that one would need to edit to solve the software development problem.

I have already find 5 relevant files. Accrording to the import relations, construct the call graph first.

### Problem Description ###
{problem_statement}
```

```
756
       ###
757
758
       ### Repository Structure ###
759
       {structure}
760
       ###
761
762
       ### Files To Be Explored ###
763
       {pre_files}
764
       ###
765
766
       ### Import Relations ###
767
       {import_content}
768
       ###
769
770
       Based on the import relationships, please analyze which files in the
771
       repository depend on which other files within the same repository. Ignore
772
       any libraries or modules that are imported from outside the current
773
      repository. Present the results in the following format:
774
       file1.py -> file6.py, file7.py
775
       file2.py ->
776
       file3.py -> xxx/file8.py
777
       file4.py -> file2.py, file9.py
778
       file5.py -> xxx/file10.py, xxx/file11.py
779
      Note: Solving the problem not only requires determining where to modify
780
      the code, but also identifying which other code to refer to in order to
781
      understand and invoke.
782
783
      Based on the files listed above and the import relations, reconfirm and
784
      provide the complete names of the top 10 most likely relevant files,
       considering both where changes need to be made and which files are
785
      important to refer to.
786
787
      Please think step by step and give you reasoning process first. Finally,
788
      provide full path and return top 10 files.
789
      The final returned files should be separated by new lines ordered by most
790
       to least important and wrapped with '''
791
      For example:
792
793
       file1.py
       file2.py
794
      file3.py
795
      file4.py
796
       file5.py
797
       file6.py
798
      file7.py
      file8.py
799
       file9.py
800
       file10.py
801
802
      Note: file1.py indicates the top-1 file, file2.py indicates the top-2
803
       file, and so on. Do not include test files.
804
```

**Phase 3: Refined File Localization.** This phase leverages the results of parsing code repository. The model is provided with the structural skeletons (file\_internal\_structure, including class and function definitions) of the files identified in the previous stage. Using the skeletons of the files, the model makes a final selection, outputting a refined list of at most 5 files as the final localization result.

807

808

```
810
811
      Please look through the following problem description and repository
       structure and provide a list of files that one would need to edit to
812
       solve the software development problem.
813
       I have already find 10 relevent files according to the file structure and
814
       dependency.
815
      I will further give the internal structures of these 10 files.
816
      Please rank them again and reflect the result according to the internal
      structures and dependency.
817
818
       ### Problem Description ###
819
      {problem_statement}
820
821
       ###
822
       ### Import Relations ###
823
       {import_content}
824
825
       ###
826
827
       ### Files To Be Ranked ###
828
       {file_internal_structure}
829
830
       ###
831
      Please think step by step and give you reasoning process first. Finally,
832
      provide full path and return top 5 files.
833
834
      The returned files should be separated by new lines ordered by most to
       least important and wrapped with '''
836
      For example:
837
      file1.py
838
      file2.py
839
       file3.py
840
      file4.py
      file5.py
841
842
      Note: file1.py indicates the top-1 file, file2.py indicates the top-2
843
       file, and so on. Do not include test files.
844
```

The CosAgentless inference pipeline is structured as a sequential execution of the three phases detailed above. However, during the construction of our ground-truth-based augmented dataset, we employ a specific self-correction process: if the expert model fails to recall all the files in the ground truth code edit in its initial attempt, we provide the ground truth files as an explicit prompt to guide a re-prediction.

#### D.2 COMPONENT LOCALIZATION

845 846

847

848

849

850851852

853

854

855

856 857

858

859 860

861

862

The prompt for our component localization process is an enhanced version of the one presented in Jiang et al. (2025), specifically adapted to generalize to a broader set of programming languages and repository-level coding tasks.

**System Prompt.** Within the system prompt, we define a set of tool-use capabilities that enable the model to inspect the code of a specified component.

```
You will be presented with a repository-level coding problem with repository file structure to access the source code of the software. Since the modification is based on the code repository, the modified locations may include files, classes, and functions, and the modifications may be in the form of addition, deletion, or update.
```

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882 883 884

885

886

887

888

889 890 891

892

893

894

895

897

898

899

900

901

902

903

904

905906907

908

909

910

911

912 913

914

915

916

917

```
Your task is to locate the top-5 most likely edit locations based on the
problem description and the information you retrieve using given
functions.
Function calls you can use are as follows:
* get_code_of_class('file_name', 'class_name') -> Get the code of a
specified class in the given file and python project. 'file_name' -> The
name of the file. 'class_name' -> The name of the class. \star
* get_code_of_class_function('file_name', 'class_name', 'func_name') ->
Get the code of a specified function in the given class, file, and python
project. 'file_name' -> The name of the file. 'class_name' -> The name
of the class. 'func_name' -> The name of the function. *
* get_code_of_file_function('file_name', 'func_name') -> Get the code of
a specified function in the given file and python project. 'file_name' ->
The name of the file. 'func_name' -> The name of the function. *
* get_toplevel_code('file_name') -> Get all the code in a given file that
is not part of a class or function definition. This is useful for
viewing imports, global variables, constants, and any top-level script
logic. 'file_name' -> The name of the file. *
* exit() -> Exit function calling to give your final answer when you are
confident of the answer. *
You have {max_try} chances to call function.
```

**Iterative Localization Initialization.** The iterative localization process begins with an initialization step where the model is prompted to identify a single, relevant component. This component serves as the initial point of inspection for the codebase.

```
### Problem Description ###
{problem_statement}
Let's locate the relevant elements (function/class) step by step using
reasoning and function calls.
I have pre-identified top-5 relevant files. There stuctures are as
follows:
{bug file list}
The formal parameter 'file_name' takes the value in "file:"
The formal parameter 'class_name' takes the value in "class:"
The formal parameter 'func_name' takes the value in "static functions:"
and "class functions: "
Avoid making multiple identical calls to save overhead.
You must strictly follow the structure I give to call different tools.
For static functions, you can use 'get_code_of_file_function', and for
class functions, you can use 'get_code_of_class_function'.
In order to locate accurately, you can pre-select {pre_select_num}
locations, then check them through function calls, and finally confirm {
top_n} file names.
Don't make the first function call in this message.
```

**Iterative Component Inspection.** This process involves the iterative inspection of components. At each step, the model leverages the full dialogue history to inform its decision. By making sequential function calls, it dynamically identifies the next component or code segment to examine, progressively refining its understanding of the problem and the codebase. file\_internal\_structure are the skeletons of found files in the file localization stage.

```
Now call a function in this format 'FunctionName(Argument)' in a single line without any other word or signal (such as ''').

Don't call the same function you've previous called, because this may waste your context length.

{file_internal_structure}
```

919

920

921

922

949

950

951

952

953

954

For each component component\_retrieved viewed during an iteration, the model is prompted to assess its relevance to the given problem. If the component is relevant to the problem description, the model is then directed to consider its internal call and dependency relationships to select the next component to be viewed. Conversely, if a component is determined to be irrelevant, no further analysis is conducted on it.

```
923
       You will be presented with a repository-level coding problem with
924
      repository file structure to access the source code of the software.
925
926
       Your task is to locate the top-5 most likely edit locations based on the
927
      problem description.
928
       ### Problem Description ###
929
       {problem_statement}
930
931
       ###
932
      Here is a result of a function/class code retrived by '{content}'.
933
      Please check if the code is related to the problem and if the code should
934
       be added into context.
935
       <code>
936
       {component_retrieved}
937
       </code>
      Return True if the code is related to the problem and should be added
938
      into context, otherwise return False.
939
      Since your answer will be processed automatically, please give your
940
      answer in the format as follows.
941
      The returned content should be wrapped with '''.
942
      True
943
       . . .
944
       or
945
946
      False
947
       . . .
948
```

**Final Component Localization Output.** The final component localization output is generated upon the termination of the iterative process. This occurs when the model either invokes the <code>exit()</code> function or reaches the maximum iteration limit. The model is then instructed to summarize its full inspection history and produce a final set of predicted components relevant to the <code>problem\_statement</code>.

```
955
      {file_internal_structure}
956
      Based on the available information, reconfirm and provide complete names
      of the most likely edit locations (10 locations at most).
957
      Before make the final decision, please check whether the function name is
958
       correct or not, for static functions, don't add class name.
959
      {bug_file_list}
960
      Please provide the complete set of locations as either a class name, a
961
      function name, or a file name.
962
      The returned files should be separated by new lines ordered by most to
963
      least important and wrapped with '''
964
      Since your answer will be processed automatically, please give your
965
      answer in the exapmle format as follows.
966
      top1_file_fullpath.py
967
      function: Class1.Function1
968
969
      top2_file_fullpath.py
970
      function: Function2
971
      top3_file_fullpath.py
```

```
972
      class: Class3
973
974
       top4_file_fullpath.py
975
      function: Class4.Function4
976
      top5_file_fullpath.py
977
      function: Function5
978
979
      top6_file_fullpath.py
980
      global
981
       top7_file_fullpath.py
982
      function: Class7.Function7
983
984
      Replace the 'Top_file_fullpath.py' with the actual file path, the 'Class'
       with the actual class name and the 'Function' with the actual function
985
      name. 'global' means the code is not in a class or function.
986
      For example,
987
988
       sklearn/linear_model/__init__.py
989
       function: LinearRegression.fit
990
991
```

During the construction of our augmented data from GitHub pull requests, we incorporate a feedback mechanism. If the final localization output does not fully recall all of the locations specified by the ground truth code edit, we explicitly introduce the omitted locations back into the iterative inspection process. This allows the model to correct its course and regenerate a complete and accurate final localization result.

#### D.3 CODE EDITING

992

993

994

995

996

997 998

999

1000

1001

1002

**Final Code Editing.** Leveraging the code at the locations identified by our component localization process, we construct a context block denoted as top\_n\_content. This contextual information is then provided to the large language model to facilitate the code editing task.

```
1003
      We are currently solving the following task within our repository. Here
1004
      is the task description.
1005
1006
       ### Task Description ###
1007
       {problem_statement}
1008
1009
       ###
1010
      Below are some code segments, each from a relevant file. One or more of
1011
      these files may need to be edited to solve the task.
1012
1013
1014
       --- BEGIN FILE ---
       . . .
1015
       {top_n_content}
1016
1017
       --- END FILE ---
1018
      Please first localize the positions to edit based on the task statement,
1019
      and then output the files that need to be deleted, modified or added.
1020
      '- file' means deleting the file;
1021
      '\star file' means modificating the file;
1022
       '+ file' means adding the file.
1023
      The file should list like below:
1024
      - file1.py
1025
      * file2.py
```

```
1026
      * file3.py
1027
      + file4.py
1028
1029
      To solve the task, you should then generate *SEARCH/REPLACE* edits.
1030
1031
      Every *SEARCH/REPLACE* edit must use this format:
1032
      1. The file path
1033
      2. The start of search block: <<<< SEARCH
1034
      3. A contiguous chunk of lines to search for in the existing source code
      4. The dividing line: =====
1035
      5. The lines to replace into the source code
1036
      6. The end of the replace block: >>>>> REPLACE
1037
1038
      Here is an example:
1039
       '''python
1040
      ### mathweb/flask/app.py
1041
      <<<<< SEARCH
1042
      from flask import Flask
1043
      import math
1044
      from flask import Flask
1045
      >>>>> REPLACE
1046
1047
1048
      Please note that the *SEARCH/REPLACE* edit REQUIRES PROPER INDENTATION.
1049
      If you would like to add the line '
                                                   print(x)', you must fully
      write that out, with all those spaces before the code!
1050
      Wrap the *SEARCH/REPLACE* edit in blocks '''python...'''.
1051
      When multiple edits should be done, please output *SEARCH/REPLACE* edit
1052
      one by one and give your reasoning process before each *SEARCH/REPLACE*
1053
      block.
1054
```

**Code Editing Output Format.** The model's output is structured according to a specific format. The code edits are presented as **SEARCH-REPLACE** pairs, where the old code is provided alongside the corresponding new code. In addition, the model is required to generate a detailed editing plan and a clear rationale for the proposed changes.

```
The plan of solving this software task:
{plan}
The files that should be edited include:
{files}

Here are my edits for code.
{search_replaces}
```

For the construction of our code editing dataset, we first transform the ground truth code edits into a structured format of **SEARCH-REPLACE** pairs. Subsequently, we utilize the expert model to generat the reasoning for these edits, populating the designated plan field of the response.

#### D.4 CODE COMPLETION

1055

1056

1057

1058

1059

1067

1068

1069

1070 1071

1072

1074

1075

1077

1078

1079

The code completion data is constructed from pull requests that lack a complete problem description (<50 words). When generating these samples, we prioritize targeting a complete, newly added function within the patch. If no such function exists, the target becomes a contiguous block of new code. The code in the context (top\_n\_content) is then composed of the target file along with other relevant components identified via our CosAgentless localization pipeline. The line requiring completion is replaced with the token [TODO], thereby converting the code completion task as an editing task centered on the [TODO] token. Consequently, the prompt and the subsequent augmented data generation steps align with those used for our code editing task.

# E DETAILS OF FILTERING

#### E.1 FILTERING PULL REQUESTS

To ensure the collection of high-quality data during the pull request (PR) scraping phase, we applied the following exclusion criteria:

- PRs originating from repositories with fewer than 100 stars or a total of fewer than 100 PRs, as this often indicates low repository quality or inconsistent contribution patterns.
- PRs that were not merged into the main branch, indicating that the proposed code changes were not accepted.
- PRs whose patches failed to apply cleanly to the codebase.
- PRs where tree-sitter encountered parsing errors in the changed files.
- PRs whose changed files exclusively consisted of non-code file extensions, including: [".json", ".png", "csv", ".txt", ".md", ".jpg", ".jpeg", ".pkl", ".yml", ".yaml", ".toml"].

# E.2 FILTERING TRAINING DATA

After the construction of the SWE-Ext training data, further filtering is required to ensure data quality. First, we remove data samples that contain expert model API call errors, which are commonly caused by exceeding the context window length. Second, we remove samples with obvious errors for each data type:

- **File Localization data**: We remove samples where the ground truth edited files are not fully contained within the final prediction or are not included in the Phase 2 results.
- Component Localization data: We remove samples where the ground truth edited components are not fully contained in the final output.
- Code Edit and Code Completion data: We remove samples where the "search" code (the code to be changed) is not present in the provided context, as well as those with incomplete search-replace pairs.

To ensure the model is trained on complete conversational data, we also filter out any samples with a length exceeding 32768 tokens. Finally, we remove any dialogue data that corresponds to pull requests included in our evaluation datasets, including SWE-bench Verified and FEA-Bench, to prevent data leakage.

# F DATA EXAMPLES

In this section, we present several short data examples for reference. Each turn of the dialogue is formatted as follows:

# F.1 FILE LOCALIZATION

```
1134
      The problem description is as follows:
1135
1136
       ### GitHub Problem Description ###
1137
       <request>
       Adding fill parameters to pairwise and windowed
1138
       I found these handy for when I needed lookahead/lookbehind, but didn't
1139
       want to add extra processing code before/after the processing loop. Here'
1140
       s a sketch of one way I've used it:
1141
       '''python
1142
       for record, next_record in pairwise_iter(records, fill=None):
1143
        if next_record is not None:
1144
            record.link(next_record)
1145
        record.process()
1146
1147
       Without fill, I have to call 'records[-1].process()' after the loop,
1148
       assuming it's a list and not an iterable.
1149
       </request>
1150
1151
       ###
1152
       ### Candidate Files ###
1153
      boltons/
1154
           __init__.py
1155
           cacheutils.py
1156
           debugutils.py
           deprutils.py
1157
           dictutils.py
1158
           easterutils.py
1159
           ecoutils.py
1160
           excutils.py
1161
           fileutils.py
           formatutils.py
1162
           funcutils.py
1163
           gcutils.py
1164
           ioutils.py
1165
           iterutils.py
1166
           jsonutils.py
           listutils.py
1167
           mathutils.py
1168
           mboxutils.py
1169
           namedutils.py
1170
           pathutils.py
1171
           queueutils.py
           setutils.py
1172
           socketutils.py
1173
           statsutils.py
1174
           strutils.py
1175
           tableutils.py
1176
           tbutils.py
           timeutils.py
1177
           typeutils.py
1178
           urlutils.py
1179
       docs/
1180
           conf.py
1181
      misc/
           bench_omd.py
1182
           linkify_changelog.py
1183
           table_html_app.py
1184
       setup.py
1185
1186
       ###
1187
```

```
1188
1189
      Let's locate the relevant file step by step using reasoning.
1190
      In order to locate accurately, you can pre-select 7 files, and finally
1191
      confirm 5 file names.
1192
1193
      Based on the available information, confirm and provide complete name of
1194
      the top-5 most likely relevant files that need to be edited for the
1195
      problem.
1196
      You should output your reasoning process first.
      Since your final answer will be processed automatically, please give your
1197
       final answer of relevant files in the format as follows.
1198
      The returned files should be separated by new lines ordered by most to
1199
      least important and wrapped with '''.
1200
      file1.py
1201
      file2.py
1202
      file3.py
1203
      file4.py
1204
      file5.py
1205
1206
      Replace the 'file1.py' with the actual file path.
      For example,
1207
1208
      sklearn/linear_model/__init__.py
1209
      sklearn/base.py
1210
1211
1212
      >>>>>> assistant:
1213
      To solve the problem of adding fill parameters to 'pairwise' and '
1214
      windowed' functions, we need to identify where these functions are likely
1215
      implemented. The problem description suggests that these functions are
1216
      related to iterating over collections, which is a common task in utility
      modules dealing with iteration or functional programming.
1217
1218
      Given the directory structure, the most relevant file is likely to be '
1219
      iterutils.py 'because:
1220
      - The name 'iterutils.py' suggests it contains utilities for iteration,
      which aligns with the functions 'pairwise' and 'windowed' that deal with
1221
      iterating over elements in a collection.
1222
1223
      Let's consider the reasoning for selecting the files:
1224
1225
      1. **boltons/iterutils.py**: This file is the most likely candidate as it
       is named to suggest it contains utilities for iteration, which is
1226
      directly related to the 'pairwise' and 'windowed' functions mentioned in
1227
      the problem.
1228
1229
      2. **boltons/funcutils.py**: This file might contain functional utilities
1230
      , which could include higher-order functions or iterators. If 'pairwise'
      and 'windowed' are implemented using functional programming paradigms,
1231
      they might be here.
1232
1233
      3. **boltons/listutils.py**: This file could potentially contain
1234
      utilities for list operations. If 'pairwise' and 'windowed' are
1235
      specifically designed for lists, they might be implemented here.
1236
      4. **boltons/typeutils.py**: This file might contain utilities related to
1237
       type handling, which could include generic functions for handling
1238
      iterables of various types.
1239
1240
      5. **boltons/strutils.py**: Although less likely, if 'pairwise' and '
      windowed' have any specific implementations for strings or string-like
1241
      iterables, they might be here.
```

```
1242
1243
      Now, let's finalize the top 5 files based on the reasoning:
1244
1245
      boltons/iterutils.py
1246
      boltons/funcutils.py
1247
      boltons/listutils.py
1248
      boltons/typeutils.py
1249
      boltons/strutils.py
1250
1251
      >>>>>> user:
1252
1253
      Please look through the following problem description and repository
1254
      structure and provide a list of files that one would need to edit to
      solve the software development problem.
1255
      I have already find 5 relevant files. Accrording to the import relations,
1256
       construct the call graph first.
1257
1258
       ### Problem Description ###
1259
      <request>
1260
      Adding fill parameters to pairwise and windowed
      I found these handy for when I needed lookahead/lookbehind, but didn't
1261
      want to add extra processing code before/after the processing loop. Here'
1262
      s a sketch of one way I've used it:
1263
1264
       '''python
      for record, next_record in pairwise_iter(records, fill=None):
1265
        if next_record is not None:
1266
            record.link(next_record)
1267
        record.process()
1268
1269
      Without fill, I have to call 'records[-1].process()' after the loop,
1270
      assuming it's a list and not an iterable.
1271
      </request>
1272
1273
1274
       ### Repository Structure ###
1275
      boltons/
1276
           \__{\text{init}}_{\text{.py}}
1277
           cacheutils.py
1278
           debugutils.py
1279
           deprutils.py
           dictutils.py
1280
           easterutils.py
1281
           ecoutils.py
1282
           excutils.py
1283
           fileutils.py
1284
           formatutils.py
           funcutils.py
1285
           gcutils.py
1286
           ioutils.py
1287
           iterutils.py
1288
           jsonutils.py
1289
           listutils.py
           mathutils.py
1290
           mboxutils.py
1291
           namedutils.py
1292
           pathutils.py
1293
           queueutils.py
1294
           setutils.py
1295
           socketutils.py
           statsutils.py
```

```
1296
           strutils.py
1297
           tableutils.py
1298
           tbutils.py
1299
           timeutils.pv
           typeutils.py
1300
           urlutils.py
1301
      docs/
1302
           conf.py
1303
      misc/
1304
           bench_omd.py
1305
           linkify_changelog.py
           table_html_app.py
1306
      setup.py
1307
1308
       ###
1309
       ### Files To Be Explored ###
1310
      ['boltons/iterutils.py', 'boltons/funcutils.py', 'boltons/listutils.py', 'boltons/typeutils.py', 'boltons/strutils.py']
1311
1312
1313
       ###
1314
       ### Import Relations ###
1315
      file: boltons/iterutils.py
1316
       ['import os', 'import math', 'import time', 'import codecs', 'import
1317
       random', 'import itertools']
1318
1319
      file: boltons/funcutils.py
        ['from __future__ import print_function', 'import sys', 'import re', '
1320
       import inspect', 'import functools', 'import itertools', 'from types
1321
       import MethodType, FunctionType']
1322
1323
      file: boltons/listutils.py
       ['from __future__ import print_function, division', 'import operator', '
1324
        from math import log as math_log', 'from itertools import chain, islice
1325
1326
1327
      file: boltons/typeutils.py
1328
       ['import sys', 'from collections import deque']
1329
      file: boltons/strutils.py
1330
        ['from __future__ import print_function', 'import re', 'import sys', '
1331
       import uuid', 'import zlib', 'import string', 'import unicodedata',
1332
       import collections', 'from gzip import GzipFile']
1333
       ###
1334
1335
1336
      Based on the import relationships, please analyze which files in the
1337
      repository depend on which other files within the same repository. Ignore
1338
       any libraries or modules that are imported from outside the current
      repository. Present the results in the following format:
1339
1340
       file1.py -> file6.py, file7.py
1341
      file2.py ->
1342
      file3.py -> xxx/file8.py
1343
      file4.py -> file2.py, file9.py
      file5.py -> xxx/file10.py, xxx/file11.py
1344
1345
      Note: Solving the problem not only requires determining where to modify
1346
      the code, but also identifying which other code to refer to in order to
1347
      understand and invoke.
1348
1349
      Based on the files listed above and the import relations, reconfirm and
      provide the complete names of the top 10 most likely relevant files,
```

```
1350
      considering both where changes need to be made and which files are
1351
      important to refer to.
1352
1353
      Please think step by step and give you reasoning process first. Finally,
      provide full path and return top 10 files.
1354
1355
      The final returned files should be separated by new lines ordered by most
1356
       to least important and wrapped with '''
1357
      For example:
1358
      file1.py
1359
      file2.py
1360
      file3.py
1361
      file4.py
1362
      file5.py
      file6.py
1363
      file7.py
1364
      file8.py
1365
      file9.py
1366
      file10.py
1367
      Note: file1.py indicates the top-1 file, file2.py indicates the top-2
1368
      file, and so on. Do not include test files.
1369
1370
1371
      >>>>>> assistant:
1372
      To solve the problem, we need to analyze the import relationships and
      determine which files are relevant for adding the "fill" parameter to the
1373
       'pairwise' and 'windowed' functions. Here's the step-by-step reasoning
1374
      process:
1375
1376
      1. **Identify the Core Functionality**: The problem is about adding a "
1377
      fill" parameter to the 'pairwise' and 'windowed' functions. These
      functions are likely to be part of the 'iterutils.py' file, as they deal
1378
      with iteration utilities.
1379
1380
      2. **Analyze Import Relations**: We need to understand the dependencies
1381
      between the files to see if any other files might be relevant for
1382
      understanding or modifying the 'pairwise' and 'windowed' functions.
1383
      3. **Check Import Statements**:
1384
         - 'boltons/iterutils.py' imports standard libraries like 'os', 'math',
1385
          'time', 'codecs', 'random', and 'itertools'. It does not import any
1386
         other files from the 'boltons' package.
1387
           'boltons/funcutils.py' imports 'itertools' and other standard
         libraries, but does not import 'iterutils'.
1388
         - 'boltons/listutils.py' imports 'itertools' and other standard
1389
         libraries, but does not import 'iterutils'.
1390
         - 'boltons/typeutils.py' and 'boltons/strutils.py' do not import '
1391
         iterutils' or any other relevant files.
1392
      4. **Determine Relevant Files**: Since 'iterutils.py' does not import or
1393
      is not imported by any other files in the list, it is likely self-
1394
      contained for the functionality of 'pairwise' and 'windowed'. However,
1395
      understanding how these functions might interact with other utilities
1396
      could be useful, so we consider the files that import 'itertools' as they
1397
      might have similar iteration-related utilities.
1398
      5. **Prioritize Files**: Based on the above analysis, 'iterutils.py' is
1399
      the most relevant file for direct modification. The other files ('
1400
      funcutils.py', 'listutils.py', 'typeutils.py', 'strutils.py') might
1401
      provide additional context or similar implementations that could be
1402
      useful for reference.
1403
```

```
1404
      6. **Top 10 Files**: Since the problem is specific to iteration utilities
1405
      , the most relevant file is 'iterutils.py'. The other files are included
1406
      for reference and understanding of similar patterns or utilities.
1407
      Based on this reasoning, here are the top 10 files, ordered by relevance:
1408
1409
1410
      boltons/iterutils.pv
1411
      boltons/funcutils.py
1412
      boltons/listutils.py
      boltons/typeutils.py
1413
      boltons/strutils.py
1414
      boltons/cacheutils.py
1415
      boltons/debugutils.py
1416
      boltons/deprutils.py
      boltons/dictutils.py
1417
      boltons/easterutils.py
1418
1419
1420
      Note: The last five files are included to fill the top 10 list, assuming
1421
      they might provide additional context or utility functions that could be
1422
      relevant indirectly.
1423
      >>>>>> 11ser:
1424
1425
      Please look through the following problem description and repository
1426
      structure and provide a list of files that one would need to edit to
1427
      solve the software development problem.
      I have already find 10 relevent files according to the file structure and
1428
      dependency.
1429
      I will further give the internal structures of these 10 files.
1430
      Please rank them again and reflect the result according to the internal
1431
      structures and dependency.
1432
      ### Problem Description ###
1433
      <reguest>
1434
      Adding fill parameters to pairwise and windowed
1435
      I found these handy for when I needed lookahead/lookbehind, but didn't
1436
      want to add extra processing code before/after the processing loop. Here'
      s a sketch of one way I've used it:
1437
1438
       '''python
1439
      for record, next_record in pairwise_iter(records, fill=None):
1440
        if next_record is not None:
1441
           record.link(next_record)
1442
        record.process()
1443
1444
      Without fill, I have to call 'records[-1].process()' after the loop,
1445
      assuming it's a list and not an iterable.
1446
      </request>
1447
      ###
1448
1449
      ### Import Relations ###
1450
      file: boltons/iterutils.py
       ['import os', 'import math', 'import time', 'import codecs', 'import
1451
       random', 'import itertools']
1452
1453
      file: boltons/funcutils.py
1454
       ['from __future__ import print_function', 'import sys', 'import re', '
1455
       import inspect', 'import functools', 'import itertools', 'from types
1456
       import MethodType, FunctionType']
1457
      file: boltons/listutils.py
```

```
1458
         ['from __future__ import print_function, division', 'import operator', '
1459
         from math import log as math_log', 'from itertools import chain, islice
         ']
1460
1461
        file: boltons/typeutils.py
1462
         ['import sys', 'from collections import deque']
1463
1464
        file: boltons/strutils.py
1465
         ['from __future__ import print_function', 'import re', 'import sys', '
1466
         import uuid', 'import zlib', 'import string', 'import unicodedata',
        import collections', 'from gzip import GzipFile']
1467
1468
1469
        ###
1470
1471
        ### Files To Be Ranked ###
1472
        file: boltons/iterutils.py
1473
                 class: ['PathAccessError', 'GUIDerator', 'SequentialGUIDerator']
1474
                 static functions: ['is_iterable', 'is_scalar', 'is_collection',
'split', 'split_iter', 'lstrip', 'lstrip_iter', 'rstrip', '
1475
                 rstrip_iter', 'strip', 'strip_iter', 'chunked', '
1476
                 _validate_positive_int', 'chunked_iter', 'chunk_ranges', '
1477
                 pairwise', 'pairwise_iter', 'windowed', 'windowed_iter', 'xfrange', 'frange', 'backoff', 'backoff_iter', 'bucketize', 'partition', 'unique', 'unique_iter', 'redundant', 'one', 'first', '
1478
1479
1480
                 flatten_iter', 'flatten', 'same', 'default_visit', 'default_enter
                 ', 'default_exit', 'remap', 'get_path', 'research', 'soft_sorted
1481
                 ', 'untyped_sorted']
1482
                 class functions: [
1483
                           PathAccessError: ['__init__', '__repr__', '__str__']
1484
                           GUIDerator: ['__init__', 'reseed', '__iter__', '__next__
1485
                           ', '__next__']
                           SequentialGUIDerator: ['reseed', 'reseed', '__next__']
1486
1487
        file: boltons/funcutils.py
1488
                 class: ['InstancePartial', 'CachedInstancePartial', '
1489
                 FunctionBuilder', 'MissingArgument', 'ExistingArgument']
                 static functions: ['inspect_formatargspec',
1490
                 get_module_callables', 'mro_items', 'dir_dict', 'copy_function',
1491
                 'partial_ordering', 'format_invocation', 'format_exp_repr', 'format_nonexp_repr', 'wraps', 'update_wrapper', '
1492
1493
                 _parse_wraps_expected', '_indent', 'total_ordering', 'noop']
1494
                 class functions: [
1495
                           InstancePartial: ['_partialmethod', '__get__']
                           CachedInstancePartial: ['_partialmethod', '__set_name__',
1496
                            '__get__']
1497
                           FunctionBuilder: ['_argspec_to_dict', '_argspec_to_dict',
1498
                           '__init__', 'get_sig_str', 'get_invocation_str',
get_sig_str', 'get_invocation_str', 'from_func', '
1499
1500
                           get_func', 'get_defaults_dict', 'get_arg_names', 'add_arg
                           ', 'add_arg', 'remove_arg', '_compile']
1501
                           MissingArgument: []
1502
                           ExistingArgument: []
1503
1504
        file: boltons/listutils.py
1505
                 class: ['BarrelList', 'SplayList']
                 static functions: []
1506
                 class functions: [
1507
                           BarrelList: ['__init__', '_cur_size_limit', '
1508
                           _translate_index', '_balance_list', 'insert', 'append', '
1509
                           extend', 'pop', 'iter_slice', 'del_slice', 'from_iterable
', '__iter__', '__reversed__', '__len__', '__contains__',
'__getitem__', '__delitem__', '__setitem__', '
1510
1511
```

```
1512
                            __getslice__', '__setslice__', '__repr__', 'sort', '
1513
                           reverse', 'count', 'index']
1514
                           SplayList: ['shift', 'swap']
1515
                 1
        file: boltons/typeutils.py
1516
                 class: ['classproperty']
1517
                 static functions: ['make_sentinel', 'issubclass', '
1518
                 get_all_subclasses'
1519
                 class functions: [
1520
                           classproperty: ['__init__', '__get__']
1521
        file: boltons/strutils.py
1522
                 class: ['DeaccenterDict', 'HTMLTextExtractor', 'MultiReplace']
1523
                 static functions: ['camel2under', 'under2camel', 'slugify',
                 split_punct_ws', 'unit_len', 'ordinalize', 'cardinalize', '
singularize', 'pluralize', '_match_case', 'find_hashtags', 'al0n
1524
1525
                  ', 'strip_ansi', 'asciify', 'is_ascii', 'bytes2human', 'html2text
1526
                 ', 'gunzip_bytes', 'gzip_bytes', 'iter_splitlines', 'indent', 'is_uuid', 'escape_shell_args', 'args2sh', 'args2cmd', 'parse_int_list', 'format_int_list', 'complement_int_list', '
1527
1528
1529
                 int_ranges_from_int_list', 'multi_replace', 'unwrap_text']
                 class functions: [
1530
                           DeaccenterDict: ['__missing__', '__getitem__']
1531
                           HTMLTextExtractor: ['__init__', 'handle_data',
1532
                           handle_charref', 'handle_entityref', 'get_text']
1533
                           MultiReplace: ['__init__', '_get_value', 'sub']
1534
1535
        file: boltons/cacheutils.py
                 class: ['RLock', 'LRI', 'LRU', '_HashedKey', 'CachedFunction', '
1536
                 CachedMethod', 'cachedproperty', 'ThresholdCounter', 'MinIDMap']
                 static functions: ['make_cache_key', 'cached', 'cachedmethod']
1538
                 class functions: [
1539
                           RLock: ['__enter__', '__exit__']
                           LRI: ['__init__', '_init_ll', '_print_ll', '
1540
                           1541
                             _set_key_and_add_to_front_of_ll', '
1542
                           _set_key_and_evict_last_in_ll', '_remove_from_ll', '
                            __setitem__', '__getitem__', 'get', '__delitem__', 'pop', 'popitem', 'clear', 'copy', 'setdefault', 'update', '
1543
1544
                           _eq_', '__ne_', '__repr__']
LRU: ['__getitem_']
1545
1546
                           _HashedKey: ['__init__', '__hash__', '__repr__']
CachedFunction: ['__init__', '__call__', '__repr__'
CachedMethod: ['__init__', '__get__', '__call__', '
1547
1548
1549
                            __repr__']
                           cachedproperty: ['__init__', '__get__', '__repr__']
ThresholdCounter: ['__init__', 'threshold', 'add', '
1550
1551
                           elements', 'most_common', 'get_common_count', '
1552
                           get_uncommon_count', 'get_commonality', '__getitem__', '
                           __len__', '__contains__', 'iterkeys', 'keys', 'itervalues', 'values', 'iteritems', 'items', 'get', 'update']
1553
1554
                           MinIDMap: ['__init__', 'get', 'drop', '_clean', '
__contains__', '__iter__', '__len__', 'iteritems']
1555
1556
1557
        file: boltons/debugutils.py
1558
                 class: []
1559
                 static functions: ['pdb_on_signal', 'pdb_on_exception', '
1560
                 trace_print_hook', 'wrap_trace']
                 class functions: [
1561
1562
        file: boltons/deprutils.py
1563
                 class: ['DeprecatableModule']
1564
                 static functions: ['deprecate_module_member']
1565
                 class functions: [
                           DeprecatableModule: ['__init__', '__getattribute__']
```

```
1566
1567
        file: boltons/dictutils.py
1568
                   class: ['OrderedMultiDict', 'FastIterOrderedMultiDict', 'OneToOne
1569
                   ', 'ManyToMany', 'FrozenHashError', 'FrozenDict']
                   static functions: ['subdict']
1570
                   class functions: [
1571
                             OrderedMultiDict: ['__new__', '__init__', '__getstate__',
    '__setstate__', '_clear_ll', '_insert', 'add', 'addlist
    ', 'get', 'getlist', 'clear', 'setdefault', 'copy', '
    fromkeys', 'update', 'update_extend', '__setitem__', '
    __getitem__', '__delitem__', '__eq__', '__ne__', '__ior__
    ', 'pop', 'popall', 'poplast', '_remove', '_remove_all',
    '.iteritoms', 'iteritoms', 'iteritoms', 'inserted
1572
1573
1574
1575
1576
                              'iteritems', 'iterkeys', 'itervalues', 'todict', 'sorted
1577
                              ', 'sortedvalues', 'inverted', 'counts', 'keys', 'values
                             ', 'items', '__iter__', '__reversed__', '__repr__', 'viewkeys', 'viewvalues', 'viewitems']
1578
1579
                             FastIterOrderedMultiDict: ['_clear_ll', '_insert', '
1580
                             _remove', '_remove_all', 'iteritems', 'iterkeys', _reversed__']
1581
1582
                             OneToOne: ['__init__', 'unique', '__setitem__', '
__delitem__', 'clear', 'copy', 'pop', 'popitem',
1583
                             setdefault', 'update', '__repr__']
ManyToMany: ['__init__', 'get', '__getitem__', '
1584
1585
                              __setitem__', '__delitem__', 'update', 'add', 'remove', '
1586
                              replace', 'iteritems', 'keys', '__contains__', '__iter_
1587
                              ', '__len__', '__eq__', '__repr__']
1588
                              FrozenHashError: []
                             FrozenDict: ['updated', 'fromkeys', '__repr__', '
__reduce_ex__', '__hash__', '__copy__', '
1589
1590
                             _raise_frozen_typeerror']
                  1
1592
        file: boltons/easterutils.py
1593
                  class: []
                  static functions: ['gobs_program']
1594
                   class functions: [
1595
1596
1597
        ###
1598
        Please think step by step and give you reasoning process first. Finally,
1599
        provide full path and return top 5 files.
1600
1601
        The returned files should be separated by new lines ordered by most to
1602
        least important and wrapped with '''
1603
        For example:
1604
        file1.py
1605
        file2.py
1606
        file3.py
1607
        file4.py
1608
        file5.py
1609
        Note: file1.py indicates the top-1 file, file2.py indicates the top-2
1610
        file, and so on. Do not include test files.
1611
1612
1613
        >>>>>> assistant:
1614
        To solve the problem of adding fill parameters to the 'pairwise' and '
1615
        windowed' functions, we need to identify the files that contain these
1616
        functions and any related utility functions that might be affected by or
1617
        need to support the new feature.
1618
1619
        ### Step-by-Step Reasoning:
```

```
1620
      1. **Identify the Functions:**
1621
          - The problem description specifically mentions 'pairwise_iter' and
1622
         implies a similar change might be needed for 'windowed' functions. We
1623
         need to locate these functions in the codebase.
1624
      2. **Locate the Functions:**
1625
         - From the provided internal structure of the files, we see that '
1626
         pairwise', 'pairwise_iter', 'windowed', and 'windowed_iter' are static
1627
           functions in 'boltons/iterutils.py'.
1628
1629
      3. **Assess Dependencies:**
          - Since these functions are in 'boltons/iterutils.py', this file is
1630
         the primary candidate for modification.
1631
         - We should also consider if any other files might be using these
1632
         functions or if there are utility functions that might need to be
         updated to support the new 'fill' parameter.
1633
1634
      4. **Check for Related Utilities:**
1635
          - The `boltons/iterutils.py` file might have other utility functions
1636
         that interact with 'pairwise' and 'windowed'. We should ensure that
1637
         any such functions are also updated if necessary.
1638
      5. **Consider Other Files:**
1639
         - While other files like 'boltons/funcutils.py', 'boltons/listutils.py
1640
         ', etc., are part of the repository, they do not contain the functions
1641
          in question or directly relate to the problem at hand based on the
1642
         provided information.
1643
      6. **Rank the Files:**
1644
          - Given the focus on 'pairwise' and 'windowed', 'boltons/iterutils.py'
1645
          is the most relevant file.
1646
         - Other files are less relevant unless they are found to directly use
1647
         or extend the functionality of 'pairwise' and 'windowed'.
1648
      ### Conclusion:
1649
1650
      Based on the above reasoning, the top file to edit is 'boltons/iterutils.
1651
      py'. Other files are not directly relevant to the problem as described,
1652
      so they are not included in the top 5 list.
1653
1654
      boltons/iterutils.py
1655
1656
```

#### F.2 COMPONENT LOCALIZATION

```
1674
       project. 'file_name' -> The name of the file. 'class_name' -> The name
1675
       of the class. 'func_name' -> The name of the function. *
1676
       * get_code_of_file_function('file_name', 'func_name') -> Get the code of
1677
       a specified function in the given file and python project. 'file_name' ->
       The name of the file. 'func_name' -> The name of the function. *
1678
       * get_toplevel_code('file_name') -> Get all the code in a given file that
1679
       is not part of a class or function definition. This is useful for
1680
       viewing imports, global variables, constants, and any top-level script
1681
       logic. 'file_name' -> The name of the file. *
1682
       * exit() -> Exit function calling to give your final answer when you are
       confident of the answer. *
1683
1684
       You have 10 chances to call function.
1685
1686
       >>>>>> iser:
1687
       ### Problem Description ###
1688
       <request>
1689
       (PUP-5548) Use 'service' in debian provider to interact with all init
1690
       systems
1691
      In the Debian family, the 'service' command provides an
1692
       abstraction around available init systems, allowing it to be used
       to properly determine service state regardless of which init system
1693
      is active on the system.
1694
1695
       This commit updates the Debian provider's 'statuscmd' method to use
1696
       'service' at all times, allowing the provider to properly query
       services regardless of OS version or init system in place.
1697
       </request>
1698
1699
       ###
1700
1701
      Let's locate the relevant elements (function/class) step by step using
       reasoning and function calls.
1702
       I have pre-identified top-5 relevant files. There stuctures are as
1703
       follows:
1704
       file: lib/puppet/provider/service/debian.rb
1705
               class: []
1706
               static functions: ['disable', 'enabled?', 'get_start_link_count
               ', 'enable', 'statuscmd']
1707
               class functions: [
1708
1709
       file: lib/puppet/provider/service/base.rb
1710
               class: []
1711
               static functions: ['getps', 'getpid', 'status', 'statuscmd', '
start', 'startcmd', 'stop', 'stopcmd']
1712
               class functions: [
1713
1714
       file: lib/puppet/type/service.rb
1715
               class: ['Puppet']
1716
               static functions: []
               class functions: [
1717
                        Puppet: ['retrieve', 'retrieve', 'sync', 'refresh']
1718
1719
       file: lib/puppet/provider/service/systemd.rb
1720
               class: []
1721
               static functions: ['instances', 'disable', 'get_start_link_count
               ', 'enabled?', 'debian_enabled?', 'status', 'enable', 'mask', '
1722
               unmask', 'restartcmd', 'startcmd', 'stopcmd']
1723
               class functions: [
1724
1725
       file: lib/puppet/provider/service/init.rb
1726
               class: []
               static functions: ['defpath', 'excludes', 'instances', '
get_services', 'hasstatus=', 'initscript', 'paths', 'search', '
1727
```

```
1728
              startcmd', 'stopcmd', 'restartcmd', 'texecute', 'statuscmd', '
1729
              is_init?']
1730
              class functions: [
1731
      The formal parameter 'file_name' takes the value in "file:"
1732
      The formal parameter 'class_name' takes the value in "class:"
1733
      The formal parameter 'func_name' takes the value in "static functions:"
1734
      and "class functions: "
1735
      Avoid making multiple identical calls to save overhead.
1736
      You must strictly follow the structure I give to call different tools.
      For static functions, you can use 'get_code_of_file_function', and for
1737
      class functions, you can use 'get_code_of_class_function'.
1738
      In order to locate accurately, you can pre-select 7 locations, then check
1739
      them through function calls, and finally confirm 5 file names.
1740
      Don't make the first function call in this message.
1741
1742
      >>>>>> assistant:
1743
      To address the problem, we need to focus on the 'statuscmd' method in the
1744
      Debian provider, as the request specifically mentions updating this
1745
      method to use the 'service' command. Based on the provided file structure
      , the relevant files and functions to consider are:
1746
1747
      1. **lib/puppet/provider/service/debian.rb**
1748
         - Function: 'statuscmd'
1749
1750
      2. **lib/puppet/provider/service/base.rb**
         - Function: 'statuscmd'
1751
1752
      3. **lib/puppet/provider/service/systemd.rb**
1753
         - Function: 'statuscmd'
1754
1755
      4. **lib/puppet/provider/service/init.rb**
         - Function: 'statuscmd'
1756
1757
      These files contain the 'statuscmd' function, which is likely related to
1758
      the service status command functionality. The problem description
1759
      specifically mentions the Debian provider, so the 'debian.rb' file is of
1760
      primary interest. However, examining the 'statuscmd' function in other
      related files might provide additional context or necessary changes.
1761
1762
      Let's proceed with checking the code of these functions to identify the
1763
      necessary modifications. I will start by examining the 'statuscmd'
1764
      function in 'lib/puppet/provider/service/debian.rb'.
1765
      >>>>> user:
1766
1767
      Now call a function in this format 'FunctionName (Argument)' in a single
1768
      line without any other word or signal (such as ''').
1769
      Don't call the same function you've previous called, because this may
1770
      waste your context length.
1771
1772
      >>>>>> assistant:
1773
      get_code_of_file_function('lib/puppet/provider/service/debian.rb', '
1774
      statuscmd')
1775
      >>>>> user:
1776
      I have already checked this function/class is related to the task.
1777
1778
      >>>>> user:
1779
        def statuscmd
1780
          os = Facter.value(:operatingsystem).downcase
1781
          if os == 'debian'
```

```
1782
            majversion = Facter.value(:operatingsystemmajrelease).to_i
1783
          else
1784
            majversion = Facter.value(:operatingsystemmajrelease).split('.')
1785
            [01.to i
          end
1786
1787
1788
          if ((os == 'debian' && majversion >= 8) || (os == 'ubuntu' &&
1789
          majversion >= 15))
1790
            # SysVInit scripts will always return '0' for status when the
            service is masked,
1791
            # even if the service is actually stopped. Use the SysVInit-Systemd
1792
             compatibility
1793
            # layer to determine the actual status. This is only necessary when
1794
             the SysVInit
            # version of a service is queried. I.e, 'ntp' instead of 'ntp.
1795
            service'.
1796
            (@resource[:hasstatus] == :true) && ["systemctl", "is-active",
1797
            @resource[:name]]
1798
          else
1799
            super
1800
          end
        end
1801
1802
      >>>>> user:
1803
1804
      Now call a function in this format 'FunctionName (Argument)' in a single
      line without any other word or signal (such as ''').
1805
      Don't call the same function you've previous called, because this may
1806
      waste your context length.
1807
1808
      You can check the functions that this part of code calls.
1809
1810
      >>>>>> assistant:
1811
      get_code_of_file_function('lib/puppet/provider/service/base.rb', '
1812
      statuscmd')
1813
1814
      >>>>> user:
      I have already checked this function/class is related to the task.
1815
1816
      >>>>>> user:
1817
        def statuscmd
1818
        end
1819
      >>>>> user:
1820
1821
      Now call a function in this format 'FunctionName (Argument)' in a single
1822
      line without any other word or signal (such as ''').
1823
      Don't call the same function you've previous called, because this may
1824
      waste your context length.
1825
      You can check the functions that this part of code calls.
1826
1827
1828
      >>>>>> assistant:
1829
      get_code_of_file_function('lib/puppet/provider/service/systemd.rb', '
      statuscmd')
1830
1831
      >>>>> user:
1832
      The call of get_code_of_file_function('lib/puppet/provider/service/
1833
      systemd.rb', 'statuscmd') encountered error. Maybe you retrieved a
1834
      location which does not exist.
1835
      >>>>> user:
```

```
1836
1837
      Now call a function in this format 'FunctionName(Argument)' in a single
1838
      line without any other word or signal (such as ''').
1839
      Don't call the same function you've previous called, because this may
      waste your context length.
1840
1841
1842
      >>>>>> assistant:
1843
      get_code_of_file_function('lib/puppet/provider/service/init.rb', '
1844
      statuscmd')
1845
      >>>>>> user:
1846
      I have already checked this function/class is related to the task.
1847
1848
      >>>>> user:
       def statuscmd
1849
          (@resource[:hasstatus] == :true) && [initscript, :status]
1850
1851
1852
      >>>>> user:
1853
1854
      Now call a function in this format 'FunctionName (Argument)' in a single
      line without any other word or signal (such as ''').
1855
      Don't call the same function you've previous called, because this may
1856
      waste your context length.
1857
1858
      You can check the functions that this part of code calls.
1859
1860
      >>>>>> assistant:
1861
      get_code_of_file_function('lib/puppet/provider/service/init.rb', '
1862
      initscript')
1863
      >>>>> user:
1864
      I have already checked this function/class is not related to the task.
1865
1866
      >>>>> user:
1867
        def initscript
1868
          @initscript ||= self.search(@resource[:name])
1869
1870
      >>>>> user:
1871
1872
      Now call a function in this format 'FunctionName(Argument)' in a single
1873
      line without any other word or signal (such as ''').
      Don't call the same function you've previous called, because this may
1874
      waste your context length.
1875
1876
      Don't check the functions that this part of code calls.
1877
1878
      >>>>>> assistant:
1879
      get_toplevel_code(lib/puppet/provider/service/debian.rb)
1880
1881
      >>>>> user:
1882
      I have already checked this function/class is related to the task.
1883
      >>>>> user:
1884
      # Manage debian services. Start/stop is the same as InitSvc, but enable/
1885
      disable
1886
      # is special.
1887
      Puppet::Type.type(:service).provide :debian, :parent => :init do
1888
        desc <<-EOT
          Debian's form of 'init'-style management.
1889
```

```
1890
           The only differences from 'init' are support for enabling and
1891
           disabling
1892
           services via 'update-rc.d' and the ability to determine enabled
1893
           status via
           'invoke-rc.d'.
1894
1895
         EOT
1896
1897
         commands :update_rc => "/usr/sbin/update-rc.d"
1898
         # note this isn't being used as a command until
         # http://projects.reductivelabs.com/issues/2538
1899
         # is resolved.
1900
         commands :invoke_rc => "/usr/sbin/invoke-rc.d"
1901
1902
         defaultfor :operatingsystem => :cumuluslinux
         defaultfor :operatingsystem => :debian, :operatingsystemmajrelease =>
1903
         ['5','6','7']
1904
1905
         # Remove the symlinks
1906
1907
      end
1908
       . . .
1909
      >>>>>> user:
1910
1911
      Now call a function in this format 'FunctionName(Argument)' in a single
1912
      line without any other word or signal (such as ''').
1913
      Don't call the same function you've previous called, because this may
      waste your context length.
1914
1915
      You can check the functions that this part of code calls.
1916
1917
      >>>>>> assistant:
1918
      exit()
1919
1920
      >>>>> user:
1921
1922
      Based on the available information, reconfirm and provide complete names
      of the most likely edit locations (10 locations at most).
1923
      Before make the final decision, please check whether the function name is
1924
       correct or not, for static functions, don't add class name.
1925
      file: lib/puppet/provider/service/debian.rb
1926
               class: []
1927
               static functions: ['disable', 'enabled?', 'get_start_link_count
               ', 'enable', 'statuscmd']
1928
               class functions: [
1929
1930
       file: lib/puppet/provider/service/base.rb
1931
               class: []
               static functions: ['getps', 'getpid', 'status', 'statuscmd', '
1932
               start', 'startcmd', 'stop', 'stopcmd']
1933
               class functions: [
1934
1935
      file: lib/puppet/type/service.rb
1936
               class: ['Puppet']
               static functions: []
1937
1938
               class functions: [
                       Puppet: ['retrieve', 'retrieve', 'sync', 'refresh']
1939
1940
      file: lib/puppet/provider/service/systemd.rb
1941
               class: []
1942
               static functions: ['instances', 'disable', 'get_start_link_count
', 'enabled?', 'debian_enabled?', 'status', 'enable', 'mask', '
1943
               unmask', 'restartcmd', 'startcmd', 'stopcmd']
```

```
1944
               class functions: [
1945
1946
      file: lib/puppet/provider/service/init.rb
               class: []
1947
               static functions: ['defpath', 'excludes', 'instances', '
1948
               get_services', 'hasstatus=', 'initscript', 'paths', 'search', '
1949
               startcmd', 'stopcmd', 'restartcmd', 'texecute', 'statuscmd',
1950
               is init?'
1951
               class functions: [
1952
1953
      Please provide the complete set of locations as either a class name, a
1954
      function name, or a file name.
1955
      The returned files should be separated by new lines ordered by most to
1956
      least important and wrapped with '''
      Since your answer will be processed automatically, please give your
1957
      answer in the exapmle format as follows.
1958
1959
      top1_file_fullpath.py
1960
      function: Class1.Function1
1961
1962
      top2_file_fullpath.py
      function: Function2
1963
1964
      top3_file_fullpath.py
1965
      class: Class3
1966
1967
      top4_file_fullpath.py
      function: Class4.Function4
1968
1969
      top5_file_fullpath.py
1970
      function: Function5
1971
      top6_file_fullpath.py
1972
      global
1973
1974
      top7_file_fullpath.py
1975
      function: Class7.Function7
1976
      Replace the 'Top_file_fullpath.py' with the actual file path, the 'Class'
1977
       with the actual class name and the 'Function' with the actual function
1978
      name. 'global' means the code is not in a class or function.
1979
      For example,
1980
1981
      sklearn/linear_model/__init__.py
      function: LinearRegression.fit
1982
1983
1984
1985
      >>>>>> assistant:
1986
      lib/puppet/provider/service/debian.rb
1987
      global
1988
1989
      lib/puppet/provider/service/debian.rb
1990
      function: statuscmd
1991
      lib/puppet/provider/service/base.rb
1992
      function: statuscmd
1993
1994
      lib/puppet/provider/service/init.rb
1995
      function: statuscmd
1996
1997
```

# F.3 CODE EDITING

1998

```
2000
2001
      >>>>> user:
2002
       We are currently solving the following task within our repository. Here
2003
      is the task description.
2004
2005
2006
       ### Task Description ###
       <request>
2007
       Add port mapping on run
2008
       CF issue : https://github.com/jolicode/docker-starter/issues/31
2009
       </request>
2010
2011
      <discussion>
2012
      Mapping ports on "run" command
2013
2014
2015
      Is there any reason to avoid mapping ports on "run" command ?
2016
      https://github.com/jolicode/docker-starter/blob/
      c319756cc1ab1de8e4a3f29bea306695752f84c2/fabfile.py#L179
2017
2018
      I just needed to add \operatorname{--service-ports} option to the fabfile.py to run my
2019
      container with port mapping embed.
2020
2021
      :kissing_heart:
2022
2023
       since we use docker-compose, it's better to set the port mapping in it.
2024
      But if you need something special, you can use port mapping via the run
2025
      command.
2026
      Feel free to modify your fabfile. This docker-starter is not a lib, just
2027
       a starter that must be customized according to your needs.
2028
       --service-ports does actually the mapping defined in docker-compose.yaml.
2029
2030
2031
       </discussion>
2032
       ###
2033
2034
       Below are some code segments, each from a relevant file. One or more of
2035
       these files may need to be edited to solve the task.
2036
2037
       --- BEGIN FILE ---
2038
       . . .
2039
       ### fabfile.py
2040
2041
       env.extra_domains = []
       # This is the host directory containing your PHP application
2042
       env.project_directory = 'application'
2043
2044
       services_to_build_first = [
2045
           'php-base',
           'builder',
2046
       1
2047
2048
2049
       def with_builder(func):
2050
           @wraps(func)
2051
           def decorated(*args, **kwargs):
               compose_files = env.compose_files[:]
```

```
2052
               env.compose_files = ['docker-compose.builder.yml'] + env.
2053
               compose_files
2054
               ret = func(*args, **kwargs)
2055
               env.compose_files = compose_files
2056
               return ret
2057
           return decorated
2058
2059
2060
       @with_builder
       def build():
2061
           11 11 11
2062
           Build the infrastructure
2063
2064
           command = 'build'
           command += ' --build-arg PROJECT_NAME=%s' % env.project_name
2065
           command += ' --build-arg USER_ID=%s' % env.user_id
2066
2067
           for service in services_to_build_first:
2068
               commandForService = '%s %s' % (command, service)
2069
               docker_compose(commandForService)
2070
           docker_compose (command)
2071
2072
2073
       @task
2074
       def up():
2075
           Build and start the infrastructure
2076
           11 11 11
2077
           build()
2078
           docker_compose('up --remove-orphans -d')
2079
2080
       . . .
2081
       @task
2082
       @with_builder
2083
       def destroy():
2084
2085
           Clean the infrastructure (remove container, volume, networks)
2086
           docker_compose('down --volumes --rmi=local')
2087
2088
2089
       def run_in_docker_or_locally_for_dinghy(command):
2090
           Mac users have a lot of problems running Yarn / Webpack on the Docker
2091
           stack so this func allow them to run these tools on their host
2092
2093
           if env.dinghy:
2094
               local('cd %s && %s' % (env.project_directory, command))
2095
           else:
               docker_compose_run(command)
2096
2097
2098
       def docker_compose(command_name):
2099
           domains = ''' + '', ''.join([env.root_domain] + env.extra_domains) +
2100
2101
2102
           localEnv = {
2103
               'PROJECT_NAME': env.project_name,
2104
               'PROJECT_DIRECTORY': env.project_directory,
2105
               'PROJECT_ROOT_DOMAIN': env.root_domain,
               'PROJECT_DOMAINS': domains,
```

```
2106
2107
2108
           with shell_env(**localEnv):
2109
               local('docker-compose -p %s %s %s' % (
                   env.project_name,
2110
                    ' '.join('-f ' + env.root_dir + '/infrastructure/docker/' +
2111
                   file for file in env.compose_files),
2112
                   command_name
2113
               ))
2114
2115
       def docker_compose_run(command_name, service="builder", user="app",
2116
       no_deps=False, workdir=None):
2117
           args = [
2118
               'run',
               '--rm '
2119
               '-u %s ' % _shell_escape(user),
2120
2121
2122
           if no_deps:
2123
               args.append('--no-deps')
2124
           if workdir is not None:
2125
               args.append('-w %s ' % _shell_escape(workdir))
2126
2127
           docker_compose('%s %s /bin/bash -c "exec %s"' % (
2128
               ' '.join(args),
2129
               _shell_escape(service),
               _shell_escape(command_name)
2130
           ))
2131
2132
2133
       def set_local_configuration():
           env.compose_files = ['docker-compose.yml']
2134
           env.user_id = int(local('id -u', capture=True))
2135
           env.root_dir = os.path.dirname(os.path.abspath(__file___))
2136
2137
           if env.user_id > 256000:
2138
               env.user_id = 1000
2139
       . . .
2140
2141
2142
2143
       --- END FILE ---
2144
2145
      Please first localize the positions to edit based on the task statement,
2146
       and then output the files that need to be deleted, modified or added.
2147
       '- file' means deleting the file;
2148
       '\star file' means modificating the file;
       '+ file' means adding the file.
2149
       The file should list like below:
2150
       * * *
2151
      - file1.py
2152
       * file2.py
2153
       * file3.py
       + file4.py
2154
2155
2156
      To solve the task, you should then generate *SEARCH/REPLACE* edits.
2157
2158
      Every *SEARCH/REPLACE* edit must use this format:
2159
       1. The file path
      2. The start of search block: <<<< SEARCH
```

```
2160
      3. A contiguous chunk of lines to search for in the existing source code
2161
      4. The dividing line: ======
2162
      5. The lines to replace into the source code
2163
      6. The end of the replace block: >>>>> REPLACE
2164
      Here is an example:
2165
2166
       '''python
2167
       ### mathweb/flask/app.py
2168
      <<<<< SEARCH
      from flask import Flask
2169
2170
      import math
2171
      from flask import Flask
2172
      >>>>> REPLACE
2173
2174
      Please note that the *SEARCH/REPLACE* edit REQUIRES PROPER INDENTATION.
2175
      If you would like to add the line '
                                                   print(x)', you must fully
2176
      write that out, with all those spaces before the code!
      Wrap the *SEARCH/REPLACE* edit in blocks '''python...'''.
2177
2178
      When multiple edits should be done, please output *SEARCH/REPLACE* edit
      one by one and give your reasoning process before each *SEARCH/REPLACE*
2179
      block.
2180
2181
2182
       >>>>> assistant:
2183
      The plan of solving this software task:
      To address the task of adding port mapping on the "run" command, we need
2184
      to modify the 'docker_compose_run' function in the 'fabfile.py'. The
2185
      discussion suggests using the '--service-ports' option to enable port
2186
      mapping as defined in the 'docker-compose.yml' file. We will add a
2187
      parameter to the 'docker_compose_run' function to allow this option to be
2188
       toggled.
2189
      The files that should be edited include:
2190
      * fabfile.py
2191
2192
      Here are my edits for code.
2193
      The motivation of this edit:
2194
      We need to modify the 'docker_compose_run' function to include an
2195
      optional parameter for port mapping. By adding a 'port_mapping' parameter
2196
       , we can control whether the '--service-ports' option is included in the
2197
      Docker run command. This change allows users to enable port mapping when
      running a service, which is useful for exposing ports as defined in the \dot{}
2198
      docker-compose.yml' file.
2199
2200
       '''python
2201
       ### fabfile.py
2202
      <<<<< SEARCH
2203
              ))
2204
2205
      def docker_compose_run(command_name, service="builder", user="app",
2206
      no_deps=False, workdir=None):
2207
          args = [
              'run',
2208
               '--rm',
2209
      _____
2210
              ))
2211
2212
2213
       def docker_compose_run(command_name, service="builder", user="app",
      no_deps=False, workdir=None, port_mapping=False):
```

```
2214
           args = [
2215
               'run ',
               '--rm ',
2216
2217
      >>>>> REPLACE
      <<<<< SEARCH
2218
           if no_deps:
2219
               args.append('--no-deps')
2220
2221
           if workdir is not None:
2222
               args.append('-w %s ' % _shell_escape(workdir))
2223
      _____
2224
           if no_deps:
2225
               args.append('--no-deps')
2226
           if port_mapping:
2227
               args.append('--service-ports')
2228
2229
           if workdir is not None:
2230
               args.append('-w %s ' % _shell_escape(workdir))
2231
2232
      >>>>> REPLACE
2233
2234
```

#### F.4 CODE COMPLETION

2235

```
2237
     >>>>>> user:
2238
2239
     We are currently solving the following task within our repository. Here
2240
     is the task description.
2241
2242
     ### Task Description ###
2243
     Please help me to complete the code at the position of '[TODO]', in the
2244
     b2/sync/scan_policies.py
2245
         def __init__(
2246
            self,
2247
            exclude_dir_regexes=tuple(),
            exclude_file_regexes=tuple(),
2248
             include_file_regexes=tuple(),
2249
         ):
2250
2251
      ###
2252
     Below are some code segments, each from a relevant file. One or more of
2253
     these files may need to be edited to solve the task.
2254
2255
2256
      --- BEGIN FILE ---
      . . .
2257
     ### b2/sync/scan_policies.py
2258
     2259
2260
     # File: b2/sync/scan_policies.py
2261
      # Copyright 2018 Backblaze Inc. All Rights Reserved.
2262
2263
      # License https://www.backblaze.com/using_b2_code.html
2264
2265
      2266
2267
      import logging
     import re
```

```
2268
2269
       logger = logging.getLogger(__name__)
2270
2271
       class RegexSet (object):
2272
2273
           Holds a (possibly empty) set of regular expressions, and knows how to
2274
            check
2275
           whether a string matches any of them.
           11 11 11
2276
2277
           def __init__(self, regex_iterable):
2278
               self._compiled_list = [re.compile(r) for r in regex_iterable]
2279
2280
           def matches(self, s):
               return any(c.match(s) is not None for c in self._compiled_list)
2281
2282
2283
       class ScanPoliciesManager(object):
2284
2285
           Policy object used when scanning folders for syncing, used to decide
2286
           which files to include in the list of files to be synced.
2287
2288
           def __init__(
2289
               self,
2290
               exclude_dir_regexes=tuple(),
2291
               exclude_file_regexes=tuple(),
               include_file_regexes=tuple(),
2292
           ):
2293
               # [TODO]
2294
2295
           def should_exclude_file(self, file_path):
2296
               Given the full path of a file, should it be excluded from the
2297
2298
2299
               :param file_path: The path of the file, relative to the root
2300
               directory
2301
                                  being scanned.
               :return: True iff excluded.
2302
2303
               return self._exclude_file_set.matches(file_path) and \
2304
                      not self._include_file_set.matches(file_path)
2305
           def should_exclude_directory(self, dir_path):
2306
2307
               Given the full path of a directory, should all of the files in it
2308
2309
               excluded from the scan?
2310
               :param dir_path: The path of the directory, relative to the root
2311
               directory
2312
                                 being scanned. The path will never end in '/'.
2313
               :return: True iff excluded.
2314
2315
               return self._exclude_dir_set.matches(dir_path)
2316
2317
      DEFAULT_SCAN_MANAGER = ScanPoliciesManager()
2318
2319
2320
2321
      --- END FILE ---
```

```
2322
2323
      Please first localize the positions to edit based on the task statement,
2324
      and then output the files that need to be deleted, modified or added.
2325
      '- file' means deleting the file;
      '* file' means modificating the file;
2326
      '+ file' means adding the file.
2327
      The file should list like below:
2328
2329
      - file1.py
2330
      * file2.py
      * file3.py
2331
      + file4.py
2332
2333
2334
      To solve the task, you should then generate *SEARCH/REPLACE* edits.
2335
      Every *SEARCH/REPLACE* edit must use this format:
2336
      1. The file path
2337
      2. The start of search block: <<<< SEARCH
2338
      3. A contiguous chunk of lines to search for in the existing source code
2339
      4. The dividing line: ======
2340
      5. The lines to replace into the source code
      6. The end of the replace block: >>>>> REPLACE
2341
2342
      Here is an example:
2343
2344
      '''python
      ### mathweb/flask/app.py
2345
      <<<<< SEARCH
2346
      from flask import Flask
2347
      _____
2348
      import math
2349
      from flask import Flask
      >>>>> REPLACE
2350
2351
2352
      Please note that the *SEARCH/REPLACE* edit REQUIRES PROPER INDENTATION.
2353
      If you would like to add the line '
                                                  print(x)', you must fully
2354
      write that out, with all those spaces before the code!
      Wrap the *SEARCH/REPLACE* edit in blocks '''python...'''.
2355
      When multiple edits should be done, please output *SEARCH/REPLACE* edit
2356
      one by one and give your reasoning process before each *SEARCH/REPLACE*
2357
      block.
2358
2359
      >>>>>> assistant:
2360
      The plan of solving this software task:
2361
      The task requires completing the '__init__' method of the '
2362
      ScanPoliciesManager' class in the 'b2/sync/scan_policies.py' file. The
2363
      method should initialize the class attributes using the provided regular
2364
      expression tuples. The class uses 'RegexSet' objects to manage and check
      regular expressions. Therefore, the '__init__' method should create '
2365
      RegexSet' instances for each of the regular expression tuples:
2366
      exclude_dir_regexes', 'exclude_file_regexes', and 'include_file_regexes'.
2367
2368
      The files that should be edited include:
2369
      * b2/sync/scan_policies.py
2370
      Here are my edits for code.
2371
2372
      The motivation of this edit:
2373
      The 'ScanPoliciesManager' class is designed to manage policies for
2374
      scanning directories and files, determining which should be included or
      excluded based on regular expressions. The '__init__' method needs to
2375
      initialize the class attributes '_exclude_dir_set', '_exclude_file_set',
```

```
2376
      and '_include_file_set' using the 'RegexSet' class. This setup allows the
2377
       'should_exclude_file' and 'should_exclude_directory' methods to function
2378
       correctly by checking file and directory paths against the compiled
2379
      regular expressions.
2380
       '''python
2381
      ### b2/sync/scan_policies.py
2382
      2383
              exclude_file_regexes=tuple(),
2384
              include_file_regexes=tuple(),
          ):
2385
               # [TODO]
2386
2387
          def should_exclude_file(self, file_path):
2388
2389
              exclude_file_regexes=tuple(),
2390
               include_file_regexes=tuple(),
2391
2392
              self._exclude_dir_set = RegexSet(exclude_dir_regexes)
2393
              self._exclude_file_set = RegexSet(exclude_file_regexes)
              self._include_file_set = RegexSet(include_file_regexes)
2394
2395
          def should_exclude_file(self, file_path):
2396
2397
      >>>>> REPLACE
2398
2399
```