

SWE-Ext: Extending and Scaling Augmented Data for Mid-Training of Repository-Level Coding Tasks

Anonymous ACL submission

Abstract

Repository-level benchmarks such as SWE-Bench have highlighted the challenges of scaling language models to complex software engineering tasks. However, the open-sourced real data with verifiable environments from which we can get trajectories to perform rejection fine-tuning, is scarce. In such a situation, we propose training by augmented data as a mid-training stage. Previous augmented data primarily focused on monolingual issue resolving and feature implementation. In this work, we introduce **SWE-Ext**, an effort to scale and extend 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. Augmented data and these extensions yield consistent improvements on Python repository-level benchmarks like SWE-Bench and FEA-Bench, and assist post-training performance on them. Our method offers a simple yet effective way to leverage more open-source 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 training strategies, which mainly fall into two categories: **augmented data** for direct **supervised fine-tuning (SFT)** and **verified data** for **rejection fine-tuning (RFT)**. 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 based on rejection sampling, 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 sandboxes 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. So, we propose training using augmented data as an individual mid-training stage.

Despite their advantages, existing augmented

084 datasets remain limited in scope, predominantly 131
085 focusing on Python repositories and monolingual 132
086 issue resolution or feature implementation (Xie 133
087 et al., 2025). This narrow focus underutilizes the 134
088 vast diversity of open-source GitHub data, restrict- 135
089 ing model generalization across languages and task 136
090 types. In this work, we introduce **SWE-Ext**, a 137
091 scalable pipeline to extend and augment repository- 138
092 level coding data along two orthogonal dimensions: 139
093 multilingual coverage and auxiliary task inclusion 140
094 (e.g., code completion). Starting from high-quality 141
095 PRs crawled from GitHub Archive ¹ across ten 142
096 programming languages (Python, Go, JavaScript, 143
097 Ruby, PHP, Java, TypeScript, C#, C++, and C), 144
098 we construct datasets of four complementary sub- 145
099 tasks: (1) file localization to identify relevant files,
100 (2) component localization to pinpoint functions
101 or methods, (3) code editing to generate patches
102 for code repositories, and (4) code completion de-
103 rived from short-description PRs to provide code
104 in-filling on fixed-positions. Our contributions are
105 as follows:

- 106 • We propose augmented data for mid-training of 146
107 repository-level coding tasks and develop a 147
108 scalable data extension pipeline that broadens 148
109 the utility of GitHub PRs, enabling extension 149
110 across multilingual and completion-based do- 150
111 mains while preserving baseline effectiveness. 151
- 112 • We demonstrate consistent performance gains 152
113 on repository-level benchmarks like SWE- 153
114 Bench and FEA-Bench, with multilingual and 154
115 completion extensions yielding up to +1.4% 155
116 and +2.5% improvements on mid-training of 156
117 32B models, respectively, and up to +5.4% 157
118 on 7B models. The mid-training stage can 158
119 raise performance by up to +5.2% after post- 159
120 training. 160
- 121 • We provide preliminary validation of strong 161
122 cross-task and cross-language transfer effects, 162
123 showing that multilingual or completion data 163
124 can boost single-language downstream perfor- 164
125 mance and highlighting the value of diverse 165
126 augmentation for generalization. 166

127 2 Problem Definition 167

128 **Notation.** Let \mathcal{R} denote a software repository 168
129 represented as a finite collection of files $\{f_i\}_{i=1}^N$ 169
130 together with auxiliary structures (ASTs, file-level 170

skeletons, import graph, etc.). Let q denote a prob- 131
132 lem statement (e.g., a PR description, an issue body, 133
134 a failing test, or a natural-language specification). 135
136 Let \mathcal{L} be the universe of addressable program po- 137
138 sitions (file identifiers \times component spans). A 139
140 *localization* is a subset $\ell \subseteq \mathcal{L}$. The space of ed- 141
142 its (patches) is denoted \mathcal{P} ; an edit (patch) $\delta \in \mathcal{P}$ 143
144 is an operator $p : \mathcal{R} \rightarrow \mathcal{R}'$ that applies insert/re- 145
146 place/delete/move operations to positions in \mathcal{P} . Fi- 147
148 nally, $L(\mathcal{R}', q)$ is a task loss measuring how well 149
150 the edited repository \mathcal{R}' satisfies q . Note that $P(\cdot)$ 151
152 denotes a probability distribution. 153

Repository-level coding task. Given (\mathcal{R}, q) , the 143
144 objective of a repository-level coding agent is to 145
146 produce a patch $\delta \in \mathcal{P}$ that minimizes the task loss: 147

$$148 \hat{\delta} = \arg \min_{\delta \in \mathcal{P}} L(\delta(\mathcal{R}), q), \quad (1) \quad 149$$

or equivalently, under a probabilistic modeling 147
148 view,

$$149 \hat{\delta} = \arg \max_{\delta \in \mathcal{P}} P(\delta | \mathcal{R}, q). \quad (2) \quad 150$$

Decomposition into localization and editing. 150
151 Because modern LLMs cannot feasibly consume 152
153 the entirety of large repositories at once, we decom- 154
155 pose the posterior over patches by marginalizing 156
157 over possible localizations: 158

$$159 P(\delta | \mathcal{R}, q) = \sum_{\ell \in \mathcal{L}} P(\ell | \mathcal{R}, q) P(\delta | \mathcal{R}, q, \ell). \quad (3) \quad 160$$

This decomposition separates (i) a *localization* 156
157 model that identifies candidate regions ℓ of inter- 158
159 est, and (ii) an *editing* model that produces a con- 160
161 crete patch conditional on the selected regions. In 162
163 practice, localization and editing are performed in 164
165 separate iterative phases: the system first iterates to 166
167 identify a sufficiently accurate set of relevant files 168
169 or components, and only then initiates an editing 169
170 phase where patches can be generated, potentially 170
171 with feedback. 171

Code completion as a special case. Code comple- 166
167 tion corresponds to the special case where the 168
169 localization ℓ^* is given *a priori* (a fixed position or 169
170 contiguous span) and the patch is constrained to be 170
171 a continuation or replacement at that position: 171

$$172 \hat{\delta} = \arg \max_{\delta \in \mathcal{P}} P(\delta | \mathcal{R}, q, \ell^*). \quad (4) \quad 173$$

¹<https://www.gharchive.org>

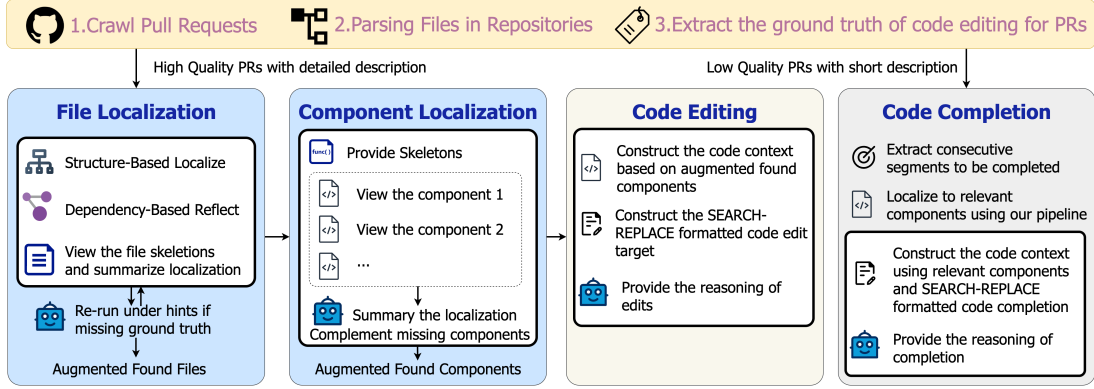


Figure 1: Data collection pipeline for 4 types of tasks in SWE-Ext. We crawl and process pull requests from GH Archive, and then process well-described PRs for the localization-editing task and process other PRs for the completion task.

3 SWE-Ext: Extend and Scale 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 \mathcal{F}_{gt} and the ground-truth modified components \mathcal{C}_{gt} . To limit tractability, we retain only PRs with $|\mathcal{F}_{\text{gt}}| \leq 5$ and $|\mathcal{C}_{\text{gt}}| \leq 10$. 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

$$s = (\mathcal{R}, q, \mathcal{F}_{\text{gt}}, \mathcal{C}_{\text{gt}}, \delta_{\text{gt}}), \quad (5)$$

where δ_{gt} is the ground-truth patch.

File localization. To approximate $P(\ell \mid \mathcal{R}, q)$ at file granularity, we construct file-localization data. We employ a three-stage pipeline driven by an expert model: (i) rank files according to the file tree of the problem statement; (ii) expand via file-level dependency to form a candidate set; and (iii) re-rank using file skeletons and the expert

model to produce a final top- n set \mathcal{S} . If $\mathcal{F}_{\text{gt}} \not\subseteq \mathcal{S}$, we re-run the expert with rationalization hints so that \mathcal{S} always covers the ground truth. The resulting dataset pairs $(\mathcal{R}, q, \mathcal{S})$ 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 $P(c \mid \mathcal{R}, q, \ell)$. 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 $(\mathcal{R}, q, \mathcal{F}, \mathcal{S}_{\text{comp}})$ thus can serve 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 δ_{gt} . For each selected localization context $\ell_{\text{ctx}} = \mathcal{S}_{\text{comp}}$, we convert patches into *before/after* code fragments (search-replace pairs δ_{sp}). To enrich supervision, an expert model generates a rationale explaining the patch. Each example $(\mathcal{R}, q, \ell_{\text{ctx}}, \delta_{\text{sp}})$ provides supervision for $P(\delta \mid \mathcal{R}, q, \ell_{\text{ctx}})$. This partition trains models to propose concrete patches given a localized context.

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

Table 1: Statistics of the dialogue data for 4 tasks in the SWE-Ext dataset. Token counts are computed using Qwen2.5-Coder. **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	4 746	56 384	9 787.8	552.1	-	1.66/4.83
Component-Localization	4 702	56 726	8 164.3	59.8	-	2.55/6.32
Code-Editing	4 319	45 858	6 447.7	955.0	32.0	-
Code-Completion	1 642	23 853	5 929.7	430.1	88.0	-

cases, we extract newly added functions or contiguous inserted regions and frame them as completion tasks with a fixed localization ℓ^* . Each example provides supervision for $P(\delta | \mathcal{R}, q, \ell^*)$. 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 $P(\delta | \mathcal{R}, q)$. By supervising different sub-tasks separately, we enable mid-training for repository-level tasks while ensuring that the final system can integrate these capabilities into a single model. For further details, please refer to Appendix D.

3.2 Dataset Characteristics

The SWE-Ext corpus is a *mid-training* dataset constructed from real-world pull requests for repository-level coding tasks. The dataset is organized into four complementary tasks: file localization, component 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 4,746 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 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 applications.

Scale, context length, and edit complexity. The dataset provides long, realistic contexts: average input contexts (**Dialogue Len**) range from \sim 5.9k tokens (completion) up to \sim 9.8k tokens (file-localization), reflecting the long context challenge for real-world 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, TypeScript, Java, C#, C++, and C. Language distributions (See Figure 3) show consistent dominance of Python and Go across modalities (roughly 27–30% and 15–21% respectively), moderate representation for JavaScript, Ruby and PHP (\approx 9–13%), 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 comprehensive utilization of the available data 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 to perform supervised fine-tuning.

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 Qwen2.5-Coder-Instruct (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.

Mid-Training Data. 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.

Post-Training Data. To validate that our mid-training recipe based on augmented data can also boost post-training performance, we use mainly DeepSeek-V3.2(Liu et al., 2025) and Claude-4.5-Sonnet(Anthropic, 2025c) to perform reject sampling on a verifiable dataset SWE-

Gym (Pan et al., 2024). An adjusted OpenHands that uses command-line and string replace tools is utilized as the scaffold. Finally, we get 3,128 OpenHands trajectories as a standard post-training dataset to perform rejection fine-tuning (RFT).

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 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. For the evaluation after RFT by OpenHands trajectories, we use OpenHands agent with maximum rounds of 100.

5 Results

5.1 Software Engineering Benchmarks

Competitive performance as a mid-training method. According to Table 2, although the absolute performance of SWE-Ext models is lower than the strongest systems, our models achieve competitive results among open-weight models even with the baseline augmented data. 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 certain improvements. The results validate our pipeline of data collection.

Consistent improvements with data extension. Table 2 and Table 3 demonstrate that extending augmented data along different dimensions yields consistent gains during mid-training. On SWE-Bench Verified, the multilingual extension improves performance from 31.2% to 32.6% (+1.4%), while the completion extension reaches 32.2%. On FEA-Bench Lite, similar trends are observed: the model trained with extended data outperforms the baseline by up to +2.5%. These results confirm that the proposed data extensions are effective in enhancing model performance.

Model	System	Expert Model	Exec	Recipe	% Resolved
<i>Closed Weight Models</i>					
GPT-4o (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%
Claude 4.5 Sonnet (Anthropic, 2025c)	-	-	-	-	77.2%
Llama3-SWE-RL-70B (Wei et al., 2025)	Agentless	-	-	RL	41.0%
<i>Open Weight Models</i>					
DeepSeek-V3-671B (Liu et al., 2024)	Agentless	-	-	-	42.0%
DeepSeek-V3.2-685B (Liu et al., 2025)	-	-	-	-	73.1%
Kimi K2-1TB (Team et al., 2025)	Agentless	-	-	-	65.8%
Lingma-SWE-GPT-72B (Ma et al., 2024)	SWE-SynInfer	-	-	-	28.8%
Qwen3-235B-A22B (Yang et al., 2025a)	OpenHands	-	-	-	34.4%
SWE-gym-32B (Pan et al., 2024)	OpenHands	Claude 3.5 Sonnet*	✓	RFT	20.6%
SWE-gym-32B (Reproduce)	OpenHands	DeepSeek-V3.2*	✓	RFT	48.0%
R2E-Gym-32B (Jain et al., 2025)	OpenHands	Claude 3.5 Sonnet	✓	RFT	34.4%
SWE-smith-32B (Yang et al., 2025b)	SWE-agent	Claude 3.7 Sonnet*	✓	RFT	40.2%
SWE-fixer-72B (Xie et al., 2025)	SWE-Fixer	-	✗	SFT	32.8%
SoRFT-Qwen-32B (Ma et al., 2025)	Agentless	Claude 3.5 Sonnet	✗	SFT+RL	30.8%
MCTS-Refine-32B (Wang et al., 2025)	Agentless	DeepSeek-V3	✗	SFT	32.4%
<i>SWE-Ext Models</i>					
SWE-Ext-Baseline-32B	CosAgentless	GPT-4o	✗	SFT	31.2%
SWE-Ext-Multilingual-32B	CosAgentless	GPT-4o	✗	SFT	32.6%
SWE-Ext-Completion-32B	CosAgentless	GPT-4o	✗	SFT	32.2%
Ext-Baseline-32B + SWE-Gym RFT	OpenHands	DeepSeek-V3.2*	✓	SFT+RFT	50.4%
Ext-Multilingual-32B + SWE-Gym RFT	OpenHands	DeepSeek-V3.2*	✓	SFT+RFT	49.6%
Ext-Completion-32B + SWE-Gym RFT	OpenHands	DeepSeek-V3.2*	✓	SFT+RFT	53.2%

Table 2: Resolve rates on **SWE-bench Verified**, collected from (Yang et al., 2025b) and Kimi-K2 (Team et al., 2025) technical reports. **Expert Model** indicates the LLMs that generated content during the data construction process. **Exec** indicates whether verifiable data is used for rejective sampling. All performance are results by single attempt of agent systems. * Indicates the primary data for training is *mainly* generated by the specified expert model.

Model	System	% Resolved
<i>Zero-Short Inference</i>		
GPT-4o (Hurst et al., 2024)	Agentless	9.0%
o1 (Jaech et al., 2024)	Agentless	14.0%
DeepSeek-V3-671B (Liu et al., 2024)	Agentless-Lite	11.0%
<i>SWE-Ext Models</i>		
SWE-Ext-Baseline-32B	CosAgentless	10.0%
SWE-Ext-Multilingual-32B	CosAgentless	11.5%
SWE-Ext-Completion-32B	CosAgentless	12.5%

Table 3: Resolve rates for existing solutions on **FEA-Bench Lite**, collected from (Li et al., 2025). All performance are results by single attempt of agent systems.

Mid-training benefits for downstream post-training. As depicted in Section 4, we report results under a standard rejection fine-tuning (RFT) pipeline to evaluate the quality of mid-training. When trained with the same SWE-Gym setup, models initialized from SWE-EXT mid-training consistently outperform the reproduced baseline, with gains of up to +5.2% (from 48.0% to 53.2%). Performance across different extensions remains broadly comparable after post-training, suggesting that the primary advantage lies in the augmented-data-based mid-training paradigm, which provides a stronger and more transferable foundation for subsequent fine-tuning.

Multilingual vs. Completion extension. From results on two benchmarks, multilingual extension brings more consistent improvements in issue resolving, suggesting that cross-lingual data helps the model capture generalizable reasoning strategies. In contrast, completion extension is particularly effective for feature implementation, where generating new functionality requires stronger code completion capabilities. The results highlight that different types of target complementary aspects of repository-level development 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 mid-training to improve the software engineering capabilities of LLMs.

Model	System	Component-level Localization				File-level Localization			
		Hit@1	Hit@3	MAP	MRR	Hit@1	Hit@3	MAP	MRR
Qwen2.5-Coder-32B	CoSIL	43.0	54.3	46.1	48.9	60.7	77.3	69.8	69.4
Qwen2.5-Coder-32B	CosAgentless	47.4	61.8	47.4	54.8	69.0	86.2	76.1	77.7
SWE-Ext-Baseline	CosAgentless	55.0	65.4	53.0	60.8	72.8	87.6	77.0	80.2
SWE-Ext-Multilingual	CosAgentless	57.2	68.2	55.3	63.1	75.8	90.2	79.4	82.8
SWE-Ext-Completion	CosAgentless	52.8	62.0	51.3	57.9	72.8	86.6	76.7	79.8

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

5.2 Localization Analysis

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 knowledge that strengthens 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 δ given a fixed ground-truth location ℓ^* , without providing supervision for selecting ℓ . 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

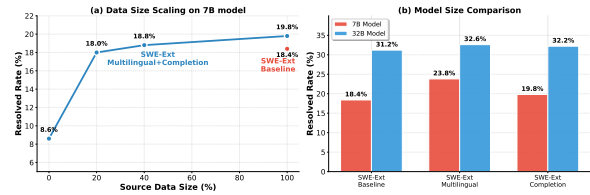


Figure 2: Scaling law analysis on SWE-bench. (a) Performance scaling with training data size follows an approximate logarithmic curve. (b) Model size comparison across three configurations demonstrates consistent improvements from 7B to 32B models.

To validate the effectiveness and robustness of our approach, we conduct experiments across both model sizes and training data volumes. We evaluate our method on Qwen2.5-Coder-7B-Instruct and systematically sample 20%, 40%, and 100% of collected data to examine scaling behaviors. The results are presented in Figure 2.

Consistent improvement across different model sizes. Our experiments reveal that although the 32B model substantially outperforms its 7B version, 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.

Extension strategies enhance data scaling. Figure 2(a) shows that applying our extension strategies yields clear logarithmic scaling patterns. 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% resolution rate, our approach lifts the scaling curve.

519	By converting previously unusable data into effective training data through multilingual and completion extensions, we improve both efficiency and the scaling potential of augmented data.	
520		
521		
522		
523	6 Related Work	
524	6.1 Repository-level Coding Task	
525	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 problems in code 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.	
532		
533		
534		
535		
536		
537		
538		
539		
540		
541		
542		
543		
544	6.2 Software Engineering Agents	
545	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. OpenHands (Wang et al., 2024) uses iterative tool calls to complete a task. 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.	
550		
551		
552		
553		
554		
555		
556		
557		
558		
559		
560		
561		
562		
563		
564		
565		
566		
567		
	6.3 Training Data for Repository-level Coding	
	Traditional instruction tuning methods enhance general code generation, such as WizardCoder (Luo et al., 2023), and WaveCoder (Yu et al., 2024). However, constructing training data for repository-level tasks differs significantly, as it requires building input-output and reasoning processes to simulate 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 to produce more verifiable tasks under constrained environments. Other methods focus on reinforcement learning. SoRFT (Ma et al., 2025) employs subtask-oriented rejection sampling and PPO.	
	Due to the diversity requirements of repositories, some works augment or synthesize 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. In relation to these efforts, while data construction for repository-level tasks predominantly focuses on Python, the proposed pipeline extends the scope by systematically gathering multilingual data of different subtasks, demonstrating that out-of-distribution data like code completion can further enhance performance.	
	7 Conclusion	
	In this work, we proposed augmented data as an explicit <i>mid-training</i> stage for repository-level coding tasks and introduced SWE-Ext , a scalable pipeline to extend such data beyond existing settings. By broadening the utilization of GitHub data along two orthogonal dimensions, multilingual coverage and an auxiliary code completion task, we expand the scope of usable augmented data without relying on execution environments. Experiments on SWE-Bench and FEA-Bench show consistent gains at the mid-training stage, with different extensions providing complementary benefits, and further indicate that augmented-data-based mid-training offers a stronger initialization for downstream fine-tuning. Our findings highlight mid-training with diverse augmented data as a simple yet effective strategy for improving repository-level code models under realistic data constraints.	

618 Limitations

619 This work focuses on augmented data as a mid-
620 training stage for repository-level coding tasks.
621 While SWE-Ext demonstrates consistent gains, sev-
622 eral limitations remain. First, due to limited API
623 resources, the scale of data synthesis is constrained,
624 and we are unable to exhaustively explore all com-
625 binations of augmentation dimensions or expand
626 to larger volumes of augmented data. Second, due
627 to restricted computational resources, our experi-
628 ments are conducted on a limited set of base mod-
629 els and training configurations, primarily at the
630 32B scale, and the observed trends may not fully
631 generalize to substantially smaller or larger mod-
632 els. Finally, although we consider multilingual and
633 completion-based extensions, other forms of aug-
634 mentation, such as refactoring-oriented supervision
635 or test generation, are not explored in this work. We
636 leave systematic investigation of larger-scale data
637 synthesis and broader augmentation strategies to
638 future work.

639 Ethics Statement

640 This work uses data derived exclusively from pub-
641 licly available open-source repositories on GitHub.
642 As such, it does not involve personal or private
643 user data and poses minimal privacy concerns. We
644 respect the original licensing terms of all source
645 repositories and only use content that is publicly
646 accessible.

647 During data construction, large language models
648 are used to assist with data augmentation and lan-
649 guage refinement. While model-generated content
650 may occasionally contain biased, insecure, or incor-
651 rect code patterns, the data is restricted to software
652 engineering contexts and does not involve sensitive
653 domains. Models trained on SWE-Ext may still
654 generate code with potential security vulnerabili-
655 ties or unintended behaviors, which is a general risk
656 for code generation systems. We therefore advise
657 users to carefully review and test generated code
658 before deployment in real-world systems.

659 Large language models were also used to assist
660 with code completion and language polishing dur-
661 ing the preparation of this manuscript. All research
662 ideas, experimental designs, and methodological
663 decisions were independently developed by the au-
664 thors, who take full responsibility for the content
665 and conclusions presented in this work.

References 666

- Anthropic. 2024. [Introducing Claude 3.5 Sonnet](#). 667
- Anthropic. 2025a. [Introducing Claude 3.7 Sonnet](#). 668
- Anthropic. 2025b. [Introducing Claude 4](#). 669
- Anthropic. 2025c. [Introducing Claude 4.5 Sonnet](#). 670
- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, and 1 others. 2021. [Program synthesis with large language models](#). *ArXiv preprint*, abs/2108.07732. 671-675
- Ramakrishna Bairi, Atharv Sonwane, Aditya Kanade, Arun Iyer, Suresh Parthasarathy, Sriram Rajamani, B Ashok, and Shashank Shet. 2024. [Codeplan: Repository-level coding using llms and planning](#). *Proceedings of the ACM on Software Engineering*, 1(FSE):675–698. 676-681
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, and 1 others. 2021. [Evaluating large language models trained on code](#). *ArXiv preprint*, abs/2107.03374. 682-687
- Zhaoling Chen, Xiangru Tang, Gangda Deng, Fang Wu, Jialong Wu, Zhiwei Jiang, Viktor Prasanna, Arman Cohan, and Xingyao Wang. 2025. [Locagent: Graph-guided llm agents for code localization](#). *arXiv preprint arXiv:2503.09089*. 688-692
- Xueying Du, Mingwei Liu, Kaixin Wang, Hanlin Wang, Junwei Liu, Yixuan Chen, Jiayi Feng, Chaofeng Sha, Xin Peng, and Yiling Lou. 2023. [Classeval: A manually-crafted benchmark for evaluating llms on class-level code generation](#). *ArXiv preprint*, abs/2308.01861. 693-698
- Binyuan Hui, Jian Yang, Zeyu Cui, Jiayi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Keming Lu, and 1 others. 2024. [Qwen2. 5-coder technical report](#). *ArXiv preprint*, abs/2409.12186. 699-703
- Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, and 1 others. 2024. [Gpt-4o system card](#). *ArXiv preprint*, abs/2410.21276. 704-708
- Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, and 1 others. 2024. [Openai o1 system card](#). *ArXiv preprint*, abs/2412.16720. 709-713
- Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. 2024. [Livecodebench: Holistic and contamination free evaluation of large language models for code](#). *ArXiv preprint*, abs/2403.07974. 714-719

720	Naman Jain, Jaskirat Singh, Manish Shetty, Liang Zheng, Koushik Sen, and Ion Stoica. 2025. R2e-gym: Procedural environments and hybrid verifiers for scaling open-weights swe agents. <i>arXiv preprint arXiv:2504.07164</i> .	776
721		777
722		
723		
724		
725	Zhonghao Jiang, Xiaoxue Ren, Meng Yan, Wei Jiang, Yong Li, and Zhongxin Liu. 2025. Cosil: Software issue localization via llm-driven code repository graph searching. <i>arXiv preprint arXiv:2503.22424</i> .	
726		
727		
728		
729	Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R Narasimhan. 2024. SWE-bench: Can language models resolve real-world github issues? In <i>The Twelfth International Conference on Learning Representations</i> .	
730		
731		
732		
733		
734		
735	Jia Li, Ge Li, Xuanming Zhang, Yihong Dong, and Zhi Jin. 2024a. Evocodebench: An evolving code generation benchmark aligned with real-world code repositories. <i>ArXiv preprint, abs/2404.00599</i> .	
736		
737		
738		
739	Jia Li, Ge Li, Yunfei Zhao, Yongmin Li, Huanyu Liu, Hao Zhu, Lecheng Wang, Kaibo Liu, Zheng Fang, Lanshen Wang, and 1 others. 2024b. Deveval: A manually-annotated code generation benchmark aligned with real-world code repositories. <i>ArXiv preprint, abs/2405.19856</i> .	
740		
741		
742		
743		
744		
745	Wei Li, Xin Zhang, Zhongxin Guo, Shaoguang Mao, Wen Luo, Guangyue Peng, Yangyu Huang, Houfeng Wang, and Scarlett Li. 2025. Fea-bench: A benchmark for evaluating repository-level code generation for feature implementation. <i>arXiv preprint arXiv:2503.06680</i> .	
746		
747		
748		
749		
750		
751	Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, and 1 others. 2024. Deepseek-v3 technical report. <i>ArXiv preprint, abs/2412.19437</i> .	
752		
753		
754		
755		
756	Aixin Liu, Aoxue Mei, Bangcai Lin, Bing Xue, Bingxuan Wang, Bingzheng Xu, Bochao Wu, Bowei Zhang, Chaofan Lin, Chen Dong, and 1 others. 2025. Deepseek-v3. 2: Pushing the frontier of open large language models. <i>arXiv preprint arXiv:2512.02556</i> .	
757		
758		
759		
760		
761	Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xubo Geng, Wenxiang Hu, Chongyang Tao, Jing Ma, Qingwei Lin, and Daxin Jiang. 2023. Wizardcoder: Empowering code large language models with evol-instruct. <i>ArXiv preprint, abs/2306.08568</i> .	
762		
763		
764		
765		
766	Yingwei Ma, Rongyu Cao, Yongchang Cao, Yue Zhang, Jue Chen, Yibo Liu, Yuchen Liu, Binhua Li, Fei Huang, and Yongbin Li. 2024. Lingma swe-gpt: An open development-process-centric language model for automated software improvement. <i>arXiv preprint arXiv:2411.00622</i> .	
767		
768		
769		
770		
771		
772	Zexiong Ma, Chao Peng, Pengfei Gao, Xiangxin Meng, Yanzhen Zou, and Bing Xie. 2025. Sorft: Issue resolving with subtask-oriented reinforced fine-tuning. <i>arXiv preprint arXiv:2502.20127</i> .	
773		
774		
775		
	Christopher D Manning. 2008. <i>Introduction to information retrieval</i> . Syngress Publishing,.	
	Jiayi Pan, Xingyao Wang, Graham Neubig, Navdeep Jaitly, Heng Ji, Alane Suhr, and Yizhe Zhang. 2024. Training software engineering agents and verifiers with swe-gym. <i>arXiv preprint arXiv:2412.21139</i> .	778
		779
		780
		781
	Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. 2024. Hybridflow: A flexible and efficient rlhf framework. <i>arXiv preprint arXiv:2409.19256</i> .	782
		783
		784
		785
		786
	Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen, Jiahao Chen, Ningxin Chen, Ruijue Chen, Yanru Chen, Yuankun Chen, Yutian Chen, and 1 others. 2025. Kimi k2: Open agentic intelligence. <i>arXiv preprint arXiv:2507.20534</i> .	787
		788
		789
		790
		791
	Xingyao Wang, Boxuan Li, Yufan Song, Frank F Xu, Xiangru Tang, Mingchen Zhuge, Jiayi Pan, Yueqi Song, Bowen Li, Jaskirat Singh, and 1 others. 2024. Openhands: An open platform for ai software developers as generalist agents. <i>arXiv preprint arXiv:2407.16741</i> .	792
		793
		794
		795
		796
		797
	Yibo Wang, Zhihao Peng, Ying Wang, Zhao Wei, Hai Yu, and Zhiliang Zhu. 2025. Mcts-refined cot: High-quality fine-tuning data for llm-based repository issue resolution. <i>arXiv preprint arXiv:2506.12728</i> .	798
		799
		800
		801
	Yuxiang Wei, Olivier Duchenne, Jade Copet, Quentin Carbonneaux, Lingming Zhang, Daniel Fried, Gabriel Synnaeve, Rishabh Singh, and Sida I Wang. 2025. Swe-rl: Advancing llm reasoning via reinforcement learning on open software evolution. <i>arXiv preprint arXiv:2502.18449</i> .	802
		803
		804
		805
		806
		807
	Chunqiu Steven Xia, Yinlin Deng, Soren Dunn, and Lingming Zhang. 2024. Agentless: Demystifying llm-based software engineering agents. <i>ArXiv preprint, abs/2407.01489</i> .	808
		809
		810
		811
	Chengxing Xie, Bowen Li, Chang Gao, He Du, Wai Lam, Difan Zou, and Kai Chen. 2025. Swe-fixer: Training open-source llms for effective and efficient github issue resolution. <i>arXiv preprint arXiv:2501.05040</i> .	812
		813
		814
		815
		816
	Frank F. Xu, Uri Alon, Graham Neubig, and Vincent Josua Hellendoorn. 2022. A systematic evaluation of large language models of code. In <i>Proceedings of the 6th ACM SIGPLAN International Symposium on Machine Programming</i> , MAPS 2022, page 1–10, New York, NY, USA. Association for Computing Machinery.	817
		818
		819
		820
		821
		822
		823
	An Yang, Anpeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025a. Qwen3 technical report. <i>arXiv preprint arXiv:2505.09388</i> .	824
		825
		826
		827
		828

829	John Yang, Carlos E Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik Narasimhan, and Ofir Press. 2024. Swe-agent: Agent-computer interfaces enable automated software engineering . <i>ArXiv preprint</i> , abs/2405.15793.	• 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.	883
830			884
831			885
832			886
833			
834	John Yang, Kilian Lieret, Carlos E Jimenez, Alexander Wettig, Kabir Khandpur, Yanzhe Zhang, Binyuan Hui, Ofir Press, Ludwig Schmidt, and Diyi Yang. 2025b. Swe-smith: Scaling data for software engineering agents. <i>arXiv preprint arXiv:2504.21798</i> .	• PRs that were not merged into the main branch, indicating that the proposed code changes were not accepted.	887
835			888
836			889
837			
838			
839	Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2023. React: Synergizing reasoning and acting in language models. In <i>International Conference on Learning Representations (ICLR)</i> .	• PRs whose patches failed to apply cleanly to the codebase.	890
840			891
841			
842		• PRs where tree-sitter encountered parsing errors in the changed files.	892
843			893
844	Zhaojian Yu, Xin Zhang, Ning Shang, Yangyu Huang, Can Xu, Yishujie Zhao, Wenxiang Hu, and Qiufeng Yin. 2024. WaveCoder: Widespread and versatile enhancement for code large language models by instruction tuning . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 5140–5153, Bangkok, Thailand. Association for Computational Linguistics.	• PRs whose changed files exclusively consisted of non-code file extensions, including: [".json", ".png", ".csv", ".txt", ".md", ".jpg", ".jpeg", ".pk1", ".yaml", ".yml", ".toml"].	894
845			895
846			896
847			897
848			898
849			
850			
851			
852			
853	Yuntong Zhang, Haifeng Ruan, Zhiyu Fan, and Abhik Roychoudhury. 2024. Autocoderover: Autonomous program improvement . In <i>Proceedings of the 33rd ACM SIGSOFT International Symposium on Software Testing and Analysis, ISSTA 2024</i> , page 1592–1604, New York, NY, USA. Association for Computing Machinery.	A.3 Filtering Training Data	899
854		After the construction of the SWE-Ext training data, further filtering is required to ensure data quality.	900
855		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:	901
856			902
857			903
858			904
859			905
860	Yanli Zhao, Andrew Gu, Rohan Varma, Liang Luo, Chien-Chin Huang, Min Xu, Less Wright, Hamid Shojanazeri, Myle Ott, Sam Shleifer, and 1 others. 2023. Pytorch fsdp: experiences on scaling fully sharded data parallel. <i>arXiv preprint arXiv:2304.11277</i> .	• 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.	906
861			907
862			908
863			909
864			910
865			
866	Terry Yue Zhuo, Minh Chien Vu, Jenny Chim, Han Hu, Wenhao Yu, Ratnadira Widayarsi, Imam Nur Bani Yusuf, Haolan Zhan, Junda He, Indraneil Paul, and 1 others. 2024. Bigcodebench: Benchmarking code generation with diverse function calls and complex instructions . <i>ArXiv preprint</i> , abs/2406.15877.	• Component Localization data: We remove samples where the ground truth edited components are not fully contained in the final output.	911
867			912
868			913
869			914
870			
871			
872	A Details of Collected Data	• 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.	915
873	A.1 Distribution of Programming Languages		916
874	The proposed augmented data for mid-training can be constructed easily across any repository using different programming languages. The distribution of programming languages of our collected data is shown in Figure 3.		917
875			918
876			919
877			
878			
879	A.2 Filtering pull requests	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.	920
880	To ensure the collection of high-quality data during the pull request (PR) scraping phase, we applied the following exclusion criteria:		921
881			922
882			923
			924
			925
			926

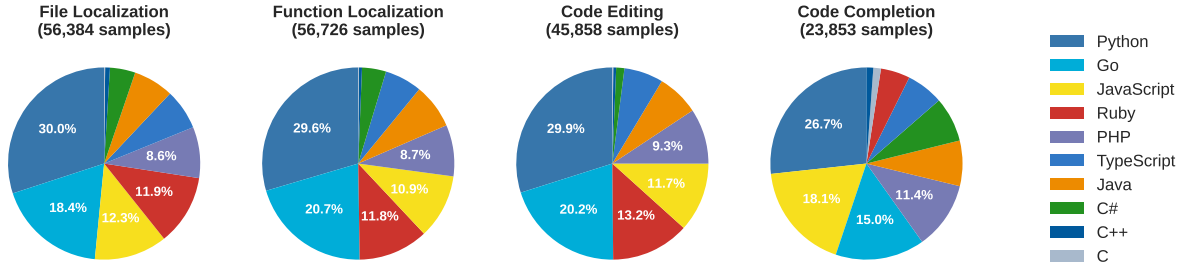


Figure 3: Language distribution of dialogues for 4 sub-tasks in SWE-Ext dataset.

B Training Settings

From our preliminary experiments, we found that using Qwen2.5-Coder 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 Qwen2.5-Coder-32B-Instruct and Qwen2.5-Coder-7B-Instruct as the backbone models for fine-tuning. For the training framework, we use the multi-turn dialogue SFT framework provided by ver1 (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 2*8 NVIDIA A100 GPUs. Training is performed in float32 precision with CPU offloading enabled. We set the maximum sequence length to 32,768. The training configuration includes a sequence parallel size of 4, a learning rate of 5×10^{-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 a 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: Initial 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 shown in Figure 4.

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, as shown in Figure 5.

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

Figure 4: Prompt for Refined File Localization in the pipeline of data construction.

```

Prompt for Initial File Localization

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
...

```

list of at most 5 files as the final localization result. The prompt template used in this phase is shown in Figure 6.

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

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 shown in 7, we define a set of tool-use capabilities that enable the model to inspect the code of a specified component.

Iterative Localization Initialization. The iterative localization process begins with an initialization step where the model is prompted to identify a single, relevant component, as shown in Figure 8. This component serves as the initial point of inspection for the codebase.

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. And the instruction is shown in Figure 9. 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.

For each component `component_retrieved` viewed during an iteration, the model is prompted to assess its relevance to the given problem, as depicted in Figure 10. 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

Figure 5: Prompt for Dependency Analysis and Reflection in the pipeline of data construction.

```
Prompt for Dependency Analysis and Reflection

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}

###

### Repository Structure ###
{structure}

###

### Files To Be Explored ###
{pre_files}

###

### Import Relations ###
{import_content}
###

Based on the import relationships, please analyze which files in the repository depend on which other files within the same repository. Ignore any libraries or modules that are imported from outside the current repository. Present the results in the following format:

file1.py -> file6.py, file7.py
file2.py ->
file3.py -> xxx/file8.py
file4.py -> file2.py, file9.py
file5.py -> xxx/file10.py, xxx/file11.py

Note: Solving the problem not only requires determining where to modify the code, but also identifying which other code to refer to in order to understand and invoke.

Based on the files listed above and the import relations, reconfirm and provide the complete names of the top 10 most likely relevant files, considering both where changes need to be made and which files are important to refer to.

Please think step by step and give you reasoning process first. Finally, provide full path and return top 10 files.

The final returned files should be separated by new lines ordered by most to least important and wrapped with ```
For example:
```
file1.py
file2.py
file3.py
file4.py
file5.py
file6.py
file7.py
file8.py
file9.py
file10.py
```

Note: file1.py indicates the top-1 file, file2.py indicates the top-2 file, and so on. Do not include test files.
```

Figure 6: Prompt for Refined File Localization in the pipeline of data construction.

```
Prompt for Refined File Localization

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 10 relevent files according to the file structure and dependency.
I will further give the internal structures of these 10 files.
Please rank them again and reflect the result according to the internal structures and dependency.

### Problem Description ###
{problem_statement}

###

### Import Relations ###
{import_content}

###

### Files To Be Ranked ###
{file_internal_structure}

###

Please think step by step and give you reasoning process first. Finally, provide full path and return top 5 files.

The returned files should be separated by new lines ordered by most to least important and wrapped with ```
For example:
```
file1.py
file2.py
file3.py
file4.py
file5.py
```

Note: file1.py indicates the top-1 file, file2.py indicates the top-2 file, and so on. Do not include test files.
```

Figure 7: System Prompt for Component localization in the pipeline of data construction.

```
System Prompt for Component localization

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.
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. *
* 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.
```

Figure 8: Prompt for the Initialization for Iterative Localization in the pipeline of data construction.

```
Prompt for the Initialization for Iterative Localization

### 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. Their structures 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.
```

Figure 9: Prompt for Iterative Component Inspection in the pipeline of data construction.

```
Prompt for Iterative Component Inspection

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}
```

Figure 10: Prompt for Judgement for Relevant Components in the pipeline of data construction.

```
Prompt for Judgement for Relevant Components

You will be presented with a repository-level coding problem with repository file structure to access the source code of
 the software.

Your task is to locate the top-5 most likely edit locations based on the problem description.

Problem Description
{problem_statement}

###

Here is a result of a function/class code retrieved by '{content}'.
Please check if the code is related to the problem and if the code should be added into context.
<code>
{component_retrieved}
</code>
Return True if the code is related to the problem and should be added into context, otherwise return False.
Since your answer will be processed automatically, please give your answer in the format as follows.
The returned content should be wrapped with ````.
```
True
```
or
```
False
```
```

Figure 11: Prompt for Final Component Localization in the pipeline of data construction.

```
Prompt for Final Component Localization

{file_internal_structure}
Based on the available information, reconfirm and provide complete names of the most likely edit locations (10 locations
at most).
Before make the final decision, please check whether the function name is correct or not, for static functions, don't add
class name.
{bug_file_list}

Please provide the complete set of locations as either a class name, a function name, or a file name.
The returned files should be separated by new lines ordered by most to least important and wrapped with ```
Since your answer will be processed automatically, please give your answer in the exapmle format as follows.
```

top1_file_fullpath.py
function: Class1.Function1

top2_file_fullpath.py
function: Function2

top3_file_fullpath.py
class: Class3

top4_file_fullpath.py
function: Class4.Function4

top5_file_fullpath.py
function: Function5

top6_file_fullpath.py
global

top7_file_fullpath.py
function: Class7.Function7
```

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 name. 'global' means the code is not in a class or function.
For example,
```
sklearn/linear_model/__init__.py
function: LinearRegression.fit
```
```

Figure 12: Prompt for Code Editing in the pipeline of data construction.

### Prompt for Code Editing

We are currently solving the following task within our repository. Here is the task description.

```
Task Description
{problem_statement}

###
```

Below are some code segments, each from a relevant file. One or more of these files may need to be edited to solve the task.

```
--- BEGIN FILE ---
{top_n_content}
--- END FILE ---
```

Please first localize the positions to edit based on the task statement, and then output the files that need to be deleted, modified or added.

- '- file' means deleting the file;
- \* file' means modifying the file;
- '+ file' means adding the file.

The file should list like below:

```

- file1.py
* file2.py
* file3.py
+ file4.py

```

To solve the task, you should then generate \*SEARCH/REPLACE\* edits.

Every \*SEARCH/REPLACE\* edit must use this format:

1. The file path
2. The start of search block: <<<<<< SEARCH
3. A contiguous chunk of lines to search for in the existing source code
4. The dividing line: =====
5. The lines to replace into the source code
6. The end of the replace block: >>>>>> REPLACE

Here is an example:

```
```python
### mathweb/flask/app.py
<<<<<< SEARCH
from flask import Flask
=====
import math
from flask import Flask
>>>>>> REPLACE
```
```

Please note that the \*SEARCH/REPLACE\* edit REQUIRES PROPER INDENTATION. If you would like to add the line ' print(x)', you must fully write that out, with all those spaces before the code!

Wrap the \*SEARCH/REPLACE\* edit in blocks ```python...```.

When multiple edits should be done, please output \*SEARCH/REPLACE\* edit one by one and give your reasoning process before each \*SEARCH/REPLACE\* block.

Figure 13: Template for Code Editing Output in the pipeline of data construction.

### Template for Code Editing Output

The plan of solving this software task:

```
{plan}
```

The files that should be edited include:

```
{files}
```

Here are my edits for code.

```
{search_replaces}
```





Figure 16: An Example of File Localization.

### An Example of File Localization: Part III

Without fill, I have to call `records[-1].process()` after the loop, assuming it's a list and not an iterable.  
</request>

```


Repository Structure ###
boltons/
 __init__.py
 cacheutils.py
 debugutils.py
 deprutils.py
 dictutils.py
 easterutils.py
 ecutils.py
 excutils.py
 fileutils.py
 formatutils.py
 funcutils.py
 gcutils.py
 ioutils.py
 iterutils.py
 jsonutils.py
 listutils.py
 mathutils.py
 mboxutils.py
 namedutils.py
 pathutils.py
 queueutils.py
 setutils.py
 socketutils.py
 statsutils.py
 strutils.py
 tableutils.py
 tbutils.py
 timeutils.py
 typeutils.py
 urlutils.py
docs/
 conf.py
misc/
 bench_omd.py
 linkify_changelog.py
 table_html_app.py
setup.py

Files To Be Explored ###
['boltons/iterutils.py', 'boltons/funcutils.py', 'boltons/listutils.py', 'boltons/typeutils.py', 'boltons/strutils.py']

Import Relations ###
file: boltons/iterutils.py
 ['import os', 'import math', 'import time', 'import codecs', 'import random', 'import itertools']

file: boltons/funcutils.py
 ['from __future__ import print_function', 'import sys', 'import re', 'import inspect', 'import functools', 'import
 itertools', 'from types import MethodType, FunctionType']

file: boltons/listutils.py
 ['from __future__ import print_function, division', 'import operator', 'from math import log as math_log', 'from
 itertools import chain, islice']

file: boltons/typeutils.py
 ['import sys', 'from collections import deque']

file: boltons/strutils.py
 ['from __future__ import print_function', 'import re', 'import sys', 'import uuid', 'import zlib', 'import string', '
 import unicodedata', 'import collections', 'from gzip import GzipFile']

Based on the import relationships, please analyze which files in the repository depend on which other files within the
same repository. Ignore any libraries or modules that are imported from outside the current repository. Present the
results in the following format:
```





















Figure 27: An Data Example of Code Editing.

### An Data Example of Code Editing: Part III

The motivation of this edit:

We need to modify the `docker_compose_run` function to include an optional parameter for port mapping. By adding a `port_mapping` parameter, we can control whether the `--service-ports` option is included in the 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 `docker-compose.yml` file.

```
```python
### fabfile.py
<<<<<<< SEARCH
))

def docker_compose_run(command_name, service="builder", user="app", no_deps=False, workdir=None):
    args = [
        'run ',
        '--rm ',
    ]
    =====
    ))

def docker_compose_run(command_name, service="builder", user="app", no_deps=False, workdir=None, port_mapping=False):
    args = [
        'run ',
        '--rm ',
    ]
    >>>>>> REPLACE
    <<<<<<< SEARCH
        if no_deps:
            args.append('--no-deps ')

        if workdir is not None:
            args.append('-w %s ' % _shell_escape(workdir))

    =====
        if no_deps:
            args.append('--no-deps ')

        if port_mapping:
            args.append('--service-ports ')

        if workdir is not None:
            args.append('-w %s ' % _shell_escape(workdir))

    >>>>>> REPLACE
    ```
```



