

SWERANK+: Multilingual, Multi-Turn Code Ranking for Software Issue Localization

Anonymous ACL submission

Abstract

Maintaining large-scale, multilingual codebases hinges on accurately localizing issues, which requires mapping natural-language error descriptions to the relevant functions that need to be modified. However, existing ranking approaches are often Python-centric and perform a single-pass search over the codebase. This work introduces SWERANK+, a framework that couples SWERANKMULTI, a cross-lingual code ranking tool, with SWERANKAGENT, an agentic search setup, for iterative, multi-turn reasoning over the code repository. SWERANKMULTI comprises a code embedding retriever and a listwise LLM reranker, and is trained using a carefully curated large-scale issue localization dataset spanning multiple popular programming languages. SWERANKAGENT adopts an agentic search loop that moves beyond single-shot localization with a memory buffer to reason and accumulate relevant localization candidates over multiple turns. Our experiments on issue localization benchmarks spanning various languages demonstrate new state-of-the-art performance with SWERANKMULTI, while SWERANKAGENT further improves localization over single-pass ranking.

1 Introduction

The maintenance of large-scale software systems constitutes a significant and ever-growing portion of the software development lifecycle. A persistent bottleneck in this process is *software issue localization* (Wong et al., 2016): the task of identifying where in a codebase a fix should be applied for a given bug report or feature request. This task requires mapping natural language descriptions, such as those found in GitHub issues, to specific code elements including files, modules, or functions. As modern code repositories grow in size and complexity to encompass thousands of files across multiple programming languages, manual localization becomes increasingly infeasible. Automating this

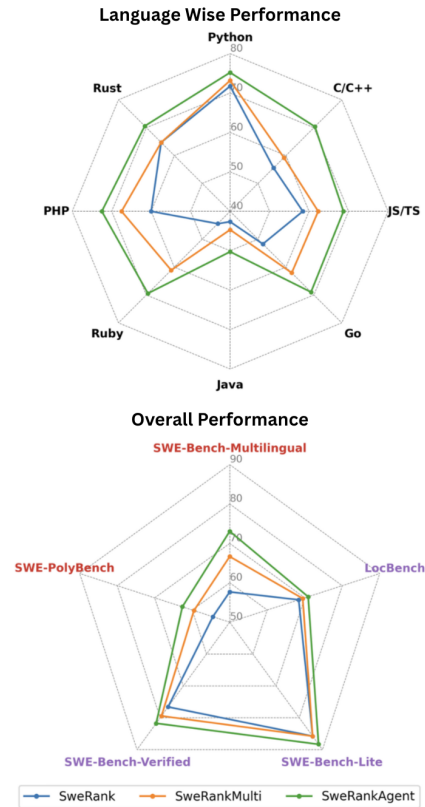


Figure 1: Comparison of function localization accuracy@10 against the SWERANK baseline. SWERANKMULTI shows significant improvement on multilingual benchmarks (in red) while maintaining strong performance on Python-specific evaluations (in violet). SWERANKAGENT further improves over single-pass ranking approaches across the board.

process can therefore accelerate issue resolution and significantly enhance developer productivity.

Recent advances in large language models (LLMs) have led to the development of sophisticated *agentic systems* (Yu et al., 2025; Chen et al., 2025) capable of navigating complex codebases. These systems operate by issuing commands to read files, search for patterns, and analyze dependencies. While powerful, such approaches often suffer from high latency and substantial computational overhead. To improve efficiency, SWER-

ANK (Reddy et al., 2025) reformulates issue localization as a *retrieve-and-rerank* problem, leveraging specialized bi-encoder retrievers and LLM rerankers trained on curated datasets (Suresh et al., 2024). This approach has achieved high precision while substantially reducing computational cost.

Despite its success, SWERANK faces two key limitations that reflect the evolving needs of modern software engineering. First, its scope is largely confined to Python-dominated repositories, whereas real-world enterprise systems are inherently *multilingual*, comprising interconnected components written in diverse programming languages. A practical localization tool must therefore generalize across languages. Second, SWERANK’s single-pass ranking design may be inadequate for complex issues that demand *iterative reasoning* or involve changes dispersed across multiple, loosely coupled functions. To overcome these challenges, we introduce SWERANK+, a comprehensive framework for multilingual, multi-turn issue localization.

Our first component, SWERANKMULTI, extends SWERANK to a multilingual setting. It includes a suite of retriever and reranker models trained on SWELOCMULTI, a newly curated, large-scale dataset containing high-quality issue–code pairs spanning multiple popular programming languages. SWERANKMULTI enables high-accuracy localization across heterogeneous repositories, marking the first such effort in this domain.

Next, we present SWERANKAGENT, a lightweight agentic framework that moves beyond single-pass localization. Instead of relying on a single retrieval pass, the agent adopts an iterative process: it begins with an initial hypothesis and progressively refines its understanding through multiple turns. Each turn allows it to gather new evidence, narrow the search space, and focus on the most relevant code regions. This multi-turn reasoning process mirrors how human developers investigate issues—starting with a broad, symptom-based exploration and converging toward the root cause. Specifically, SWERANKAGENT is motivated by three key observations: (A) **Initial searches can be misleading**—the most obvious function in a traceback is often not the true source of the bug; (B) **Context is built over turns**—early retrievals reveal important structural clues about the codebase; (C) **Iterative refinement leads to precision**—repeated reasoning helps transition from symptom-level to root-cause localization.

SWERANKAGENT integrates SWERANKMULTI

as a specialized retrieval tool within this iterative process. It can issue multiple search queries, maintain a memory buffer of intermediate results, and reason over aggregated evidence to produce more accurate localizations. This hybrid design combines the efficiency of high-quality retrievers with the depth of agentic reasoning, enabling it to solve complex localization problems that remain intractable under single-pass methods.

Through extensive experiments, we show that SWERANKMULTI establishes new state-of-the-art performance for issue localization across multiple programming languages (Rashid et al., 2025; Zan et al., 2025; Yang et al., 2025b), while maintaining competitive results on Python benchmarks (Jimenez et al., 2024; Chen et al., 2025). Furthermore, the iterative reasoning in SWERANKAGENT consistently outperforms single-pass ranking for issue localization. Figure 1 summarizes our experimental results. Our contributions are:

- We present the first framework to address issue localization in a multilingual setting.
- We introduce SWERANKMULTI, trained on 10 programming languages, achieving state-of-the-art multilingual code ranking performance for issue localization.
- We propose SWERANKAGENT, an iterative, multi-turn localization framework that further improves over single-pass ranking.

2 Related Work

2.1 Software Issue Localization

Software issue localization (fault localization) aims to identify the specific code regions responsible for software defects. Traditional approaches include spectrum-based methods (Jones and Harrold, 2005) and state-based methods (Zeller, 2002), which analyze program executions and state differences to isolate potential fault sites. Learning-based methods such as DeepFL (Li et al., 2019) extended these ideas by integrating multiple fault signals using deep representations. Recent approaches to fault localization (Torun et al., 2025; Chang et al., 2025; Yaraghi et al., 2025) have leveraged LLM-based techniques by incorporating semantic reasoning and retrieval-augmented inference, moving beyond traditional spectrum- and mutation-based metrics.

The advances in AI for software engineering have spurred the development of LLM-based agentic frameworks designed to perform complex software engineering tasks (He et al., 2025; Dong et al.,

156	2025). Recent LLM-based agentic approaches tackle	3.1 SWELOCMULTI Training Dataset	204
157	issue localization as a planning and searching	We create SWELOCMULTI, a large-scale multilin-	205
158	problem. Yang et al. (2024a) proposed LLM-based	gual dataset curated specifically for issue localiza-	206
159	framework that leverages multi-turn reasoning and	tion. While SWELOC (Reddy et al., 2025) provides	207
160	tool invocation for issue resolution and code repair.	a high-quality resource for Python, modern soft-	208
161	Chen et al. (2025) incorporates graph-guided rea-	ware systems are inherently multilingual. SWE-	209
162	soning for precise function localization, while Yu	LOCMULTI extends the data collection and filter-	210
163	et al. (2025) designs a framework specialized for	ing pipeline of SWELOC to encompass JavaScript,	211
164	bug localization using search and read actions.	Java, TypeScript, Ruby, Rust, Go, PHP, C, and	212
165	While these approaches demonstrate the power	C++. Table 4 in Appendix provides the detailed	213
166	of agentic architectures, they often involve com-	language-wise distribution of training instances.	214
167	plex interactions with the codebase, which can	Issue Collection: Following SWERANK’s	215
168	incur high latency (Chen et al., 2025; Yu et al.,	methodology, we identify popular open-source	216
169	2025). SWERANK+ introduces a distinct, hybrid	repositories on GitHub for each language, filter-	217
170	agentic model, SWERANKAGENT, that combines	ing for repositories with at least 40% code in the target	218
171	the efficiency of a specialized retrieval tool with	language, over 1,000 stars, and at least one commit	219
172	a lightweight, iterative reasoning loop. This en-	in the preceding six months. From this curated set,	220
173	ables tackling complex issues that are intractable	we extract pull requests (PRs) explicitly linked to	221
174	for single-pass ranking systems while avoiding the	GitHub issues that include test file modifications.	222
175	high overhead of complex agentic frameworks.	Contrastive Data: Each issue description serves	223
176	2.2 Multilingual Code Understanding	as a query, with modified functions in the corre-	224
177	Recent code LLMs have significantly expanded	sponding PR treated as positive examples. To en-	225
178	their code understanding capabilities across	sure the model learns to distinguish fine-grained	226
179	programming languages, with models such	semantic differences, we employ <i>consistency filter-</i>	227
180	as Qwen3-Coder (Yang et al., 2025a), and	<i>ing</i> and <i>hard-negative mining</i> techniques (Suresh	228
181	CodeGemma (Team et al., 2024) exhibiting strong	et al., 2024). Specifically, we use a pretrained em-	229
182	zero-shot transfer between languages. However,	bedding model (SWERANKEMBED-Small) to per-	230
183	the integration of multilinguality into software	form consistency filtering, retaining only examples	231
184	issue localization remains comparatively nascent.	where the positive function ranks within the top-40	232
185	Most localization systems remain monolingual, of-	semantically similar functions repository-wide. We	233
186	ten trained exclusively on Python (e.g., SWER-	then mine hard negatives, which are unmodified	234
187	ANK (Reddy et al., 2025), LocAgent (Chen et al.,	functions from the same repository that are seman-	235
188	2025)), limiting their applicability to real-world	tically similar to the query, creating challenging	236
189	software ecosystems. SWERANK+ bridges this gap	instances for the contrastive training process.	237
190	as the first localization framework explicitly trained	3.2 SWERANKEMBEDMULTI Retriever	238
191	and evaluated on multilingual repositories.	The retriever component, SWERANKEMBED-	239
192	3 SWERANKMULTI	MULTI, is a bi-encoder model that maps issue	240
193	We introduce SWERANKMULTI, a framework for	descriptions and code functions to dense vector	241
194	effective software issue localization across diverse	representations in a shared multilingual embed-	242
195	programming languages. SWERANKMULTI adapts	ding space. We initialize the model with Qwen3-	243
196	the efficient retrieve-and-rerank methodology of	Embedding (Zhang et al., 2025), a state-of-the-art	244
197	SWERANK to the multilingual setting, compris-	text embedding model, and finetune it on SWELOC-	245
198	ing two key components: the SWERANKEMBED-	MULTI using the InfoNCE contrastive loss (Oord	246
199	MULTI retriever (§3.2), which pre-selects a small	et al., 2018), identical to SWERANKEMBED’s train-	247
200	set of candidate functions from large codebases	ing objective. Given an issue description as the	248
201	most likely relevant to a given issue, and the SWER-	query, the loss function encourages the query em-	249
202	ANKLLMMULTI reranker (§3.3), which produces	beddings to be closer to that of the corresponding	250
203	a refined final ranking from the top candidates.	positive code function while pushing away from	251
		other functions in the training batch, including	252

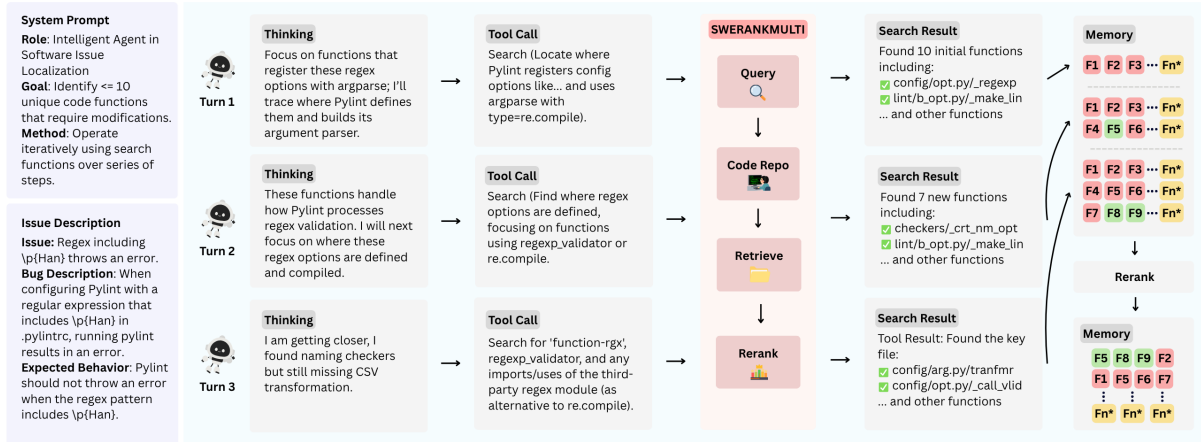


Figure 2: SWERANKAGENT interleaves *Search* actions, which retrieve candidate functions, and *Aggregate* actions, which selectively adds the most relevant functions to a persistent memory. The diagram illustrates a qualitative example of the agent solving a complex issue. The visualization details the agent’s multi-turn trajectory (turns 1-3), showing how it iterates through thinking and tool calls to SWERANKMULTI to progressively narrow down the root cause from a broad search to the target function.

mined hard negatives. During inference, all the repository functions are ranked by cosine similarity between their embeddings and the issue embedding, to identify the most probable fix locations.

3.3 SWERANKLLMMULTI Reranker

The reranker component, SWERANKLLMMULTI, refines the initial candidate list from the retriever using listwise reranking. This approach leverages the instruction-following capabilities of large language models and has proven more effective than pointwise (Zhuang et al., 2023, 2024) or pairwise ranking (Qin et al., 2024).

Training SWERANKLLMMULTI leverages the novel weakly supervised strategy introduced in SWERANK, which enables listwise reranking learning even when only a single positive example is known, without requiring a fully ordered ground-truth list. Each candidate function receives a unique identifier, and the model is trained to generate the identifier of the true positive function as its first token using standard language modeling loss. This objective effectively aligns model generation with ranking behavior, teaching it to select the most relevant function from the candidate set. Fine-tuning on the multilingual SWELOCMULTI dataset substantially enhances cross-language generalization and code ranking accuracy, setting a new state of the art for multilingual issue localization (§5.2).

4 SWERANKAGENT

While SWERANKMULTI models effectively localize relevant code functions, they operate on a single, static query—the original issue description. This

approach performs well for self-contained issues where the initial report provides sufficient signal to identify the fix location. However, many software maintenance tasks involve complex bugs or feature requests whose critical details are not fully captured in the initial description. For example, a bug report may describe only a high-level symptom, and it is only after examining the functions related to that symptom that a developer can pinpoint the downstream component responsible for the fault.

To emulate this human-like, iterative exploratory process, we introduce SWERANKAGENT, a lightweight and generic agentic framework that extends SWERANKMULTI with multi-turn search capabilities. Rather than treating issue localization as a single-shot ranking task, SWERANKAGENT performs iterative searches across multiple turns to progressively accumulate and refine localization evidence. This design enables the agent to decompose complex localization problems into a sequence of smaller, more tractable steps, effectively overcoming the limitations of single-pass retrieval.

SWERANKAGENT operates as a ReAct-style (Yao et al., 2023) reasoning agent, alternating between reasoning and action-taking. Its iterative operational loop, illustrated in Figure 2, consists of four key steps: Search, Reasoning, Reformulation, and Aggregate.

- Search with *issue description*: The agent’s primary interface with the codebase. Given an issue description as the query, this action invokes the SWERANKMULTI retriever and reranker as a tool to obtain the top- k most relevant functions from the code repository.

- 319 • Reasoning over *retrieved functions*: The
320 agent analyzes the retrieved functions to as-
321 sess their relevance to the issue and determine
322 whether further exploration is needed. This
323 step guides subsequent query refinement.
- 324 • Reformulation for *query update*: After the
325 reasoning step, the agent updates the query
326 based on the previous search and its assess-
327 ment of the retrieved functions’ relevance.
328 This enables the agent to progressively refine
329 its search and improve contextual understand-
330 ing across iterations.
- 331 • Aggregate on *candidate functions*: Across
332 turns, the agent accumulates identifiers of rel-
333 evant functions in an internal memory buffer
334 that serves as a global candidate pool. This en-
335 ables the agent to capture issues whose fixes
336 span multiple functions that may not co-occur
337 within a single retrieval. Finally, the SWER-
338 ANKLLMMULTI reranker is applied to the
339 aggregated pool to produce a globally ranked
340 list of candidate functions.

341 The agent iteratively cycles through Search,
342 Reasoning, and Reformulation until a stopping
343 condition is met—either the maximum number of it-
344 erations is reached or no new relevant functions are
345 retrieved. The accumulated candidates are then con-
346 solidated and reranked during the final Aggregate
347 step, with the top-10 functions returned as the final
348 localization output. Throughout this process, the
349 underlying LLM guides the agent’s reasoning, en-
350 abling it to refine queries, evaluate search results,
351 and dynamically adapt its strategy based on evolv-
352 ing evidence from the codebase.

353 5 Experiments

354 Our experiments are designed to address the fol-
355 lowing research questions: **RQ1**: *How effective*
356 *is training on SWELOCMULTI compared to the*
357 *Python-specific SWELOC?* and **RQ2**: *Can issue*
358 *localization benefit from multi-turn search?*. To ex-
359 amine **RQ1**, Section §5.2 reports the performance
360 of the SWERANKMULTI retriever and reranker
361 across multilingual and Python-specific issue lo-
362 calization benchmarks. To investigate **RQ2**, Sec-
363 tion §5.3 analyzes the effectiveness of multi-turn
364 issue localization using SWERANKAGENT com-
365 pared to the single-turn SWERANKMULTI.

5.1 Datasets & Metrics 366

367 Our multilingual evaluation data comprises three
368 datasets: SWE-PolyBench (Rashid et al., 2025),
369 SWE-Bench-Multilingual (Yang et al., 2025b), and
370 Multi-SWE-Bench (Zan et al., 2025). Follow-
371 ing (Suresh et al., 2024) and (Reddy et al., 2025),
372 we transform each (PR, codebase) instance pair
373 from these datasets into the localization format.
374 The PR’s corresponding github issue description
375 serves as the retrieval query. The Tree-sitter pars-
376 ing tool is employed to extract all functions from
377 the codebase, creating candidate corpus. Func-
378 tions modified within the PR are labeled are con-
379 sidered as positives. Moreover, we also consider
380 python-specific benchmarks, specifically SWE-
381 Bench-Lite (Jimenez et al., 2024), LocBench (Chen
382 et al., 2025) and SWE-Bench-Verified (Chowdhury
383 et al., 2024). We employ Accuracy at k (Acc@k)
384 for evaluation. This metric deems localization suc-
385 cessful if all relevant code locations are correctly
386 identified within the top-k results.

5.2 SweRankMulti 387

5.2.1 Setup 388

Model Training: We train the SWERANKEM-
389 BED model in two sizes: *small* and *large*. Both
390 models are trained on the SWELOCMULTI dataset
391 (§3.1), with small and large variants initial-
392 ized with the 0.6B and 8B variants of Qwen3-
393 Embedding (Zhang et al., 2025) respectively. Fol-
394 lowing SWERANK (Reddy et al., 2025), the SWER-
395 ANKLLMMULTI small and large rerankers are ini-
396 tialized with CodeRankLLM (7B) (Suresh et al.,
397 2024) and Qwen-2.5-32B-Instruct (Yang et al.,
398 2024b) respectively, and finetuned with SWELOC-
399 MULTI for listwise reranking. 400

SWERANKEMBEDMULTI Baselines: For the
401 retriever evaluation, we compare against exist-
402 ing code embedding models such SWERANKEM-
403 BED (Reddy et al., 2025), and CODERANKEM-
404 BED (Suresh et al., 2024). Since SWERANKEM-
405 BED was finetuned on GTE-Qwen2-7B-Instruct,
406 we finetune the Qwen3-Embedding models on the
407 python-specific SWELOC dataset (Reddy et al.,
408 2025) to get the SWERANKEMBEDPYTHON vari-
409 ants for better comparison. We also include
410 Gemini-Embedding (Lee et al., 2025), the current
411 top model on the MTEB leaderboard, as a general-
412 purpose closed-source baseline. 413

Model	Multilingual				Python					
	SWE-PolyBench		SWE-Bench-Multilingual		SWE-Bench-Lite		LocBench		SWE-Bench-Verified	
	Acc@5	Acc@10	Acc@5	Acc@10	Acc@5	Acc@10	Acc@5	Acc@10	Acc@5	Acc@10
OpenHands (Claude-3.5)	-	-	-	-	68.25	70.07	-	59.11	-	-
LocAgent (Claude-3.5)	-	-	-	-	73.36	77.37	-	59.29	-	-
Gemini-Embedding (unknown)	-	-	36.75	47.44	61.31	72.26	43.04	51.43	59.74	66.96
CodeRankEmbed (137M)	27.11	33.40	26.50	35.04	51.82	58.76	38.93	47.86	50.98	57.99
Qwen3-Embedding-0.6B	30.69	37.17	30.34	38.03	52.55	62.77	39.64	47.32	49.45	60.39
SWERANKEMBED-Small (137M)	35.91	42.79	33.33	44.02	63.14	74.45	51.79	58.57	59.74	68.49
SWERANKEMBEDPYTHON-Small (0.6B)	37.75	43.95	36.75	47.01	66.79	75.18	50.00	58.04	61.93	71.12
SWERANKEMBEDMULTI-Small (0.6B)	39.01	47.14	43.16	52.56	66.79	76.28	51.25	58.93	64.99	71.77
GTE-Qwen2-7B-Instruct (7B)	33.40	39.40	34.19	42.31	63.14	70.44	42.50	51.79	57.77	65.21
Qwen3-Embedding-8B	39.50	46.47	36.75	46.15	60.95	71.53	44.46	55.00	60.18	65.21
SWERANKEMBED-Large (7B)	41.92	49.18	39.74	50.85	71.90	82.12	55.18	63.21	65.65	74.18
SWERANKEMBEDPYTHON-Large (8B)	44.24	51.98	46.15	55.98	73.72	83.94	55.71	65.00	66.96	75.05
SWERANKEMBEDMULTI-Large (8B)	46.56	53.73	50.43	62.39	77.37	86.86	56.43	65.36	68.71	76.37

Table 1: Function localization performance of different retrievers in comparable sizes across benchmarks with various programming languages. Our SWERANKEMBEDMULTI models considerably improve performance on both multilingual and python-specific benchmarks.

Model	Multilingual				Python					
	SWE-PolyBench		SWE-Bench-Multilingual		SWE-Bench-Lite		LocBench		SWE-Bench-Verified	
	Acc@5	Acc@10	Acc@5	Acc@10	Acc@5	Acc@10	Acc@5	Acc@10	Acc@5	Acc@10
SweRankEmbedMulti-Small (0.6B)	39.01	47.14	43.16	52.56	66.79	76.28	51.25	58.93	64.99	71.77
+ CodeRankLLM (7B)	42.40	51.11	48.72	58.55	72.99	80.29	55.54	63.93	66.52	75.93
+ Qwen3-Instruct-8B	43.66	51.98	50.85	59.83	72.27	80.29	56.96	64.46	68.49	75.27
+ SweRankLLM-Small (7B)	51.21	56.92	55.13	63.68	77.37	85.04	63.39	69.64	73.30	77.68
+ SweRankLLMMulti-Small (7B)	53.15	59.54	56.41	66.67	80.29	85.77	63.04	69.46	73.30	79.43
SweRankEmbedMulti-Large (8B)	46.56	53.73	50.43	62.39	77.37	86.86	56.43	65.36	68.71	76.37
+ GPT-4.1	55.66	61.76	62.39	70.94	79.93	88.69	65.89	70.89	76.81	81.62
+ SweRankLLM-Large (32B)	56.73	62.73	57.69	70.51	83.58	89.42	64.64	71.25	75.49	80.96
+ SweRankLLMMulti-Large (32B)	58.28	63.21	64.10	71.37	85.40	89.78	66.43	71.96	78.12	81.18

Table 2: Function localization performance of different rerankers on multilingual and python-specific benchmarks.

SWERANKLLMMULTI Baselines: For the reranker evaluation, we compare against listwise reranker models such as SWERANKLLM (Reddy et al., 2025), CODERANKLLM (Suresh et al., 2024), and the zero-shot Qwen3-Instruct-8B model (Yang et al., 2025a) along with GPT-4.1.

5.2.2 Results

Table 1 shows the function localization performance different embedding models when compared at similar size ranges. We see that SWERANKEMBEDMULTI achieves SOTA performance at both size variants, considerably improving upon the existing SWERANKEMBED. Moreover, we see that a better base embedding model (SWERANKEMBED vs SWERANKEMBEDPYTHON) and training with multilingual data (SWERANKEMBEDPYTHON vs SWERANKEMBEDMULTI) both contribute to improvements in performance. Interestingly, SWERANKEMBEDMULTI even improves over SWERANKEMBEDPYTHON on python-specific benchmarks, despite the latter being trained on the python-exclusive SWELOC. This could be attributed to better gen-

eralization performance from including more languages in the training data.

Table 2 shows the function localization performance of the reranker when used with retrievers of different sizes. We see that SWERANKLLMMULTI reranker consistently improve localization performance over the SWERANKLLM models, with the large variant even outperforming GPT-4.1.

Language-wise Performance: Figure 3 shows performance of the retrievers separately for each of the languages in SWE-PolyBench and SWE-Bench-Multilingual. We observe that both the small and large variants of SWERANKEMBED consistently improve performance on most languages, compared to SWERANKEMBEDPYTHON.

5.3 SweRankAgent

5.3.1 Setup

SWERANKMULTI tool: We equip the SWERANKAGENT framework with a search tool that, when given with a query in the form of an issue description, returns the top-10 localized functions within the codebase. The SWERANKMULTI-

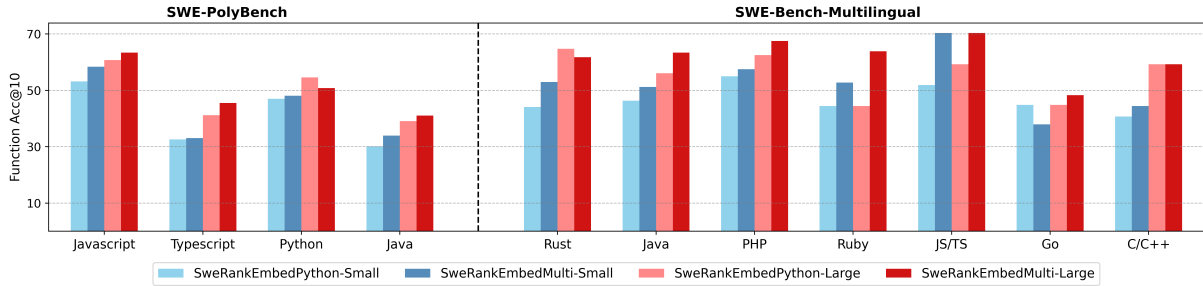


Figure 3: Language-wise function localization accuracy for SWE-PolyBench and SWE-Bench-Multilingual.

Approach	Multilingual		Python		
	SWE-PolyBench	SWE-Bench-Multilingual	SWE-Bench-Lite	LocBench	SWE-Bench-Verified
SWERANKMULTI (Single-Query)	59.54	66.67	85.77	69.46	79.43
Reformulate (Single-Turn, Multi-Query)	58.08 (-1.46)	64.96 (-1.71)	85.77	69.86 (+0.40)	78.99 (-0.44)
SWERANKAGENT (Multi-Turn)	62.63 (+3.09)	73.08 (+6.41)	88.32 (+2.55)	70.91 (+1.45)	81.74 (+2.31)

Table 3: Table comparing SWERANKAGENT’s multi-turn function localization performance (Acc@10) against the single-turn SWERANKMULTI tool in addition to a multi-query reformulation approach.

EMBED retriever first obtains the top-100 results which are then passed to the SWERANKMULTI-LLM reranker. For efficiency considerations, the codebase function embeddings for the retriever are pre-cached. Given compute considerations, we only use the small variants of the retriever and reranker for the SWERANKMULTI tool.

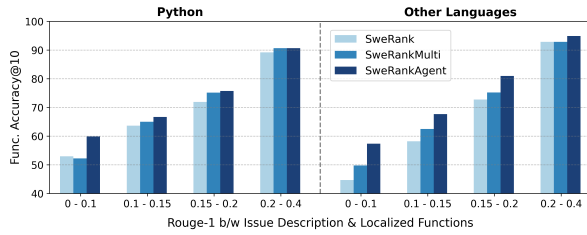
Baselines: To investigate the performance of a multi-query setup, we design a *Reformulate* baseline that does query-reformulation to generate multiple queries from the original issue description. In total, we use the original issue description as the query plus five reformulated queries generated by GPT-5. For each query, results from the SWERANKMULTI tool are collected, and the aggregated output is then again passed to the SWERANKMULTI-LLM reranker to obtain final top-10 function localization result. Moreover, we also use the SWERANKMULTI tool as a single-query baseline.

5.3.2 Results

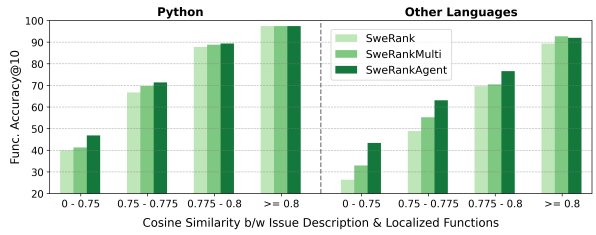
Table 3 shows the function localization performance of SWERANKAGENT using GPT-5 as the underlying LLM. Firstly, we observe that the multi-query *Reformulate* baseline yields lower performance than the single-query baseline. This suggests that directly using reformulated queries introduces additional noise, highlighting the inherent challenges in improving coverage of localization candidates with multiple queries. In contrast, SWERANKAGENT consistently improves localization performance across all datasets, with more than a 3-point boost on SWE-PolyBench and over a 6-point boost on SWE-Bench-Multilingual.

Qualitative Example: Figure 2 illustrates the advantage of multi-turn issue localization through a qualitative example where iterative search succeeds while single-pass ranking fails. At first glance, the issue appears straightforward: the pylint *re.error* traceback points to *re.compile*. Single-pass ranking would likely point at *_regex_validator*, the immediate crash site. However, the real bug lies upstream, in how a comma-separated list of regular expressions is parsed before validation. Modifying *_regex_validator* would thus be incorrect; the fix belongs in the function handling the comma-separated-value (CSV) input. SWERANKAGENT’s multi-turn trajectory below demonstrates how iterative refinement enables correct localization.

- **Turn 1: Broad Exploration.** The agent begins with a broad query to locate where configuration options are defined and validated. This identifies relevant configuration modules in *arguments_manager.py* and the symptom function *_regex_validator*, but not the root cause. A single-pass approach would likely stop here.
- **Turn 2: Contextual Refinement.** Based on the initial results, the agent infers that the issue relates to how naming options are validated, refining its search to focus on functions involved in regular expression validation. This uncovers key functions such as *_regex_csv_validator*, revealing the CSV-handling pathway.
- **Turn 3: Pinpointing the Target.** The results point to a transformer function that processes CSV inputs prior to validation. A focused search on transformer-validator interactions returns the correct target: *_regex_csv_transformer*.



(a) Performance breakdown by lexical overlap.



(b) Performance breakdown by semantic overlap.

Figure 4: Comparison of SWERANK, SWERANKMULTI, and SWERANKAGENT performance across varying levels of difficulty. The plots show Function Accuracy@10 separately for Python and other languages, broken down by (a) lexical overlap using ROUGE-1 scores, and (b) semantic overlap using cosine similarity.

5.3.3 Analysis by Issue Complexity

Aggregate performance metrics can mask how localization methods behave under varying levels of problem difficulty. We thereby analyze performance across two orthogonal dimensions of complexity: (1) the degree of lexical and semantic overlap between the issue description and target code, and (2) the number of ground-truth localized functions. We primarily aim to investigate whether the agentic search process is particularly beneficial for harder instances with multiple target functions or requiring matching beyond surface-level similarity.

Performance by Overlap: We bucket instances based on lexical overlap (ROUGE-1) and semantic overlap (cosine similarity) between the issue description and the ground-truth localized functions, with results shown in Figure 4. We can see that performance generally improves with higher overlap, indicating that keyword- or semantically-aligned issues are easier to localize. Nonetheless, SWERANKAGENT demonstrates clear advantages in low-overlap settings. In buckets with minimal lexical or semantic overlap, where naive keyword matching is insufficient, the agent consistently shows bigger improvements compared to buckets with more overlap. The agent’s ability to reformulate queries and reason over intermediate results allows it to bridge the gap when the issue description lacks specific keywords present in the target functions. Hence, multi-turn reasoning is crucial for traversing the “semantic gap” in harder instances where the failure description is distant from the root cause.

Performance by # of Target Functions: We stratify test instances by the number of functions modified in the ground-truth patch, which serves as a proxy for issue complexity. As seen in Figure 5, localization accuracy degrades as the number of target localization functions increases, reflecting the increased difficulty of multi-function

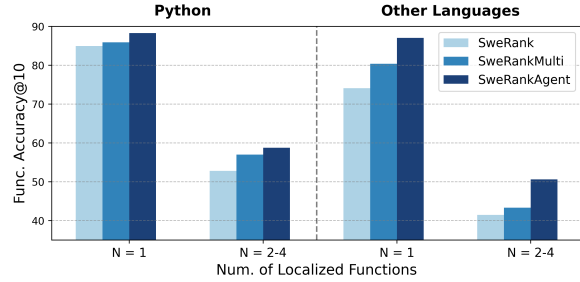


Figure 5: Function Accuracy@10 breakdown by number of target localization functions.

bugs and feature requests. While the multilingual training of SWERANKMULTI does considerably improve performance for other languages, SWERANKMULTI particularly shines for more complex multi-function localization settings.

6 Conclusion

In this work, we introduced SWERANK+, a comprehensive software issue localization framework. Our contributions are twofold. First, we developed SWERANKMULTI, a multilingual retrieve-and-rerank system trained on the newly curated SWELOCMULTI dataset, that establishes new state-of-the-art performance on multilingual benchmarks while maintaining strong efficacy on Python-specific tasks. Second, we proposed SWERANKAGENT, a lightweight agentic framework that employs iterative, multi-turn reasoning to refine localization candidates, which consistently outperforms single-pass ranking. Our analysis highlights that this agentic approach is particularly effective for complex issues characterized by low overlap between issue description and target functions or those requiring modifications across multiple functions. Our work emphasizes a hybrid design that bridges the gap between efficient retrieval systems and computationally heavy agentic frameworks, demonstrating that lightweight, iterative reasoning can effectively solve complex issue localization.

593 Limitations

594 Our pull-request selection process follows Jimenez
595 et al. (2024) to exclusively include only those that
596 add test cases, as this gives a reliable way to ensure
597 that the issue was resolved. However, in practical
598 software development contexts, numerous bug-
599 fixing instances do not involve the addition of test
600 cases. Consequently, this filtering strategy may
601 exclude a significant portion of bug-fix commits.
602 Further, we leave for future work to explore in-
603 tegrating SWERANK+ into end-to-end automated
604 program repair.

605 References

606 Jianming Chang, Xin Zhou, Lulu Wang, David Lo, and
607 Bixin Li. 2025. Bridging bug localization and is-
608 sue fixing: A hierarchical localization framework
609 leveraging large language models. *arXiv preprint*
610 *arXiv:2502.15292*.

611 Zhaoling Chen, Xiangru Tang, Gangda Deng, Fang
612 Wu, Jialong Wu, Zhiwei Jiang, Viktor Prasanna, Ar-
613 man Cohan, and Xingyao Wang. 2025. Locagent:
614 Graph-guided llm agents for code localization. *arXiv*
615 *preprint arXiv:2503.09089*.

616 Neil Chowdhury, James Aung, Chan Jun Shern, Oliver
617 Jaffe, Dane Sherburn, Giulio Starace, Evan Mays,
618 Rachel Dias, Marwan Aljubeh, Mia Glaese, Carlos E.
619 Jimenez, John Yang, Leyton Ho, Tejal Patwardhan,
620 Kevin Liu, and Aleksander Madry. 2024. [Introducing](#)
621 [SWE-bench verified](#).

622 Yihong Dong, Xue Jiang, Jiaru Qian, Tian Wang, Kechi
623 Zhang, Zhi Jin, and Ge Li. 2025. [A survey on](#)
624 [code generation with llm-based agents](#). *Preprint*,
625 *arXiv:2508.00083*.

626 Junda He, Christoph Treude, and David Lo. 2025. Llm-
627 based multi-agent systems for software engineering:
628 Literature review, vision, and the road ahead. *ACM*
629 *Transactions on Software Engineering and Method-*
630 *ology*, 34(5):1–30.

631 Neel Jain, Ping yeh Chiang, Yuxin Wen, John Kirchen-
632 bauer, Hong-Min Chu, Gowthami Somepalli, Brian R.
633 Bartoldson, Bhavya Kailkhura, Avi Schwarzschild,
634 Aniruddha Saha, Micah Goldblum, Jonas Geiping,
635 and Tom Goldstein. 2023. [Neftune: Noisy em-](#)
636 [beddings improve instruction finetuning](#). *Preprint*,
637 *arXiv:2310.05914*.

638 Carlos E Jimenez, John Yang, Alexander Wettig,
639 Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R
640 Narasimhan. 2024. [SWE-bench: Can language mod-](#)
641 [els resolve real-world github issues?](#) In *The Twelfth*
642 *International Conference on Learning Representa-*
643 *tions*.

James A Jones and Mary Jean Harrold. 2005. Em-
644 pirical evaluation of the tarantula automatic fault-
645 localization technique. In *Proceedings of the 20th*
646 *IEEE/ACM international Conference on Automated*
647 *software engineering*, pages 273–282. 648

Jinhyuk Lee, Feiyang Chen, Sahil Dua, Daniel
649 Cer, Madhuri Shanbhogue, Iftekhar Naim, Gus-
650 tavo Hernández Ábrego, Zhe Li, Kaifeng Chen, Hen-
651 rique Schechter Vera, and 1 others. 2025. Gemini
652 embedding: Generalizable embeddings from gemini.
653 *arXiv preprint arXiv:2503.07891*. 654

Xia Li, Wei Li, Yuqun Zhang, and Lingming Zhang.
655 2019. Deepfl: Integrating multiple fault diagnosis di-
656 mensions for deep fault localization. In *Proceedings*
657 *of the 28th ACM SIGSOFT international symposium*
658 *on software testing and analysis*, pages 169–180. 659

Zach Nussbaum, John X. Morris, Brandon Duderstadt,
660 and Andriy Mulyar. 2024. [Nomic embed: Training a](#)
661 [reproducible long context text embedder](#). *Preprint*,
662 *arXiv:2402.01613*. 663

Aaron van den Oord, Yazhe Li, and Oriol Vinyals. 2018.
664 Representation learning with contrastive predictive
665 coding. In *Advances in Neural Information Process-*
666 *ing Systems*, pages 10203–10213. 667

Ronak Pradeep, Sahel Sharifymoghaddam, and Jimmy
668 Lin. 2023. [Rankvicuna: Zero-shot listwise document](#)
669 [reranking with open-source large language models](#).
670 *Preprint*, *arXiv:2309.15088*. 671

Zhen Qin, Rolf Jagerman, Kai Hui, Honglei Zhuang,
672 Junru Wu, Le Yan, Jiaming Shen, Tianqi Liu, Jialu
673 Liu, Donald Metzler, and 1 others. 2024. Large lan-
674 guage models are effective text rankers with pairwise
675 ranking prompting. In *Findings of the Association*
676 *for Computational Linguistics: NAACL 2024*, pages
677 1504–1518. 678

Muhammad Shihab Rashid, Christian Bock, Zhuang
679 Yuan, Alexander Buccholz, Tim Esler, Simon
680 Valentin, Luca Franceschi, Martin Wistuba,
681 Prabhu Teja Sivaprasad, Woo Jung Kim, Anoop
682 Deoras, Giovanni Zappella, and Laurent Callot. 2025.
683 [Swe-polybench: A multi-language benchmark for](#)
684 [repository level evaluation of coding agents](#). *ArXiv*,
685 *abs/2504.08703*. 686

Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and
687 Yuxiong He. 2020. Deepspeed: System optimiza-
688 tions enable training deep learning models with over
689 100 billion parameters. In *Proceedings of the 26th*
690 *ACM SIGKDD International Conference on Knowl-*
691 *edge Discovery & Data Mining*, pages 3505–3506. 692

Revanth Gangi Reddy, Tarun Suresh, JaeHyeok Doo,
693 Ye Liu, Xuan Phi Nguyen, Yingbo Zhou, Semih
694 Yavuz, Caiming Xiong, Heng Ji, and Shafiq Joty.
695 2025. Swerank: Software issue localization with
696 code ranking. *arXiv preprint arXiv:2505.07849*. 697

698	Tarun Suresh, Revanth Gangi Reddy, Yifei Xu, Zach Nussbaum, Andriy Mulyar, Brandon Duderstadt, and Heng Ji. 2024. Cornstack: High-quality contrastive data for better code ranking. <i>arXiv preprint arXiv:2412.01007</i> .	755
699		756
700		757
701		758
702		759
703	CodeGemma Team, Heri Zhao, Jeffrey Hui, Joshua Howland, Nam Nguyen, Siqi Zuo, Andrea Hu, Christopher A. Choquette-Choo, Jingyue Shen, Joe Kelley, Kshitij Bansal, Luke Vilnis, Mateo Wirth, Paul Michel, Peter Choy, Pratik Joshi, Ravin Kumar, Sarmad Hashmi, Shubham Agrawal, and 8 others. 2024. Codegemma: Open code models based on gemma . <i>Preprint</i> , arXiv:2406.11409.	760
704		761
705		762
706		763
707		764
708		765
709		766
710		767
711	Utku Boran Torun, Mehmet Taha Demircan, Mahmut Furkan Gön, and Eray Tüzün. 2025. Past, present, and future of bug tracking in the generative ai era. <i>arXiv preprint arXiv:2510.08005</i> .	768
712		769
713		770
714		771
715	W Eric Wong, Ruizhi Gao, Yihao Li, Rui Abreu, and Franz Wotawa. 2016. A survey on software fault localization. <i>IEEE Transactions on Software Engineering</i> , 42(8):707–740.	772
716		773
717		774
718		775
719	An Yang, Anfeng 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> .	776
720		777
721		778
722		779
723		780
724	John Yang, Carlos E Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik Narasimhan, and Ofir Press. 2024a. Swe-agent: Agent-computer interfaces enable automated software engineering. <i>arXiv preprint arXiv:2405.15793</i> .	781
725		782
726		783
727		784
728		785
729	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>Preprint</i> , arXiv:2504.21798.	
730		
731		
732		
733		
734	Qwen An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxin Yang, Jingren Zhou, Junyang Lin, and 25 others. 2024b. Qwen2.5 technical report . <i>ArXiv</i> , abs/2412.15115.	
735		
736		
737		
738		
739		
740		
741	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> .	
742		
743		
744		
745		
746	Ahmadreza Saboor Yaraghi, Golnaz Gharachorlu, Sakina Fatima, Lionel C Briand, Ruiyuan Wan, and Ruifeng Gao. 2025. Black-box test code fault localization driven by large language models and execution estimation. <i>arXiv preprint arXiv:2506.19045</i> .	
747		
748		
749		
750		
751	Zhongming Yu, Hejia Zhang, Yujie Zhao, Hanxian Huang, Matrix Yao, Ke Ding, and Jishen Zhao. 2025. Orcaloca: An llm agent framework for software issue localization. <i>arXiv preprint arXiv:2502.00350</i> .	
752		
753		
754		
	Daoguang Zan, Zhirong Huang, Wei Liu, Hanwu Chen, Linhao Zhang, Shulin Xin, Lu Chen, Qi Liu, Xiaojian Zhong, Aoyan Li, Siyao Liu, Yongsheng Xiao, Liangqiang Chen, Yuyu Zhang, Jing Su, Tianyu Liu, Rui Long, Kai Shen, and Liang Xiang. 2025. Multi-swe-bench: A multilingual benchmark for issue resolving . <i>ArXiv</i> , abs/2504.02605.	
	Andreas Zeller. 2002. Isolating cause-effect chains from computer programs. <i>ACM SIGSOFT Software Engineering Notes</i> , 27(6):1–10.	
	Yanzhao Zhang, Mingxin Li, Dingkun Long, Xin Zhang, Huan Lin, Baosong Yang, Pengjun Xie, An Yang, Dayiheng Liu, Junyang Lin, and 1 others. 2025. Qwen3 embedding: Advancing text embedding and reranking through foundation models. <i>arXiv preprint arXiv:2506.05176</i> .	
	Honglei Zhuang, Zhen Qin, Kai Hui, Junru Wu, Le Yan, Xuanhui Wang, and Michael Bendersky. 2024. Beyond yes and no: Improving zero-shot llm rankers via scoring fine-grained relevance labels. In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 2: Short Papers)</i> , pages 358–370.	
	Honglei Zhuang, Zhen Qin, Rolf Jagerman, Kai Hui, Ji Ma, Jing Lu, Jianmo Ni, Xuanhui Wang, and Michael Bendersky. 2023. Rankt5: Fine-tuning t5 for text ranking with ranking losses. In <i>Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval</i> , pages 2308–2313.	

A Training Details

SWERANKEMBEDMULTI: Our data filtering, negative mining, and model finetuning are implemented using the contrastors package (Nussbaum et al., 2024). The SWERANKEMBEDMULTI-SMALL encoder is finetuned for two epochs with a learning rate of $2e-5$, a batch size of 64 and 15 hard negatives per example. The SWERANKEMBEDMULTI-LARGE encoder is finetuned for 1 epoch with a learning rate of $8e-6$, a batch size of 64 and 7 hard negatives per example. Both models are finetuned using 8 GH200 GPUs.

SWERANKLLMMULTI: For the LLM reranker training, we trained for one epoch with a global batch size of 128, an initial learning rate of $5e-6$ with 50 warmup steps, cosine learning rate scheduler, bfloat16 precision, and noisy embeddings (Jain et al., 2023) with a noise scale $\alpha = 5$. For efficient long-context, multi-gpu training, we used DeepSpeed (Rasley et al., 2020) ZeRO stage 3 with 16 GH200 GPUs. To prevent the positional bias from affecting the reranker and ensure model robustness (Pradeep et al., 2023), we shuffle the order of candidate codes for each training example.

Language	# Repos	# PRs	# Instances
JavaScript	104	1513	4254
Java	130	5518	19239
TypeScript	129	3882	11410
Ruby	308	4244	9048
Rust	269	5879	22255
Go	114	2985	11242
PHP	206	4591	16608
C	74	1023	4013
C++	278	2359	7621
Python	2448	24285	49973
Total	4060	56279	155663

Table 4: Distribution of repositories, pull requests (PRs), and training instances across different programming languages in the SWELOCMULTI dataset.

SWERANKLLMMULTI Prompt

System Prompt

You are CodeRanker, an intelligent code reviewer that can analyze GitHub issues and rank code functions based on their relevance to containing the faults causing the GitHub issue.

User Prompt

I will provide you with 10 code functions, each indicated by a numerical identifier []. Rank the code functions based on their relevance to containing the faults causing the following GitHub issue: <Issue Description>

Code Functions

[1]: <Function 1>
[2]: <Function 2>
...
[10]: <Function 10>

Response Format

All the code functions should be included and listed using identifiers, in descending order of relevance. The output format should be [] > [], e.g., [2] > [1]. Only respond with the ranking results, do not give any explanation.

SWERANKAGENT Prompt

System Prompt

You are an intelligent assistant specializing in software issue localization. Your primary goal is to identify the functions from a given codebase that are most likely to require modification to fix a given issue. You must operate by iteratively using the search tool.

Rules and Guidelines

- Iterative Search:** perform sequential 'search' calls; number of rounds is configurable.
- Review & Reflect after each 'search':** use results to inform your next query. Avoid duplicate queries.
- Explain & Reformulate:** explain relevance for new functions, then reason about how to refine next query.
- Termination:** once coverage is sufficient or rounds are done, call 'finish' with up to 10 functions.

Available Tools

"name": "search", "description": "Searches the codebase for functions relevant to the query. Returns a list of candidate functions found based on the description of the issue passed to the tool.", "parameters": "issue_description"

"name": "finish", "description": "Call this tool when you are confident you have identified all the top relevant functions.", "parameters": null

Expected Response Format

THOUGHT: Summarize what you just learned from the latest search results. For EACH newly added function, provide a brief relevance explanation describing why it may relate to the issue description.

REFORMULATION: Explain how you will adjust the next search query to improve coverage/diversity.

ACTION: {"name": "...", "arguments": { ... }}

User Prompt

< Github Issue Description >

786

787

788

789

790

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813