

RepoAttention: Focusing on Relevant Context for Repository-Level Code Completion

Anonymous ACL submission

Abstract

Repository-level code completion, which leverages the entire codebase to generate suggestions, is crucial for enhancing developer productivity. While retrieval-augmented generation (RAG) with code large language models (code LLMs) has become the standard approach for this task, existing methods still suffer from several problems: they struggle with uninformative query based on incomplete code, overlook the impact of prompt structure and snippet ordering, and retrieve textually similar but functionally different code. These issues collectively introduce noise and lead to ineffective context for code LLMs. To address these problems, we introduce RepoAttention, a training-free RAG-based framework for repository-level code completion. It improves retrieval and completion performance through dual-path query refinement, relevance-aware reranking of retrieved snippets, and a dynamic relevance-guided attention mechanism. Experiments on CCEval and RepoEval show that RepoAttention surpasses state-of-the-art methods by 23.9% in Exact Match accuracy and generalizes well across multiple code LLMs and programming languages.

1 Introduction

With the rapid development of code large language models (code LLMs) (Chen et al., 2021; Guo et al., 2024; Li et al., 2023; Roziere et al., 2023), code completion has become a pivotal capability integrated into modern programming environments, aiming to enhance developer productivity (Izadi et al., 2024; Mastropaolo et al., 2023; Nam et al., 2024; Tang et al., 2023; Xu et al., 2022). While contemporary models excel at local code synthesis within a single file, repository-level code completion remains challenging, as it requires leveraging cross-file dependencies and structural patterns across an entire codebase. This difficulty stems from the limited context window of code LLMs

and the fact that naively incorporating entire repository is both computationally inefficient and prone to introducing irrelevant noise that interferes with completion. Therefore, recent research has adopted retrieval-augmented generation (RAG) frameworks (Yu et al., 2024; Liang et al., 2024; Liu et al., 2024a; Wang et al., 2024c; Zhang et al., 2023, 2025a,b), which leverage the incomplete code within the current file as a query to retrieve relevant code snippets from the whole repository, providing cross-file context. These retrieved code snippets are then concatenated with the incomplete code to form the prompt for code LLMs. While these approaches have showed encouraging results, several problems remain in both the retrieval and completion stages.

P1 Deficient Query Representation. In most prior methods (Wu et al., 2025; Zhang et al., 2025a; Phan et al., 2024), the query to retrieve code snippets is simply the incomplete code within the current file, and the prompt is constructed by directly concatenating the query and the retrieved code snippets. RepoCoder (Zhang et al., 2023) uses an iterative retrieval method that builds new queries by combining intermediate completions with incomplete code, but incorrect generations may propagate to later iterations and multiple retrieval rounds greatly reduce efficiency. The absence of the correct code completion and the presence of noise in the query reduce the likelihood of retrieving truly helpful code snippets, and interfere with both the retrieval and completion stages.

P2 Prompt Order Sensitivity. Prior works (Jin et al., 2024; Jiang et al., 2024; Liu et al., 2023) in natural language processing tasks have demonstrated that position of retrieved context in prompt significantly impacts model performance. However, existing code completion research (Wang et al., 2024c; Zhang et al., 2025a) often overlooks the influence of code snippet ordering within the prompt on generation outcomes. Even when retrieved snippets contain useful information, the ar-

084 arbitrary or suboptimal arrangement of these snippets
085 can mislead the code LLMs, resulting in unsatis-
086 factory completion quality.

087 **P3 Similar Snippet Confusion.** Many prior re-
088 trieval approaches tend to retrieve multiple code
089 snippets that are highly similar to each other in
090 textual form but differ in functionality (Wang et al.,
091 2024c; Liang et al., 2024). These similar code
092 snippets can introduce confusion for code LLMs,
093 potentially leading to inaccurate or semantically
094 incorrect completions. The inability to effectively
095 mitigate the interference from superficial textual
096 similarity while preserving true functional rele-
097 vance remains a major challenge in guiding the
098 completion phase towards correct solutions.

099 To address these problems, we introduce Re-
100 poAttention, a novel RAG-based framework for
101 repository-level code completion. First, we adopt a
102 dual-path query refinement strategy that combines
103 the raw incomplete code with a distilled, semanti-
104 cally enriched version to capture both lexical pat-
105 terns and completion intent. We further improve
106 retrieval quality by using diverse candidate com-
107 pletions to expand the queries for dense retrieval
108 (addressing P1). Second, we develop a relevance-
109 aware reranking mechanism that mimics how pro-
110 grammers prioritize auxiliary code according to
111 its relationship with the completion task (address-
112 ing P2). Finally, during completion, we employ
113 a dynamic relevance-guided attention mechanism
114 to focus on the most pertinent code snippets while
115 reducing interference from similar but functionally
116 different code snippets (addressing P3).

117 We evaluate RepoAttention through extensive ex-
118 periments using several representative code LLMs
119 on CCEval (Ding et al., 2023) and RepoEval
120 (Zhang et al., 2023) benchmarks. Experiment re-
121 sults show that RepoAttention achieves 23.9% im-
122 provement of Exact Match metric compared with
123 state-of-the-art (SOTA) methods. Furthermore,
124 RepoAttention demonstrates high generalizability,
125 showing effectiveness across various code LLMs
126 and programming languages.

127 The main contributions of this paper are:

- 128 • We introduce the RepoAttention framework for
129 repository-level code completion. By adopting
130 a dual-path query refinement strategy that com-
131 bines semantic enrichment with diverse hypo-
132 theoretical completions, RepoAttention effectively
133 addresses the problem of deficient query repre-
134 sentation and bridges the semantic gap between
135 incomplete code and the completion target.

- 136 • We develop a relevance-aware reranking mecha-
137 nism that orders retrieved code snippets accord-
138 ing to their functional relevance to the completion
139 target. This approach reduces prompt order sen-
140 sitivity and leads to more accurate completions.
- 141 • We propose a dynamic relevance-guided atten-
142 tion mechanism for the completion stage. This
143 mechanism enables RepoAttention to robustly fo-
144 cus on the most pertinent code context, while mit-
145 igating the impact of textually similar but func-
146 tionally different code snippets.
- 147 • A comprehensive evaluation shows that RepoAt-
148 tention outperforms the SOTA approaches and
149 consistently generalizes across different code
150 LLMs and programming languages. To facili-
151 tate further research, we have made our source
152 code publicly available¹.

153 2 Related Works

154 **Code Large Language Models.** LLMs have sig-
155 nificantly advanced software engineering by au-
156 tomating tasks like code completion (Wang et al.,
157 2023a). To better capture the properties of pro-
158 gramming languages, recent work has focused on
159 code LLMs. These models are either fine-tuned
160 from general LLMs (Touvron et al., 2023) on code
161 datasets (Wang et al., 2024d; Luo et al., 2023) or
162 trained from scratch on vast amounts of source code
163 (Nijkamp et al., 2023; Wang et al., 2023b). Re-
164 cent mainstream code LLMs, such as Code Llama
165 series (Roziere et al., 2023), Code-Qwen series
166 (Hui et al., 2024), and DeepSeek-Coder series
167 (Guo et al., 2024), are typically trained on massive,
168 diverse code corpora spanning multiple programming
169 languages, primarily utilizing self-supervised ob-
170 jectives like next-token prediction, sometimes aug-
171 mented with fill-in-the-middle objectives (Bavarian
172 et al., 2022) or repository-level context modeling
173 to improve code completion performance.

174 **Repository-level Code Completion.** Repository-
175 level code completion has attracted growing at-
176 tention in recent years. RLCoder (Wang et al.,
177 2024c) introduces reinforcement learning to opti-
178 mize RAG, training retrievers without labeled data.
179 STALL++ (Liu et al., 2024a), A³-CodGen (Liao
180 et al., 2024), RepoFuse (Liang et al., 2024), and
181 CodePlan (Bairi et al., 2024) leverage static anal-
182 ysis to obtain relevant retrieved code. DRACO
183 (Cheng et al., 2024), RepoHyper (Phan et al., 2024),

¹<https://anonymous.4open.science/r/RepoAttention-38EE/>

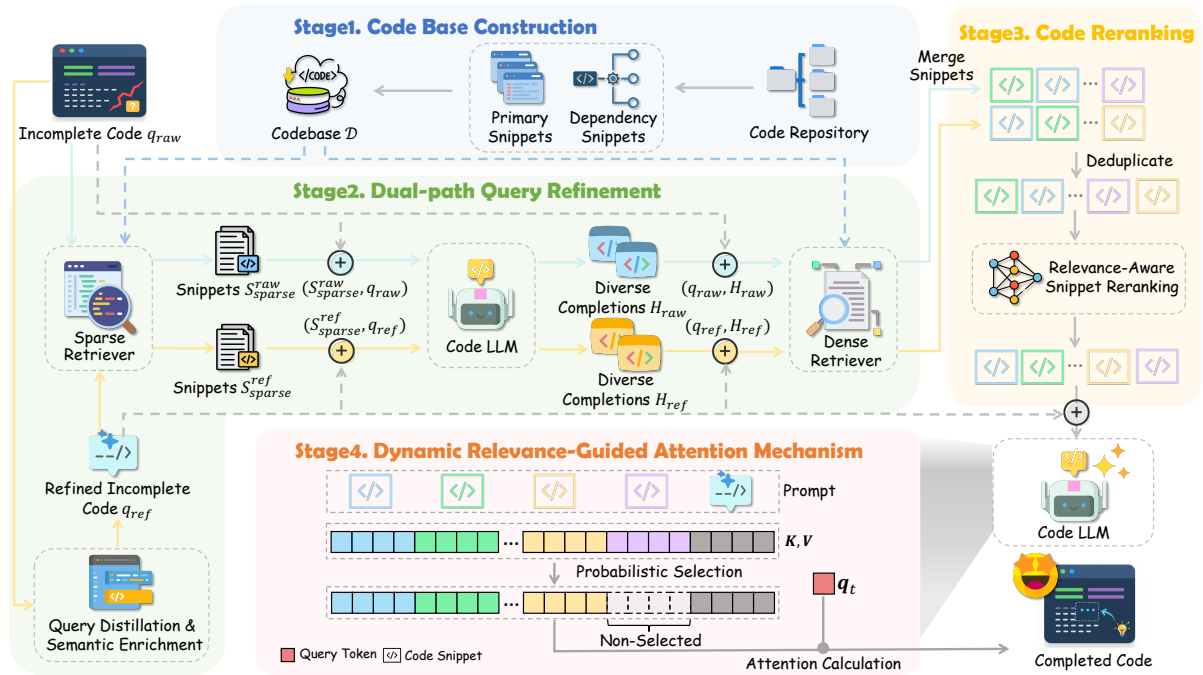


Figure 1: Overview of RepoAttention.

and CoCoMIC (Ding et al., 2022) utilize dependency analysis to retrieve relevant code snippets. RRG (Gao et al., 2024) and HCP (Zhang et al., 2025a) refine the raw retrieved code through transformation or pruning. CodeRAG (Zhang et al., 2025b) optimizes code selection by integrating multi-path retrieval and preference-aligned reranking. RepoGenReflex (Wang et al., 2024b), De-Hallucinator (Eghbali and Pradel, 2024), and RepoCoder (Zhang et al., 2023) employ iterative retrieval and generation using feedback. Repoforner (Wu et al., 2024) improves efficiency through selective retrieval, while ProCC (Tan et al., 2024) combines multiple retrieval results to better utilize context. ToolGen (Wang et al., 2024a) and CodeAgent (Zhang et al., 2024) leverage external tools like agents to tackle code completion tasks. GraphCoder (Liu et al., 2024b) and ContextModule (Guan et al., 2024) construct graphs to assist retrieval but incur high construction costs and face challenges in formulating accurate graph queries. Despite their effectiveness, these approaches do not fully solve the problems described in Section 1.

3 Preliminary Studies

To further investigate the challenges of repository-level code completion, we conduct a series of preliminary experiments under the setup detailed in Section 5.1. The results highlight two key factors that motivate our proposed approach: (1) the content of the query is essential for identifying useful

information, given that enriching the incomplete code with predicted sequences can help bridge the semantic gap and facilitate the retrieval of more relevant cross-file context; (2) the ordering of retrieved code snippets in the prompt plays a vital role in model performance, as inconsistent or random placement leads to a noticeable decrease in completion accuracy. These findings highlight the importance of both query refinement and prompt structuring. Detailed experimental results and comprehensive analysis are provided in Appendix A.

4 Methodology

As shown in Figure 1, RepoAttention involves four parts: codebase construction, dual-path query refinement, relevance-aware snippet reranking and dynamic relevance-guided attention mechanism.

4.1 Codebase Construction

Constructing the repository-level codebase involves transforming raw source code into structured snippets to facilitate effective retrieval. Traditional fixed-window slicing approaches (Zhang et al., 2023; Guan et al., 2024) often compromise code continuity and structural integrity, while dependency-parsing methods (Liu et al., 2024a,b) capture limited context and struggle with complex repository dependencies. To address these limitations, we construct a retrieval codebase composed of primary and dependency code snippets.

For primary code snippets, we align with typical

programming practices by initially segmenting the repository code based on blank lines (Wang et al., 2024c). Subsequently, we merge consecutive short code fragments into primary code snippets, ensuring each snippet maintains an appropriate length and improved semantic coherence.

For dependency code snippets, we utilize tree-sitter² to parse import statements from in-file context to identify intra-repository dependencies. We then access the referenced files to extract only the signatures of classes, methods, and functions. This process captures critical semantic details regarding class inheritance and API interplay (Zhang et al., 2025a). Distinct from simple listing these signatures, we adopt a hybrid aggregation strategy: the dependency context for each imported class is encapsulated as a standalone snippet, whereas the signatures of all imported functions and methods are consolidated into a single unified snippet. The final codebase is constructed by integrating these primary and dependency code snippets.

4.2 Dual-Path Query Refinement

Repository-level code completion relies on retrieving snippets that are semantically aligned with the intended ground truth. However, directly using incomplete code as the retrieval query introduces two issues. First, incomplete code often contains noise that weakens semantic signals and misleads both the retriever and the code LLM. Second, the retrieval process is inherently biased toward the query rather than the actual completion target, resulting in snippets similar to the query but irrelevant to the completion target. To overcome these issues, we introduce a dual-path query refinement strategy.

As shown in Appendix D, for the first path, we retain the original incomplete code c_{in} as the raw query q_{raw} to preserve exact lexical patterns. For the second path, we construct a refined query q_{ref} by applying query distillation and semantic enrichment to c_{in} . Specifically, we remove the following elements from c_{in} that are likely to introduce noise:

- **Import statements of unused modules** within the context near the completion point. These statements increase prompt length without contributing useful information, potentially distracting the retriever or model.
- **Print and logging statements**, which typically serve debugging purposes and are unrelated to the functional semantics required for completion.

- **Irrelevant variable definitions and statements** that are not semantically or data-flow related to the completion point. Only statements with strong relevance should be preserved.

To replicate the semantic clarity provided by programmer comments, we generate descriptive comment when it is missing. We utilize a lightweight LLM to analyze surrounding identifiers (e.g., function names) and generate a concise natural language comment describing the intended behavior. This comment is inserted before the completion point. By applying these modifications, q_{ref} becomes structurally concise while highlighting information vital for code completion. We use Qwen3-0.6B (Yang et al., 2025) for comment generation.

To capture both explicit keyword matches and implicit semantic relationships, we first use BM25 as the sparse retriever for both queries to independently retrieve code snippets from the codebase \mathcal{D} . This yields two initial sets of code snippets, S_{sparse}^{raw} and S_{sparse}^{ref} . These snippets may include both primary and dependency code snippets, offering information that ranges from basic code attributes to semantic details of class structures and API interplay. We construct two distinct prompts by prepending the sparse retrieval results to their respective queries and input them into the code LLM \mathcal{M} to generate hypothetical completions. For each path, we sample m diverse completions:

$$H_{raw} = \{h_1, h_2, \dots, h_m\} \sim P_{\mathcal{M}}(\cdot | S_{sparse}^{raw} \oplus q_{raw}), \quad (1)$$

$$H_{ref} = \{h'_1, h'_2, \dots, h'_m\} \sim P_{\mathcal{M}}(\cdot | S_{sparse}^{ref} \oplus q_{ref}), \quad (2)$$

where \oplus denotes concatenation. Although these completions may not be correct, they can serve as candidate completions for understanding the likely functionality and form of the correct completion.

By appending the generated completions to their corresponding queries, we form expanded queries that explicitly represent the potential completion target. Both expanded queries are then input into a dense retriever to independently retrieve the top- k most relevant code snippets from the codebase \mathcal{D} . This process can be formalized as:

$$S_{dense}^{raw} = \text{DenseRetrieve}(q_{raw} \oplus H_{raw}, \mathcal{D}) = \{s_1, s_2, \dots, s_k\}, \quad (3)$$

$$S_{dense}^{ref} = \text{DenseRetrieve}(q_{ref} \oplus H_{ref}, \mathcal{D}) = \{s'_1, s'_2, \dots, s'_k\}. \quad (4)$$

²<https://tree-sitter.github.io/tree-sitter>

4.3 Relevance-Aware Snippet Reranking

According to the experimental results presented in Appendix A.4, the position of retrieved code snippets within the prompt has a non-negligible impact on the quality of code completion. For most code LLMs, placing cross-file code snippets that are more relevant to the query either earlier or later in the prompt tends to improve completion accuracy (Jiang et al., 2024; Liu et al., 2023). This motivates the need for explicitly ordering retrieved snippets based on their relevance to the query.

Although the context retrieved in the previous stage attempts to include all potentially relevant snippets, it does not effectively assess the degree of relevance between each snippet and the query. To address this, we perform a reranking process to evaluate the relevance and arrange them in an optimal order within the prompt.

RepoAttention introduces a simple yet effective relevance-aware reranking strategy, inspired by the typical workflow of human programmers during code completion. When completing code, programmers primarily rely on the local context near the completion point to infer the intended functionality. Code snippets that are more similar to this local context in terms of functionality and semantics are typically more helpful.

This intuition suggests that focusing on code lines near the completion point is better than treating the entire query uniformly. Concretely, we extract the t lines before the completion point of q_{ref} as the *key code* segment c_k , and use a lightweight reranker to compute a relevance score r_i between c_k and each retrieved code snippet. Snippets with higher relevance scores are placed earlier in the prompt. This process is formalized as:

$$C_v = \text{Rerank}(c_k, S_{dense}^{raw} \cup S_{dense}^{ref}), \quad (5)$$

where C_v denotes the final retrieved valuable context, and r_i denotes the relevance score of the i -th snippet. We use Qwen3-Reranker-0.6B (Zhang et al., 2025c) as the reranker.

4.4 Dynamic Relevance-Guided Attention Mechanism

With the integration of the methods introduced above, the retrieved valuable context C_v is generally sufficient to guide the code LLM toward generating the correct completion. Nevertheless, in a non-trivial number of cases, C_v contains code snippets that are textually similar yet functionally

```

# protor1/agents/dueling.py
...
def update(self):
    if not self.memory.ready():
        return
    self.optimizer.zero_grad()
    self.replace_target_network()
    states, actions, rewards, states_, dones = self.sample_memory()
    indices = np.arange(len(states))

# protor1/agents/base.py
def sample_memory(self, mode='uniform'):
    memory_batch = self.memory.sample_buffer(mode=mode)
    memory_tensors = convert_arrays_to_tensors(memory_batch, self.device)
    return memory_tensors

...

# protor1/agents/ppo.py
for epoch in range(self.n_epochs):
    for batch in batches:
        indices, states, actions, rewards, states_, dones, old_probs = \
            convert_arrays_to_tensors(batch, device=self.device)
        if self.action_type == 'continuous':
            alpha, beta = self.actor(states)
            _, new_probs, entropy = self.policy(alpha, beta,
                                                old_action=actions,
                                                with_entropy=True)
            last_dim = int(len(new_probs.shape) - 1)
            prob_ratio = T.exp(
                new_probs.sum(last_dim, keepdims=True) -
                old_probs.sum(last_dim, keepdims=True))
            # a = adv[indices]

# protor1/agents/dqn.py
from protor1.agents.base import Agent
import numpy as np
import torch as T

class DQNAgent(Agent):
    ...
    def update(self):
        if not self.memory.ready():
            return
        self.optimizer.zero_grad()
        self.replace_target_network()
        if self.prioritized:
            sample_idx, states, actions, rewards, states_, dones, weights = \
                self.
memory.sample_buffer(mode='prioritized')
sample_memory(mode='prioritized')

```

Figure 2: An example where similar code snippets within the retrieved valuable context introduce noise and result in incorrect code completion.

different, which can introduce noise and lead to incorrect completions, as exemplified in Figure 2.

To mitigate this issue, it is insufficient to simply feed the retrieved results directly into the code LLM. Instead, we propose that the model dynamically adjust its reliance on different parts of the retrieved context during the generation phase.

To this end, we introduce Dynamic Relevance-Guided Attention (DRGA), which extends the standard Transformer model by dynamically selecting relevant code snippets for attention computation. By replacing the traditional attention mechanism (Vaswani et al., 2017) with DRGA, the code LLM can focus more on the relevant snippets and ignore distractors, leading to more efficient and effective reasoning during code generation.

Unlike traditional attention mechanisms, where each query token attends to the entire context, DRGA restricts each query token to attend only to a subset of selected keys and values:

$$\text{DRGA}(q, \mathbf{K}, \mathbf{V}) = \text{Softmax}\left(q\mathbf{K}[E]^\top\right)\mathbf{V}[E], \quad (6)$$

where $q \in \mathbb{R}^{1 \times d}$ denotes a single query token,

and $\mathbf{K}, \mathbf{V} \in \mathbb{R}^{N \times d}$ represent the keys and values over N tokens, respectively. $E \subseteq [N]$ is the set of selected key and values. The dimension d corresponds to the size of a single attention head. Multi-head attention extends this formulation by computing multiple such heads in parallel and concatenating their outputs.

In our code completion task, the prompt provided to the code LLM represents the entire context, consisting of the refined query q_{ref} and the retrieved valuable context C_v , which includes multiple code snippets. We enable each query token to attend only to a subset of tokens from q_{ref} and the most relevant code snippets in C_v as follows:

$$E = \left(\bigcup_{r_i \geq I} \{T_i, T_q\} \right) \cup \left(\bigcup_{r_i < I} \text{Select}(T_i) \right), \quad (7)$$

where T_i denotes the tokens of the i -th code snippet in C_v , T_q denotes the tokens of q_{ref} , and r_i is the relevance score of the i -th code snippet and c_k (defined in Section 4.3), which quantifies the relevance between q_{ref} and the i -th code snippet. I denotes the relevance threshold. The operation $\text{Select}(T_i)$ refers to the process of proportional mapping of r_i to a corresponding probability, and then decides whether to retain T_i based on that probability.

5 Experiment

5.1 Experimental Setup

Benchmarks. We utilize two widely adopted benchmarks for evaluation: CCEval (Ding et al., 2023) and RepoEval (Zhang et al., 2023). CCEval is a benchmark for repository-level code completion across multiple programming languages, reflecting real-world development scenarios. RepoEval evaluates line-level, API-level, and function-level completion on full-scale repositories. In our experiments, we use the Python and Java subsets of CCEval, which contain 2,665 and 2,139 test cases respectively, as these languages are widely used and provide sufficient diversity. For RepoEval, we focus on line-level and API-level tasks, each with 1,600 test cases, since our work targets code completion rather than function-level code generation.

Metrics. The completion results were evaluated using two widely adopted metrics: *Exact Match (EM)* and *Edit Similarity (ES)* (Levenshtein et al., 1966). Specifically, EM measures whether the predicted code completion exactly matches the correct

code without any differences. ES evaluates the textual similarity between the predicted and correct code completions based on the edit distance.

Baselines. To evaluate the effectiveness of RepoAttention, we compared it with the simple RAG method (Parvez et al., 2021) and the SOTA repository-level code completion approaches:

- **Simple RAG.** In repository-level code completion, simple RAG is a basic retrieval and generation baseline. It uses the left context of incomplete code as the query, retrieves relevant snippets from the repository with BM25, and concatenates them with the in-file context to form the prompt.
- **HCP (Zhang et al., 2025a).** HCP constructs high-quality code completion prompts by modeling the code repository at the function level, retaining topological dependencies between files while pruning irrelevant content from the retrieved code, without requiring model training.
- **RLCoder (Wang et al., 2024c).** RLCoder employs a reinforcement learning framework to train a retriever without labeled data, iteratively evaluating and retrieving useful code snippets based on the perplexity of the target code.
- **CodeRAG (Zhang et al., 2025b).** CodeRAG utilizes a multi-path retrieval strategy followed by a preference-aligned reranking mechanism to select necessary code snippets.

Experimental Details. To prevent potential data leakage, we select four popular open-source code LLMs that were released prior to the collection of benchmarks for our experiments: StarCoderBase-7B (Li et al., 2023), CodeLlama-7B (Roziere et al., 2023), DeepSeekCoder-1.3B (Guo et al., 2024), and DeepSeekCoder-6.7B. We employ the popular dense retriever RLRetriever (Wang et al., 2024c), and set the `max_new_tokens` to 64, `max_context_tokens` to 2048 for all models. We set m to 3 in Section 4.2. Following the vLLM (Kwon et al., 2023) documentation, we set the temperature to 0.8 and `top-p` to 0.95 by default during sampling. All experiments were conducted on NVIDIA A100 GPUs.

5.2 Evaluation Results

Effectiveness of RepoAttention. Table 1 presents the repository-level code completion performance of RepoAttention in comparison to the baselines. Note that both HCP and CodeRAG are primarily designed for Python programs on CCEval, with

Table 1: Repository-level code completion performance comparison.

Code LLM	Method	CCEval (Python)		CCEval (Java)		RepoEval (Line)		RepoEval (API)	
		EM	ES	EM	ES	EM	ES	EM	ES
StarCoderBase-7B	Simple RAG	14.78	66.30	16.36	66.41	44.19	65.48	32.69	55.91
	HCP	19.87	68.75	/	/	/	/	/	/
	CodeRAG	24.01	60.89	/	/	/	/	/	/
	RLCoder	25.89	72.12	25.81	68.91	47.75	68.48	35.06	58.09
	RepoAttention	31.52 ^{↑21.7%}	75.53 ^{↑4.7%}	28.63 ^{↑10.9%}	70.21 ^{↑1.9%}	49.11 ^{↑2.8%}	68.89 ^{↑0.6%}	36.68 ^{↑4.6%}	60.24 ^{↑3.7%}
CodeLlama-7B	Simple RAG	15.76	66.30	16.78	65.37	44.25	65.65	37.00	63.43
	HCP	20.06	68.18	/	/	/	/	/	/
	CodeRAG	23.11	60.44	/	/	/	/	/	/
	RLCoder	26.57	71.43	26.23	66.37	46.62	67.88	37.94	64.32
	RepoAttention	32.13 ^{↑20.9%}	75.38 ^{↑5.5%}	29.56 ^{↑12.7%}	68.79 ^{↑3.6%}	48.08 ^{↑3.1%}	68.25 ^{↑0.5%}	40.12 ^{↑5.7%}	66.17 ^{↑2.9%}
DeepSeekCoder-1.3B	Simple RAG	14.18	65.12	13.70	61.63	41.31	63.60	35.31	62.07
	HCP	17.95	66.77	/	/	/	/	/	/
	CodeRAG	21.24	58.78	/	/	/	/	/	/
	RLCoder	23.98	70.44	20.80	63.39	44.12	66.46	36.06	62.72
	RepoAttention	29.71 ^{↑23.9%}	74.24 ^{↑5.4%}	24.28 ^{↑16.7%}	65.13 ^{↑2.7%}	45.78 ^{↑3.8%}	67.04 ^{↑0.9%}	37.87 ^{↑5.0%}	64.81 ^{↑3.3%}
DeepSeekCoder-6.7B	Simple RAG	18.31	68.38	17.48	65.23	45.94	66.68	38.25	64.99
	HCP	22.62	70.49	/	/	/	/	/	/
	CodeRAG	25.85	62.70	/	/	/	/	/	/
	RLCoder	30.24	73.58	26.13	66.11	48.75	69.42	39.88	66.24
	RepoAttention	36.78 ^{↑21.6%}	78.19 ^{↑6.3%}	29.09 ^{↑11.3%}	68.78 ^{↑4.0%}	50.63 ^{↑3.9%}	70.48 ^{↑1.5%}	42.64 ^{↑6.9%}	68.34 ^{↑3.2%}

Table 2: Ablation study results on CCEval and RepoEval.

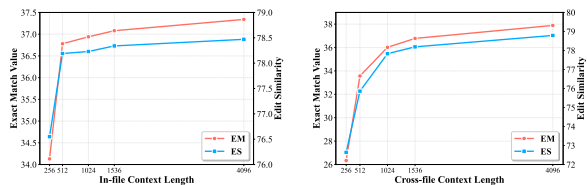
Method	CCEval (Python)		CCEval (Java)		RepoEval (Line)		RepoEval (API)	
	EM	ES	EM	ES	EM	ES	EM	ES
RepoAttention	36.78	78.19	29.09	68.78	50.63	70.48	42.64	68.34
w/o QR	35.13 ^{↓4.5%}	77.42 ^{↓1.0%}	28.31 ^{↓2.7%}	67.43 ^{↓2.0%}	50.26 ^{↓0.7%}	70.28 ^{↓0.3%}	41.62 ^{↓2.4%}	67.57 ^{↓1.1%}
w/o RO	35.88 ^{↓2.4%}	77.76 ^{↓0.5%}	28.91 ^{↓0.6%}	68.48 ^{↓0.4%}	50.18 ^{↓0.9%}	70.25 ^{↓0.3%}	41.39 ^{↓2.9%}	67.41 ^{↓1.4%}
w/o DA	34.56 ^{↓6.0%}	77.16 ^{↓1.3%}	28.14 ^{↓3.3%}	67.19 ^{↓2.3%}	50.06 ^{↓1.1%}	70.19 ^{↓0.4%}	41.02 ^{↓3.8%}	67.13 ^{↓1.8%}

no available adaptation for Java or repository-level evaluation on RepoEval. It can be observed that RepoAttention consistently outperforms all baseline methods across different code LLMs and benchmarks, demonstrating the robustness of RepoAttention for repository-level code completion. The largest improvement is observed with DeepSeekCoder-1.3B, where RepoAttention achieves a 23.9% gain in EM on CCEval Python compared to RLCoder.

Compared with HCP, which is also a training-free method, RepoAttention shows substantially stronger performance. With DeepSeekCoder-6.7B, RepoAttention improves EM by 62.6% on CCEval Python and yields consistent ES improvements across all benchmarks. These results indicate that RepoAttention can more effectively focus on relevant code snippets and enhance completion quality without additional training.

Ablation Study. To understand the contributions of each component in RepoAttention, we conduct ablation studies with DeepSeekCoder-6.7B. The results are shown in Table 2. “w/o QR” means

adopting a single-path retrieval setting that uses the incomplete code as the query, without query refinement or candidate completion generation. “w/o RO” means not performing reranking operations on the retrieved code snippets. “w/o DA” means using the traditional attention mechanism instead of our dynamic relevance-guided attention mechanism. We observe that removing any component leads to performance degradation, showing that each component contributes to the effectiveness of RepoAttention. Notably, the performance drops the most significantly for “w/o DA”, with EM decreasing by up to 6.0% on CCEval Python, indicating that the dynamic relevance-guided attention mechanism plays a critical role in RepoAttention. This highlights the importance of prioritizing relevant code snippets while discarding irrelevant or noisy information that may interfere with completion quality. Furthermore, the performance decrease in “w/o QR” and “w/o RO” suggests that both query refinement and snippet reranking consistently improve retrieval quality. These findings collectively confirm that this multi-stage approach is vital for



(a) Performance on different in-file context lengths. (b) Performance on different cross-file context lengths.

Figure 3: Performance under different context lengths.

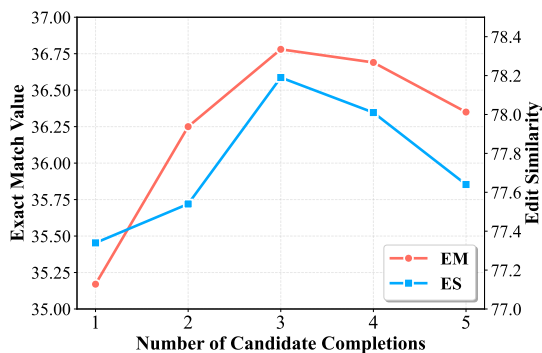


Figure 4: RepoAttention’s performance by the number of candidate completions.

effective repository-level code completion.

Impact of Context Length. We conduct experiments on CCEval Python with DeepSeekCoder-6.7B. The prompt for code LLMs includes both in-file and cross-file context, with default lengths of 512 and 1536 tokens. We vary the length of one type of context while keeping the other fixed at its default value. Figures 3(a) and 3(b) show RepoAttention’s performance as the context length increases from 256 to 4096 tokens.

Increasing either type of context consistently improves performance across all metrics, but the gains diminish as the context grows. For in-file context, expanding the length from 256 to 512 tokens yields a clear improvement, with EM increasing by 7.8% and ES by 2.1%, while further increases bring smaller gains. Cross-file context follows a similar trend, where early expansion is highly beneficial but additional tokens are less likely to introduce new relevant information, leading to a performance plateau. Overall, longer cross-file context is more effective than longer in-file context, as it provides more diverse and completion-relevant snippets. Based on these results, we recommend using around 512 tokens for in-file context and 1536 tokens for cross-file context in practice.

Impact of Candidate Completion Count. Figure 4 shows how RepoAttention’s performance varies with different numbers of candidate completions. As the number of candidate completions increases, RepoAttention initially achieves consis-

tent improvements across all metrics. However, when this number exceeds three, performance starts to degrade. This suggests that when the number of candidate completions becomes excessive, there is an increased likelihood that some candidate completions are irrelevant to the correct completion, thereby introducing noise into the retrieval process and consequently degrading the quality of the final completion. Therefore, we adopt three candidate completions as the default setting.

Computational Cost. Table 7 reports the average cost per case on CCEval Python. RepoAttention achieves the lowest latency, outperforming other baselines. Although CodeRAG incurs higher latency due to its multi-path retrieval and probing process, RepoAttention remains efficient despite using a dual-path refinement strategy. These results indicate that RepoAttention achieves a good balance between retrieval quality and computational cost, making it suitable for practical code completion. More details are provided in Appendix B.

Case Study. We present a case study to illustrate the effectiveness of RepoAttention. Figure 5 shows an example where the model is required to complete the function `pred_topk_with_label` in the `CleanSpeechDetector` class. RepoAttention first produces a refined query that supports useful candidate sampling. It then reranks the retrieved snippets, placing the most relevant one containing the target function at the top of the context, which enables correct completion. In contrast, baselines such as CodeRAG, HCP, and RLCoder retrieve similar but irrelevant snippets, resulting in incorrect outputs. More details are provided in Appendix C.

6 Conclusion

In this paper, we presented RepoAttention, a robust framework designed to enhance repository-level code completion. Our approach addresses the limitations of existing RAG methods through three key innovations: a dual-path query refinement strategy that bridges the semantic gap between incomplete code and retrieval targets, a relevance-aware reranking mechanism that optimizes the arrangement of context, and a dynamic relevance-guided attention mechanism that allows the model to focus on the most pertinent information while ignoring noise. Comprehensive experiments demonstrate that RepoAttention significantly outperforms SOTA methods and exhibits strong generalization across different programming languages and code LLMs.

635 Limitations

636 In this section, we discuss the current limitations
637 in RepoAttention.

638 **Hyperparameter Sensitivity.** As indicated in
639 our experimental setup and analysis, our frame-
640 work involves several hyperparameters that require
641 tuning, such as the sampling temperature and the
642 number of retrieved code snippets. Identifying the
643 optimal combination of these parameters for differ-
644 ent programming languages or base models can be
645 laborious and computationally expensive.

646 **Evaluation Metrics and Human Utility.** Our
647 evaluation mainly relies on automatic metrics such
648 as EM and ES. While these metrics are widely used,
649 they may not fully reflect developers’ subjective
650 judgments of code correctness or usefulness in real
651 scenarios. Future work could include human evalu-
652 ation, user studies, or downstream task-based eval-
653 uations to better assess the practical effectiveness
654 of RepoAttention.

655 References

656 Ramakrishna Bairi, Atharv Sonwane, Aditya Kanade,
657 Arun Iyer, Suresh Parthasarathy, Sriram Rajamani,
658 B Ashok, and Shashank Shet. 2024. Codeplan:
659 Repository-level coding using llms and planning.
660 *Proceedings of the ACM on Software Engineering*,
661 1(FSE):675–698.

662 Mohammad Bavarian, Heewoo Jun, Nikolas Tezak,
663 John Schulman, Christine McLeavey, Jerry Tworek,
664 and Mark Chen. 2022. Efficient training of lan-
665 guage models to fill in the middle. *arXiv preprint*
666 *arXiv:2207.14255*.

667 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan,
668 Henrique Ponde De Oliveira Pinto, Jared Kaplan,
669 Harri Edwards, Yuri Burda, Nicholas Joseph, Greg
670 Brockman, and 1 others. 2021. Evaluating large
671 language models trained on code. *arXiv preprint*
672 *arXiv:2107.03374*.

673 Wei Cheng, Yuhan Wu, and Wei Hu. 2024. Dataflow-
674 guided retrieval augmentation for repository-level
675 code completion. *arXiv preprint arXiv:2405.19782*.

676 Yangruibo Ding, Zijian Wang, Wasi Ahmad, Hantian
677 Ding, Ming Tan, Nihal Jain, Murali Krishna Ra-
678 manathan, Ramesh Nallapati, Parminder Bhatia, Dan
679 Roth, and 1 others. 2023. Crosscodeeval: A diverse
680 and multilingual benchmark for cross-file code com-
681 pletion. *Advances in Neural Information Processing*
682 *Systems*, 36:46701–46723.

683 Yangruibo Ding, Zijian Wang, Wasi Uddin Ahmad,
684 Murali Krishna Ramanathan, Ramesh Nallapati,
685 Parminder Bhatia, Dan Roth, and Bing Xiang.

2022. Cocomic: Code completion by jointly mod- 686
687 eling in-file and cross-file context. *arXiv preprint*
688 *arXiv:2212.10007*.

689 Aryaz Eghbali and Michael Pradel. 2024. De-
690 hallucinator: Mitigating llm hallucinations in code
691 generation tasks via iterative grounding. *arXiv*
692 *preprint arXiv:2401.01701*.

693 Xinyu Gao, Yun Xiong, Deze Wang, Zhenhan Guan,
694 Zejian Shi, Haofen Wang, and Shanshan Li. 2024.
695 Preference-guided refactored tuning for retrieval aug-
696 mented code generation. In *Proceedings of the 39th*
697 *IEEE/ACM International Conference on Automated*
698 *Software Engineering*, pages 65–77.

699 Zhanming Guan, Junlin Liu, Jierui Liu, Chao Peng,
700 Dexin Liu, Ningyuan Sun, Bo Jiang, Wenchao Li, Jie
701 Liu, and Hang Zhu. 2024. Contextmodule: Improv-
702 ing code completion via repository-level contextual
703 information. *arXiv preprint arXiv:2412.08063*.

704 Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai
705 Dong, Wentao Zhang, Guanting Chen, Xiao Bi,
706 Yu Wu, YK Li, and 1 others. 2024. Deepseek-
707 coder: When the large language model meets
708 programming—the rise of code intelligence. *arXiv*
709 *preprint arXiv:2401.14196*.

710 Binyuan Hui, Jian Yang, Zeyu Cui, Jiayi Yang,
711 Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun
712 Zhang, Bowen Yu, Keming Lu, and 1 others. 2024.
713 Qwen2. 5-coder technical report. *arXiv preprint*
714 *arXiv:2409.12186*.

715 Maliheh Izadi, Jonathan Katzy, Tim Van Dam, Marc
716 Otten, Razvan Mihai Popescu, and Arie Van Deursen.
717 2024. Language models for code completion: A prac-
718 tical evaluation. In *Proceedings of the IEEE/ACM*
719 *46th International Conference on Software Engineer-*
720 *ing*, pages 1–13.

721 Huiqiang Jiang, Qianhui Wu, Xufang Luo, Dong-
722 sheng Li, Chin-Yew Lin, Yuqing Yang, and Lili Qiu.
723 2024. LongLLMLingua: Accelerating and enhanc-
724 ing LLMs in long context scenarios via prompt com-
725 pression. In *Proceedings of the 62nd Annual Meeting*
726 *of the Association for Computational Linguistics (Vol-*
727 *ume 1: Long Papers)*, pages 1658–1677, Bangkok,
728 Thailand. Association for Computational Linguistics.

729 Bowen Jin, Jinsung Yoon, Jiawei Han, and Sercan O
730 Arik. 2024. Long-context llms meet rag: Overcom-
731 ing challenges for long inputs in rag. *arXiv preprint*
732 *arXiv:2410.05983*.

733 Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying
734 Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gon-
735 zalez, Hao Zhang, and Ion Stoica. 2023. Efficient
736 memory management for large language model serv-
737 ing with pagedattention. In *Proceedings of the 29th*
738 *symposium on operating systems principles*, pages
739 611–626.

740 Vladimir I Levenshtein and 1 others. 1966. Binary
741 codes capable of correcting deletions, insertions, and

742	reversals. In <i>Soviet physics doklady</i> , volume 10, pages 707–710. Soviet Union.		
743			
744	Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou, Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, and 1 others. 2023. Starcoder: may the source be with you! <i>arXiv preprint arXiv:2305.06161</i> .		
745			
746			
747			
748			
749	Ming Liang, Xiaoheng Xie, Gehao Zhang, Xunjin Zheng, Peng Di, Hongwei Chen, Chengpeng Wang, Gang Fan, and 1 others. 2024. Repofuse: Repository-level code completion with fused dual context. <i>arXiv preprint arXiv:2402.14323</i> .		
750			
751			
752			
753			
754	Dianshu Liao, Shidong Pan, Xiaoyu Sun, Xiaoxue Ren, Qing Huang, Zhenchang Xing, Huan Jin, and Qinying Li. 2024. A 3-codgen: A repository-level code generation framework for code reuse with local-aware, global-aware, and third-party-library-aware. <i>IEEE Transactions on Software Engineering</i> .		
755			
756			
757			
758			
759			
760	Junwei Liu, Yixuan Chen, Mingwei Liu, Xin Peng, and Yiling Lou. 2024a. Stall+: Boosting llm-based repository-level code completion with static analysis. <i>arXiv preprint arXiv:2406.10018</i> .		
761			
762			
763			
764	Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2023. Lost in the middle: How language models use long contexts. <i>arXiv preprint arXiv:2307.03172</i> .		
765			
766			
767			
768			
769	Wei Liu, Ailun Yu, Daoguang Zan, Bo Shen, Wei Zhang, Haiyan Zhao, Zhi Jin, and Qianxiang Wang. 2024b. Graphcoder: Enhancing repository-level code completion via code context graph-based retrieval and language model. <i>arXiv preprint arXiv:2406.07003</i> .		
770			
771			
772			
773			
774	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 evolve-instruct. <i>arXiv preprint arXiv:2306.08568</i> .		
775			
776			
777			
778			
779	Antonio Mastropaolo, Luca Pascarella, Emanuela Guglielmi, Matteo Ciniselli, Simone Scalabrino, Rocco Oliveto, and Gabriele Bavota. 2023. On the robustness of code generation techniques: An empirical study on github copilot. In <i>Proceedings of the IEEE/ACM 45th International Conference on Software Engineering</i> , pages 2149–2160.		
780			
781			
782			
783			
784			
785			
786	Daye Nam, Andrew Macvean, Vincent Hellendoorn, Bogdan Vasilescu, and Brad Myers. 2024. Using an llm to help with code understanding. In <i>Proceedings of the IEEE/ACM 46th International Conference on Software Engineering</i> , pages 1–13.		
787			
788			
789			
790			
791	Erik Nijkamp, Bo Pang, Hiroaki Hayashi, Lifu Tu, Huan Wang, Yingbo Zhou, Silvio Savarese, and Caiming Xiong. 2023. Codegen: An open large language model for code with multi-turn program synthesis. <i>ICLR</i> .		
792			
793			
794			
795			
		Md Rizwan Parvez, Wasi Uddin Ahmad, Saikat Chakraborty, Baishakhi Ray, and Kai-Wei Chang. 2021. Retrieval augmented code generation and summarization. <i>arXiv preprint arXiv:2108.11601</i> .	796 797 798 799
		Huy Nhat Phan, Hoang Nhat Phan, Tien N Nguyen, and Nghi DQ Bui. 2024. Repohyper: Better context retrieval is all you need for repository-level code completion. <i>CoRR</i> .	800 801 802 803
		Stephen Robertson, Hugo Zaragoza, and 1 others. 2009. The probabilistic relevance framework: Bm25 and beyond. <i>Foundations and Trends® in Information Retrieval</i> , 3(4):333–389.	804 805 806 807
		Stephen E Robertson and Steve Walker. 1994. Some simple effective approximations to the 2-poisson model for probabilistic weighted retrieval. In <i>SIGIR'94: Proceedings of the Seventeenth Annual International ACM-SIGIR Conference on Research and Development in Information Retrieval, organised by Dublin City University</i> , pages 232–241. Springer.	808 809 810 811 812 813 814
		Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, and 1 others. 2023. Code llama: Open foundation models for code. <i>arXiv preprint arXiv:2308.12950</i> .	815 816 817 818 819
		Hanzhuo Tan, Qi Luo, Ling Jiang, Zizheng Zhan, Jing Li, Haotian Zhang, and Yuqun Zhang. 2024. Prompt-based code completion via multi-retrieval augmented generation. <i>ACM Transactions on Software Engineering and Methodology</i> .	820 821 822 823 824
		Ze Tang, Jidong Ge, Shangqing Liu, Tingwei Zhu, Tongtong Xu, Liguang Huang, and Bin Luo. 2023. Domain adaptive code completion via language models and decoupled domain databases. In <i>2023 38th IEEE/ACM International Conference on Automated Software Engineering (ASE)</i> , pages 421–433. IEEE.	825 826 827 828 829 830
		Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrutu Bhosale, and 1 others. 2023. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> .	831 832 833 834 835 836
		Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. <i>Advances in neural information processing systems</i> , 30.	837 838 839 840 841
		Chaozheng Wang, Junhao Hu, Cuiyun Gao, Yu Jin, Tao Xie, Hailiang Huang, Zhenyu Lei, and Yuetang Deng. 2023a. Practitioners' expectations on code completion. <i>arXiv preprint arXiv:2301.03846</i> .	842 843 844 845
		Chong Wang, Jian Zhang, Yebo Feng, Tianlin Li, Weisong Sun, Yang Liu, and Xin Peng. 2024a. Teaching code llms to use autocompletion tools in repository-level code generation. <i>ACM Transactions on Software Engineering and Methodology</i> .	846 847 848 849 850

851	Jicheng Wang, Yifeng He, and Hao Chen. 2024b.	Kechi Zhang, Jia Li, Ge Li, Xianjie Shi, and Zhi Jin.	908
852	Repogenreflex: Enhancing repository-level	2024. Codeagent: Enhancing code generation with	909
853	code completion with verbal reinforcement and	tool-integrated agent systems for real-world repo-	910
854	retrieval-augmented generation. <i>arXiv preprint</i>	level coding challenges. In <i>Proceedings of the 62nd</i>	911
855	<i>arXiv:2409.13122</i> .	<i>Annual Meeting of the Association for Computational</i>	912
856	Yanlin Wang, Yanli Wang, Daya Guo, Jiachi Chen,	<i>Linguistics (Volume 1: Long Papers)</i> , pages 13643–	913
857	Ruikai Zhang, Yuchi Ma, and Zibin Zheng. 2024c.	13658.	914
858	RlCoder: Reinforcement learning for repository-level	Lei Zhang, Yunshui Li, Jiaming Li, Xiaobo Xia, Jiaxi	915
859	code completion. <i>arXiv preprint arXiv:2407.19487</i> .	Yang, Run Luo, Minzheng Wang, Longze Chen, Jun-	916
860	Yejie Wang, Keqing He, Guanting Dong, Pei Wang,	hao Liu, Qiang Qu, and 1 others. 2025a. Hierarchical	917
861	Weihao Zeng, Muxi Diao, Weiran Xu, Jingang Wang,	context pruning: Optimizing real-world code com-	918
862	Mengdi Zhang, and Xunliang Cai. 2024d. Dolphin-	pletion with repository-level pretrained code llms.	919
863	coder: Echo-locating code large language models	In <i>Proceedings of the AAAI Conference on Artificial</i>	920
864	with diverse and multi-objective instruction tuning.	<i>Intelligence</i> , volume 39, pages 25886–25894.	921
865	In <i>Proceedings of the 62nd Annual Meeting of the</i>	Sheng Zhang, Yifan Ding, Shuquan Lian, Shun Song,	922
866	<i>Association for Computational Linguistics (Volume</i>	and Hui Li. 2025b. Coderag: Finding relevant	923
867	<i>1: Long Papers)</i> , pages 4706–4721.	and necessary knowledge for retrieval-augmented	924
868	Yue Wang, Hung Le, Akhilesh Gotmare, Nghi Bui, Jun-	repository-level code completion. In <i>Proceedings</i>	925
869	nan Li, and Steven Hoi. 2023b. Codet5+: Open code	<i>of the 2025 Conference on Empirical Methods in</i>	926
870	large language models for code understanding and	<i>Natural Language Processing</i> , pages 23289–23299.	927
871	generation. In <i>Proceedings of the 2023 Conference</i>	Yanzhao Zhang, Mingxin Li, Dingkun Long, Xin Zhang,	928
872	<i>on Empirical Methods in Natural Language Process-</i>	Huan Lin, Baosong Yang, Pengjun Xie, An Yang,	929
873	<i>ing</i> , pages 1069–1088.	Dayiheng Liu, Junyang Lin, and 1 others. 2025c.	930
874	Di Wu, Wasi Uddin Ahmad, Dejiào Zhang, Murali Kr-	Qwen3 embedding: Advancing text embedding and	931
875	ishna Ramanathan, and Xiaofei Ma. 2024. Repo-	reranking through foundation models. <i>arXiv preprint</i>	932
876	former: selective retrieval for repository-level code	<i>arXiv:2506.05176</i> .	933
877	completion. In <i>Proceedings of the 41st International</i>		
878	<i>Conference on Machine Learning</i> , pages 53270–		
879	53290.		
880	Peiyang Wu, Nan Guo, Junliang Lv, Xiao Xiao, and		
881	Xiaochun Ye. 2025. Rtlrepecoder: Repository-level		
882	rtl code completion through the combination of fine-		
883	tuning and retrieval augmentation. <i>arXiv preprint</i>		
884	<i>arXiv:2504.08862</i> .		
885	Frank F Xu, Bogdan Vasilescu, and Graham Neubig.		
886	2022. In-ide code generation from natural language:		
887	Promise and challenges. <i>ACM Transactions on</i>		
888	<i>Software Engineering and Methodology (TOSEM)</i> ,		
889	31(2):1–47.		
890	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang,		
891	Binyuan Hui, Bo Zheng, Bowen Yu, Chang		
892	Gao, Chengen Huang, Chenxu Lv, and 1 others.		
893	2025. Qwen3 technical report. <i>arXiv preprint</i>		
894	<i>arXiv:2505.09388</i> .		
895	Xinran Yu, Chun Li, Minxue Pan, and Xuandong Li.		
896	2024. Droidcoder: Enhanced android code comple-		
897	tion with context-enriched retrieval-augmented gen-		
898	eration. In <i>2024 39th IEEE/ACM International Con-</i>		
899	<i>ference on Automated Software Engineering (ASE)</i> ,		
900	pages 681–693. IEEE.		
901	Fengji Zhang, Bei Chen, Yue Zhang, Jacky Keung, Jin		
902	Liu, Daoguang Zan, Yi Mao, Jian-Guang Lou, and		
903	Weizhu Chen. 2023. Repocoder: Repository-level		
904	code completion through iterative retrieval and gen-		
905	eration. In <i>Proceedings of the 2023 Conference on</i>		
906	<i>Empirical Methods in Natural Language Processing</i> ,		
907	pages 2471–2484.		

934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982

A Preliminary Studies

A.1 Setup

To prevent potential data leakage, we selected four popular open-source code LLMs that were released prior to the collection of benchmark datasets for our experiments: StarCoderBase-7B (Li et al., 2023), CodeLlama-7B (Roziere et al., 2023), DeepSeekCoder-1.3B (Guo et al., 2024), and DeepSeekCoder-6.7B. We adopt the default greedy decoding strategy without sampling, and set the `max_new_tokens` to 64, `max_context_tokens` to 2048 for all models. All experiments were conducted on NVIDIA A100 GPUs.

A.2 Baseline Evaluation

A.2.1 In-file Context Only

We first provided code LLMs with in-file context, specifically the left context of the incomplete code without any cross-file context, so that the models could rely solely on their own prior knowledge to perform the completion task. The results are presented in Table 3. The overall performance of the completion was unsatisfactory. Even the top-performing model demonstrated a limited accuracy, reaching just 9.46% on the Python portion of the CCEval benchmark. This suggests that, without access to broader cross-file context, code LLMs struggle to accurately complete code that depends on definitions or information residing outside the current file.

A.2.2 Simple RAG

RAG is a widely adopted strategy for enhancing LLM-based repository-level code completion. In the preliminary study, we employ a simple RAG approach in which the left context of the incomplete code serves as a query to retrieve relevant code snippets from the repository. These retrieved snippets are then concatenated with the in-file context to construct the prompt. Following the setup of RLCoder (Wang et al., 2024c) and reflecting typical programming practices, we segment the repository code into code snippets based on blank lines, and merge consecutive short code fragments to ensure each snippet has an appropriate length and improved semantic coherence. We utilize BM25 (Robertson et al., 2009; Robertson and Walker, 1994) as a sparse similarity metric and also employ the popular dense retriever RLRetriever (Wang et al., 2024c) to retrieve relevant code snippets. The results are shown in Table 3. It can be ob-

served that code snippets retrieved from other files within the repository contribute to more accurate code completion, as they provide relevant contextual information for the completion task.

A.3 Query Content Analysis

In repository-level code completion tasks, the query provided to the retriever is typically an incomplete code segment. Due to the absence of the correct completion or similar code, the retrieved snippets are often limited in utility, frequently containing redundant or irrelevant content and lacking critical information associated with the desired completion. To address this problem, the query should be modified in a specific way. For example, by inputting the original query into a code LLM to generate multiple lines of completed code and appending these lines to the original query, a refined query can be constructed. In this preliminary analysis, we appended only the first generated line to the original query, as the initial line typically represents the most probable and locally relevant continuation. The results are presented in Table 4, indicate that this augmented query incorporates code snippets that are similar to the correct completion, thus enabling the retrieval of more precise and relevant code snippets.

A.4 Retrieved Code Snippets Position Analysis

Additionally, we investigate the impact of the position of retrieved code snippets within the prompt on code completion performance by varying their order before inputting the prompt to the code LLMs. In previous section, the default arrangement places code snippets in descending order of relevance to the query. In this experiment, we also consider reversing this order as well as randomizing the placement of the snippets, in order to examine how the position of retrieved code snippets within the prompt affects the final completion results. The results are presented in Table 5. We observe that placing relevant code snippets at the beginning or end of the prompt has a comparable impact on completion performance. However, when retrieved code snippets are randomly positioned within the prompt, the completion accuracy drops significantly. This indicates that the placement of retrieved code snippets within the prompt affects the completion results, which is consistent with findings from previous studies on natural language tasks (Liu et al., 2023).

983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1030
1031

Table 3: The completion results of the baseline methods. **EM** denotes Exact Match, and **ES** denotes Edit Similarity.

		Baseline Evaluation							
Code LLM	Method	CCEval (Python)		CCEval (Java)		RepoEval (Line)		RepoEval (API)	
		EM	ES	EM	ES	EM	ES	EM	ES
StarCoderBase-7B	In-file Only	6.79	60.69	9.44	63.83	36.25	59.95	26.19	51.83
	Simple RAG (BM25)	14.78	66.30	16.36	66.41	44.19	65.48	32.69	55.91
	Simple RAG (RLRetriever)	25.89	72.12	25.81	68.91	47.75	68.48	35.06	58.09
CodeLlama-7B	In-file Only	7.39	60.78	9.35	63.83	37.00	60.47	29.19	57.43
	Simple RAG (BM25)	15.76	66.30	16.78	65.37	44.25	65.65	37.00	63.43
	Simple RAG (RLRetriever)	26.57	71.43	26.23	66.37	46.63	67.88	37.94	64.32
DeepSeekCoder-1.3B	In-file Only	5.85	58.21	6.92	58.37	33.88	58.46	26.94	55.86
	Simple RAG (BM25)	14.18	65.12	13.70	61.63	41.31	63.60	35.31	62.07
	Simple RAG (RLRetriever)	24.05	70.46	20.85	63.43	44.12	66.46	36.06	62.69
DeepSeekCoder-6.7B	In-file Only	9.46	61.60	11.41	62.54	39.63	61.95	30.44	59.40
	Simple RAG (BM25)	18.31	68.38	17.48	65.23	45.94	66.68	38.25	64.99
	Simple RAG (RLRetriever)	30.24	73.58	26.13	66.11	48.75	69.42	39.88	66.24

Table 4: Comparison of completion results using different methods of changing query.

		Query Content Analysis							
Code LLM	Method	CCEval (Python)		CCEval (Java)		RepoEval (Line)		RepoEval (API)	
		EM	ES	EM	ES	EM	ES	EM	ES
StarCoderBase-7B	Origin	25.89	72.12	25.81	68.91	47.75	68.48	35.06	58.09
	Append	26.34	72.40	26.65	69.44	47.56	68.34	35.13	58.19
CodeLlama-7B	Origin	26.57	71.43	26.23	66.37	46.63	67.88	37.94	64.32
	Append	27.24	72.75	26.79	68.05	46.69	67.92	37.81	64.22
DeepSeekCoder-1.3B	Origin	24.05	70.46	20.85	63.43	44.12	66.46	36.06	62.69
	Append	24.84	70.28	21.37	62.73	44.13	66.48	35.88	62.61
DeepSeekCoder-6.7B	Origin	30.24	73.58	26.13	66.11	48.75	69.42	39.88	66.24
	Append	30.66	74.79	26.37	67.28	49.63	69.62	40.25	66.46

Table 5: Comparison of completion results using different sort order of retrieved code snippets in prompt.

		Relevant Information Position Analysis							
Code LLM	Method	CCEval (Python)		CCEval (Java)		RepoEval (Line)		RepoEval (API)	
		EM	ES	EM	ES	EM	ES	EM	ES
StarCoderBase-7B	Default	25.89	72.12	25.81	68.91	47.75	68.48	35.06	58.09
	Reverse	26.68	72.24	25.90	69.71	47.81	68.22	34.56	58.07
	Random	24.47	71.31	24.87	68.74	45.94	66.83	33.38	57.10
CodeLlama-7B	Default	26.57	71.43	26.23	66.37	46.63	67.88	37.94	64.32
	Reverse	26.15	72.12	26.41	68.83	46.94	68.20	38.63	64.54
	Random	24.20	70.98	23.94	66.83	44.88	66.79	36.38	63.09
DeepSeekCoder-1.3B	Default	24.05	70.46	20.85	63.43	44.12	66.46	36.06	62.69
	Reverse	24.92	70.60	21.46	64.06	44.75	66.42	36.56	62.94
	Random	21.88	68.27	19.92	61.92	42.00	65.18	34.44	61.46
DeepSeekCoder-6.7B	Default	30.24	73.58	26.13	66.11	48.75	69.42	39.88	66.24
	Reverse	29.91	74.28	26.88	67.72	48.75	69.52	40.56	66.76
	Random	26.75	73.07	24.68	66.62	46.56	68.35	38.06	65.01

Table 6: Average cost (seconds) for each step of RepoAttention on CCEval Python.

Dual-Path Query Refinement				Rerank	Completion
Distillation & Enrichment	Sparse Retrieval	Sample	Dense Retrieval		
0.045	0.022	0.549	0.662	0.159	0.185

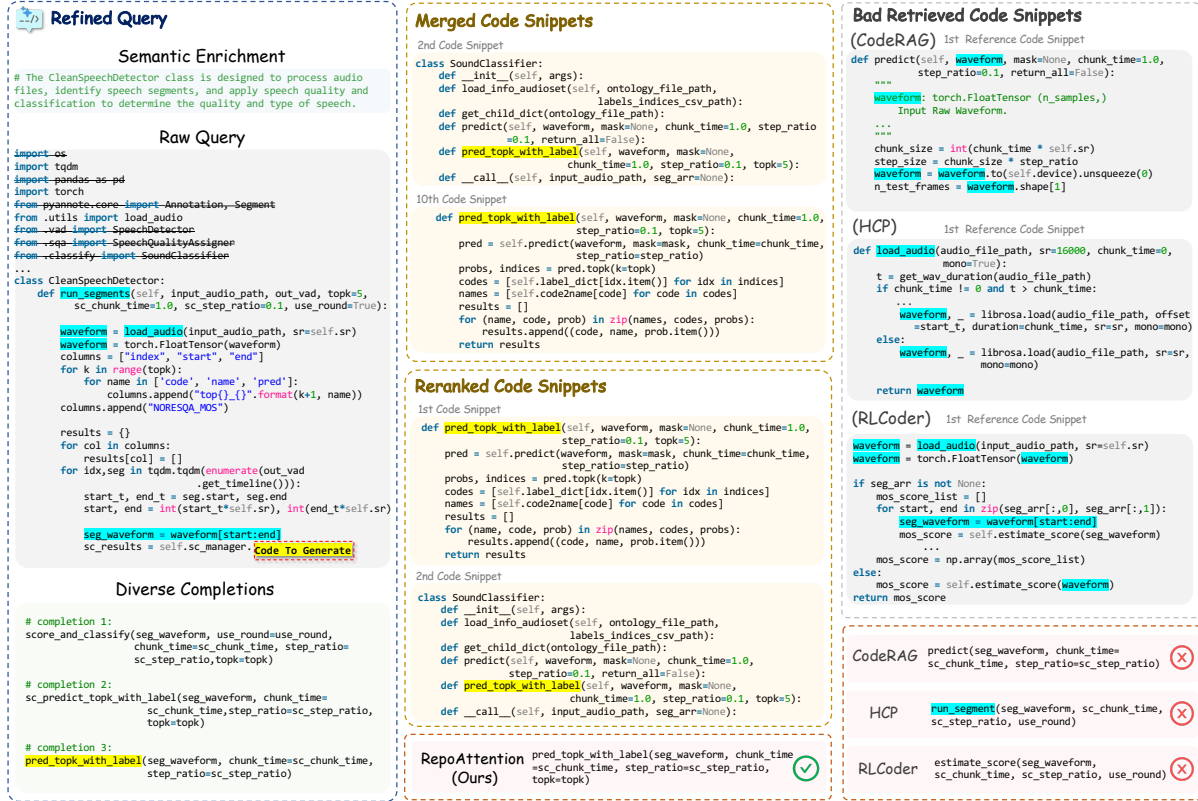


Figure 5: A case (*project_cc_python/1474*) from CCEval Python. The highlighted words indicate overlaps between the query and the code snippet. Strikethrough text marks code that is removed by the query distillation step.

Table 7: Average cost (seconds) per case on CCEval Python.

HCP	CodeRAG	RLCoder	RepoAttention
1.95	2.53	1.67	1.62

B Computational Cost Details

Table 6 reports the average time consumption for each step of RepoAttention. We observe that the primary computational costs stem from dense retrieval and sampling. This is consistent with expectations. The sampling step involves LLM inference which is naturally resource-intensive, while dense retrieval requires similarity searches in a high-dimensional vector space.

It is worth noting that our current implementation of RepoAttention prioritizes accuracy and

has not been heavily optimized for speed. The latency could be further reduced through engineering improvements. For instance, model quantization could be applied to reduce memory bandwidth pressure. Additionally, employing accelerated vector search libraries (such as Faiss³) would likely reduce the overhead of the dominant stages, further enhancing the practical deployment of RepoAttention.

C Case Study

We present a case study to demonstrate the effectiveness of RepoAttention. Figure 5 shows an example where the model is required to complete the function `pred_topk_with_label` in the `CleanSpeechDetector` class. As illustrated by the raw query, the query distillation step removes

³<https://github.com/facebookresearch/faiss>

1059 irrelevant code, indicated by the strikethrough
 1060 lines. The resulting refined query enables the
 1061 model to sample useful candidate completions.
 1062 The reranking component then prioritizes the
 1063 most relevant snippet containing the definition of
 1064 `pred_topk_with_label` and places it at the top of
 1065 the context. This precise retrieval enables RepoAt-
 1066 tention to correctly complete the target function.
 1067 In contrast, baselines like CodeRAG, HCP and
 1068 RLCoder retrieve irrelevant snippets that appear
 1069 textually similar to the incomplete code, leading to
 1070 incorrect completions.

D Algorithm of RepoAttention

Algorithm 1: RepoAttention

Input: Incomplete code c_{in} , Repository \mathcal{R} , Code
 LLM \mathcal{M} , Reranker \mathcal{R}_{rank} , Relevance
 threshold I , Number of samples m , Key code
 lines t , Top- k retrieval

Output: Code completion result \hat{c}

1 **Step 1: Codebase Construction**
 2 Parse \mathcal{R} into primary snippets via blank-line split;
 3 Extract dependency snippets via tree-sitter parsing;
 4 $\mathcal{D} \leftarrow \{\text{primary snippets}\} \cup \{\text{dependency snippets}\}$;
 5 **Step 2: Dual-Path Query Refinement**
 6 $q_{raw} \leftarrow c_{in}$;
 7 $q_{ref} \leftarrow \text{DISTILLANDENRICH}(q_{raw})$
 // Sparse retrieval via BM25
 8 $S_{sparse}^{raw} \leftarrow \text{BM25}(q_{raw}, \mathcal{D})$;
 9 $S_{sparse}^{ref} \leftarrow \text{BM25}(q_{ref}, \mathcal{D})$;
 // Generate hypothetical completions
 10 $H_{raw} \leftarrow \{h_1, \dots, h_m\} \sim P_{\mathcal{M}}(\cdot | S_{sparse}^{raw} \oplus q_{raw})$;
 11 $H_{ref} \leftarrow \{h'_1, \dots, h'_m\} \sim P_{\mathcal{M}}(\cdot | S_{sparse}^{ref} \oplus q_{ref})$;
 // Dense retrieval with expanded queries
 12 $S_{dense}^{raw} \leftarrow \text{DenseRetrieve}(q_{raw} \oplus H_{raw}, \mathcal{D})$;
 13 $S_{dense}^{ref} \leftarrow \text{DenseRetrieve}(q_{ref} \oplus H_{ref}, \mathcal{D})$;
 14 **Step 3: Relevance-Aware Snippet Reranking**
 15 $c_k \leftarrow$ Extract last t lines from q_{ref} as key code;
 // Compute relevance score
 16 **foreach** snippet $s_i \in S_{dense}^{raw} \cup S_{dense}^{ref}$ **do**
 17 $r_i \leftarrow \mathcal{R}_{rank}(c_k, s_i)$
 18 $C_v \leftarrow$ Sort snippets by r_i in descending order;
 19 **Step 4: Dynamic Relevance-Guided Attention**
 20 $E \leftarrow \emptyset$;
 21 **foreach** snippet $s_i \in C_v$ with tokens T_i **do**
 22 **if** $r_i \geq I$ **then**
 23 $E \leftarrow E \cup T_i$;
 24 **else**
 25 $p_i \leftarrow \text{MAPTOPROBABILITY}(r_i)$;
 26 **if** $\text{BERNOULLI}(p_i) = I$ **then**
 27 $E \leftarrow E \cup T_i$;
 28 $E \leftarrow E \cup T_q$; // Always include query tokens
 // Generate completion with DRGA
 29 $\hat{c} \leftarrow \mathcal{M}.\text{Complete}(C_v \oplus q_{ref})$ with attention
 restricted to E ;
 30 **return** \hat{c} ;
