

Hierarchical Context Pruning: Optimizing Real-World Code Completion with Repository-Level Pretrained Code LLMs

Lei Zhang^{1,2} Yunshui Li^{1,2} Jiaming Li^{1,2} Xiaobo Xia³ Jiayi Yang^{1,2} Run Luo^{1,2}
Minzheng Wang^{2,5} Longze Chen^{1,2} Junhao Liu⁴ Min Yang^{1,2†}

¹Shenzhen Institute of Advanced Technology, Chinese Academy of Sciences

²University of Chinese Academy of Sciences

³The University of Sydney ⁴University of California, Irvine

⁵MAIS, Institute of Automation, Chinese Academy of Sciences
{lei.zhang2, min.yang}@siat.ac.cn

Abstract

Some recently developed code large language models (Code LLMs) have been pretrained on repository-level code data (Repo-Code LLMs), enabling these models to recognize repository structures and utilize cross-file information for code completion. However, in real-world development scenarios, simply concatenating the entire code repository often exceeds the context window limits of these Repo-Code LLMs, leading to significant performance degradation. In this study, we conducted extensive preliminary experiments and analyses on six Repo-Code LLMs. The results indicate that maintaining the topological dependencies of files and increasing the code file content in the completion prompts can improve completion accuracy; pruning the specific implementations of functions in all dependent files does not significantly reduce the accuracy of completions. Based on these findings, we proposed a strategy named **Hierarchical Context Pruning (HCP)** to construct completion prompts with high informational code content. The **HCP** models the code repository at the function level, maintaining the topological dependencies between code files while removing a large amount of irrelevant code content, significantly reduces the input length for repository-level code completion. We applied the **HCP** strategy in experiments with six Repo-Code LLMs, and the results demonstrate that our proposed method can significantly enhance completion accuracy while substantially reducing the length of input. Our code and data are available at <https://github.com/Hambaobao/HCP-Coder>.

adopted in daily development practices and have significantly enhanced the productivity of developers. As research (Bavarian et al., 2022; Sun et al., 2024) on code large language models (Code LLMs) continues to evolve, some recently developed Code LLMs (Guo et al., 2024; Lozhkov et al., 2024; Team et al., 2024) have been trained on repository-level code data (Repo-Code LLMs) to overcome the limitations of previous models trained on file-level data, which struggled to recognize repository structures and integrate code across multiple files for completion tasks. However, in real-world development scenarios, simply concatenating the entire code repository often exceeds the context window size of these Repo-Code LLMs, leading to significant performance degradation and increased inference latency. How to effectively utilize the capabilities of these Repo-Code LLMs to integrate cross-file information and construct high-quality completion prompts within the model’s context window limits remains an area for further exploration.

In this study, we initially evaluated six Repo-Code LLMs on the CrossCodeEval (Ding et al., 2023) benchmark and conducted a detailed analysis of completion errors (Appendix A). The errors identified were categorized into eight distinct classes (Section 4.2). Subsequently, considering the characteristics of the decoder architecture in Code LLMs, we analyzed the impact of topological dependencies among code files on completion accuracy (Section 4.3). We found that maintaining the dependencies between code files and including more file information leads to higher accuracy. Additionally, we conducted experiments to analyze the impact of content from files at different dependency levels on completion accuracy (Section 4.4). We discovered that even pruning away the specific implementations of functions in all dependent files does not significantly reduce the accuracy of completions. Based on the results of these preliminary experiments, we proposed a strategy named **Hi-**

1 Introduction

Code completion tools based on code large language models (Chen et al., 2021; Nijkamp et al., 2023b; Li et al., 2023; Fried et al., 2023; Allal et al., 2023), such as *GitHub Copilot*¹, have been widely

[†]Min Yang is the corresponding author.

¹<https://github.com/features/copilot>

erarchical Context Pruning (HCP) to construct high-quality completion prompts. The **HCP** models the code repository at the function level, retaining the topological dependencies between files while eliminating a large amount of irrelevant code content. In our experiments, the **HCP** successfully reduced the input from over 50,000 tokens to approximately 8,000 tokens, and significantly enhanced the accuracy of completions.

In summary, our contributions are threefold:

- We conducted experiments on six Repo-Code LLMs and found that: maintaining the topological dependencies of files and increasing the content of code files in the completion prompts can enhance completion accuracy; pruning the specific implementations of functions in all dependent files does not significantly reduce the accuracy of completions.
- Based on the results of preliminary experiments, we proposed a strategy named **Hierarchical Context Pruning (HCP)** for constructing high-quality completion prompts, which models the code repository at the function level, retaining the topological dependencies between files while eliminating a large amount of irrelevant code content.
- We applied the **HCP** strategy in experiments with six Repo-Code LLMs, and the results demonstrate that our proposed method can significantly enhance completion accuracy while substantially reducing the length of input.

2 Related Work

2.1 Code Large Language Models

2.1.1 Infilling Code LLMs

Infilling scenarios constitute the majority of code completion tasks in the real world. [Bavarian et al. \(2022\)](#) demonstrates that pre-training Code LLMs with a certain proportion of fill-in-the-middle format code data can enable the Code LLMs to fill in middle code based on the surrounding context, without compromising their original left-to-right generation performance. Based on the findings of [Bavarian et al. \(2022\)](#), many subsequent Code LLMs ([Fried et al., 2023](#); [Allal et al., 2023](#); [Nijkamp et al., 2023a](#); [Li et al., 2023](#); [Rozière et al., 2024](#); [Guo et al., 2024](#); [Pinnaparaju et al., 2024](#); [Lozhkov et al., 2024](#)) have emerged with the capability to perform infilling.

2.1.2 Instruction Code LLMs

Pretrained Code LLMs are traditionally used only for continuation tasks such as code completion. Inspired by works on instruction tuning large language models ([Ouyang et al., 2022](#); [Li et al., 2024b](#)), many studies ([Wang et al., 2023a](#); [Luo et al., 2023](#); [Muennighoff et al., 2024](#); [Xu et al., 2023](#); [Wang et al., 2024b](#); [Zheng et al., 2024](#)) have attempted to finetune Code LLMs using code instruction data. This finetuning unlocks the potential of Code LLMs, enabling them to perform more complex coding tasks based on user instructions.

2.2 Code Benchmarks

2.2.1 Code Completion Benchmarks

HumanEval ([Chen et al., 2021](#)) consists of 164 manually crafted Python code problems, with an average of 7.7 tests each test case. MBPP ([Austin et al., 2021](#)) is designed for individuals with entry-level programming skills. It comprises 974 concise Python functions, each with an accompanying description in English, a specified function signature, and three manually crafted test cases for verification. MultiPL-E ([Cassano et al., 2022](#)) introduces itself as a novel benchmarking framework designed for multilingual contexts, building upon HumanEval ([Chen et al., 2021](#)) and MBPP ([Austin et al., 2021](#)). APPS ([Hendrycks et al., 2021](#)) is a benchmark including 10K less-restricted problems for code generation. CodeContests ([Li et al., 2022](#)) is a dataset specifically for competitive programming problems.

2.2.2 Infilling Code Benchmarks

[Fried et al. \(2023\)](#) constructed *single-line* and *multi-line* infilling completion tasks based on HumanEval, and [Bavarian et al. \(2022\)](#) expanded upon it to create *randomspan* infilling completion tasks, ultimately resulting in the current HumanEval-Infilling benchmark. [Allal et al. \(2023\)](#) created an Infilling benchmark that includes languages from Java, JavaScript, and Python 3, utilizing a line exactly match method for evaluation. [Lai et al. \(2022\)](#) presents a benchmark for evaluating the performance of Code LLMs in completing tasks related to Python scientific computing libraries, encompassing both regular completion and insertion (infilling) tasks.

2.3 Repo-level Code Completion

Some benchmarks for repository-level code completion have been proposed to evaluate the perfor-

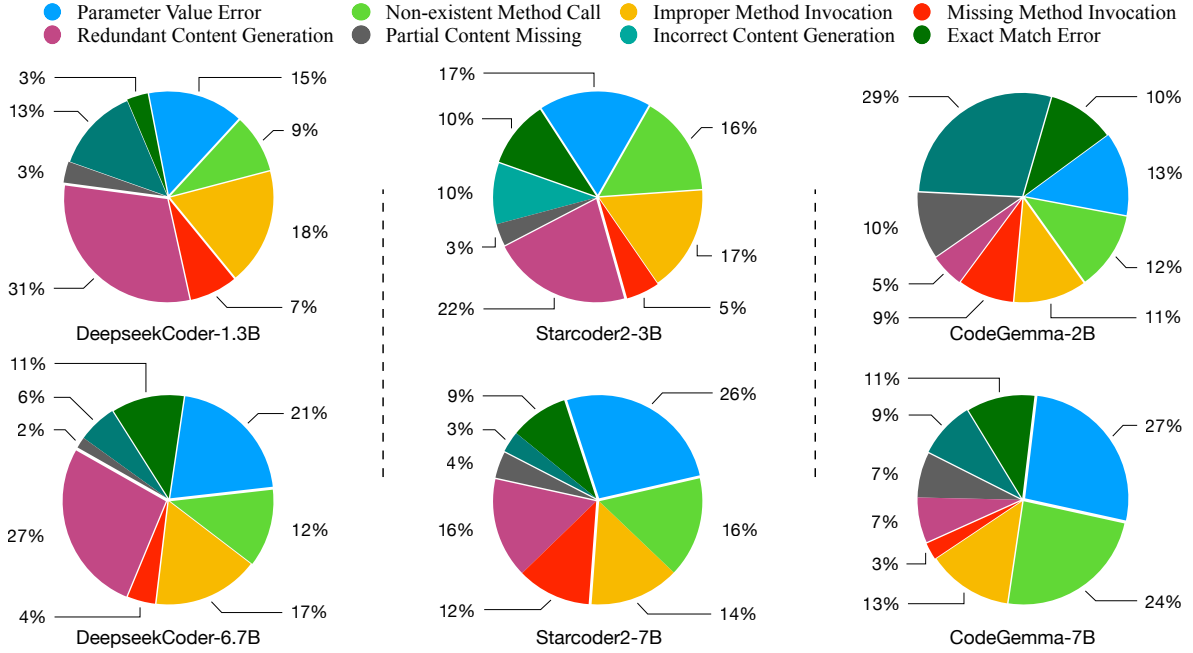


Figure 1: The error class distribution of the completion results of the DeepseekCoder, Starcoder2 and CodeGemma models on the CrossCodeEval: Python benchmark.

mance of code models in real-world completion tasks, such as CrossCodeEval (Ding et al., 2023), Repo-Bench (Liu et al., 2023), CoderEval (Zhang et al., 2024b), and EvoCodeBench (Li et al., 2024a). A lot of studies (Shrivastava et al., 2023; Zhang et al., 2023a; Bi et al., 2024; Phan et al., 2024; Liang et al., 2024) have focused on improving the accuracy of repository-level code completion tasks. However, most of these studies overlook the unique aspects of their Fill-in-the-Middle (FIM) capacities. Furthermore, despite the recent development of repository-level pretrained Code LLMs designed to process large-scale repository data, research on these models remains relatively limited.

3 Experiments Setup

3.1 Dataset & Evaluation Metrics

To assess the code completion performance of Code LLMs in real development scenarios, we utilized CrossCodeEval (Ding et al., 2023) as the evaluation dataset. The CrossCodeEval (Ding et al., 2023) benchmark provides test cases that require the use of cross-file code information for completion. Without loss of generality, in this study, we have chosen Python language as the primary language for our research.

We used the original data from CrossCodeEval, retaining the original repository structure. For each test case, we first identified the file for comple-

tion and the cursor’s position (the line and column where the completion occurs). We then removed the code after the cursor in that line to form authentic completion test cases. Ultimately, we obtained 2,655 real-world completion tests. Following the CrossCodeEval evaluation protocol, we evaluated the completion results using two metrics: *Exact Match* (EM) and *Edit Similarity* (ES).

3.2 Models & Prompt Templates

The code large language models pretrained with repository-level code data include specific tokens used to describe the repository structure in the prompt. Table 6 in appendix displays the special tokens used by DeepseekCoder, Starcoder2 and CodeGemma. The specific prompt templates used by DeepseekCoder, Starcoder2 and CodeGemma are shown in Table 7.

3.3 Hardware & Hyperparameters

All the experiments were conducted on NVIDIA A100 GPUs. We employ greedy decoding strategy for all the models, and set `max_new_tokens` to 32. The `model_max_length` of DeepseekCoder, Starcoder2 and CodeGemma is set to 16, 352, 16, 352 and 8, 160, respectively. All the prompts longer than the `model_max_length` are truncated from the left.

XF-Context	Baseline Evaluation											
	DScoder-1.3B		DScoder-6.7B		Starcoder2-3B		Starcoder2-7B		CodeGemma-2B		CodeGemma-7B	
	EM	ES	EM	ES	EM	ES	EM	ES	EM	ES	EM	ES
Infile-Only	16.72	56.58	28.14	68.36	21.92	61.49	22.98	63.58	20.64	56.26	30.58	70.36
RAG-BM25	17.28	58.18	32.65	71.78	24.45	63.84	26.26	65.32	22.89	57.73	32.89	70.81
Random-All	6.18	46.19	33.94	70.98	28.32	66.87	31.45	69.09	26.93	62.13	36.69	74.42

Table 1: The completion results of the baseline methods. **EM** denotes Exact Match, and **ES** denotes Edit Similarity.

XF-Context	Topological Dependency Analysis											
	DScoder-1.3B		DScoder-6.7B		Starcoder2-3B		Starcoder2-7B		CodeGemma-2B		CodeGemma-7B	
	EM	ES	EM	ES	EM	ES	EM	ES	EM	ES	EM	ES
D-Level: 1	15.44	55.03	33.03	70.77	26.18	64.15	28.51	66.91	24.37	58.79	34.65	73.01
D-Level: 2	13.63	53.45	33.56	70.74	26.70	64.58	29.45	67.03	25.31	59.27	35.67	73.26
D-Level: 3	13.26	53.17	33.07	70.51	26.82	64.56	29.23	67.01	25.35	59.30	35.93	73.34
D-Level: 4	13.37	53.20	33.22	70.57	26.59	64.46	29.53	67.07	25.54	59.42	36.12	73.54
D-Level: ∞	5.76	46.22	35.29	71.51	30.43	67.34	33.03	69.57	29.08	62.91	39.32	75.35

Table 2: Comparison of completion results using different context dependency levels across 6 models. All the prompts is truncated to the max context window of the Code LLMs from the left. ∞ denotes the prompt including all files in the repository.

4 Preliminary Studies

4.1 Baseline Evaluation

4.1.1 Infile Only

We initially evaluated the model’s completion ability using only information from the current file, with results presented in Table 1 under the *Infile-Only* row. The completion results are less than satisfactory. Even the best-performing model achieved an accuracy of only about 30%.

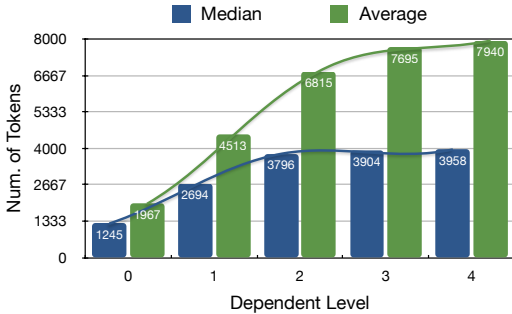


Figure 2: The distribution of tokenized prompt lengths in the CrossCodeEval benchmark. The x-axis represents the dependent level, and the y-axis represents the number of tokens.

4.1.2 RAG-BM25

We subsequently evaluated the effect of using Retrieval Augmented Generation (RAG) method to

retrieve relevant code snippets to assist with completion. Following the setup of CrossCodeEval, we chunk the repository code into units of 10 lines, and use BM25 as similarity metric for retrieving relevant code snippets. We select the top-5 relevant snippets as cross-file information, which are placed at the beginning of the prompt to assist with code generation. The results are shown in Table 1 under the *RAG-BM25* row.

4.1.3 Randomly Concatenating All Files

Additionally, we concatenated all repository code files randomly according to the pre-trained formats of various Repo-Code LLMs to create completion prompts, which were then input into the models for completion. The evaluation results are shown in Table 1 under the *Random-All* row. We observed that supplying the model with more information from the repository’s code led to superior performance compared to RAG. However, the input length of the model is limited by its context window, thereby transforming this scenario into a constrained optimization problem. The constrained optimization goal is expressed as follows:

$$\max_{\mathcal{P}} \text{Quality}(\mathcal{P}) \quad \text{s.t.} \quad \text{Length}(\mathcal{P}) \leq L \quad (1)$$

where \mathcal{P} represents the coconstructed prompt, $\text{Quality}(\mathcal{P})$ represents the quality of the coon-

XF-Context	Cross-File Content Analysis											
	DScoder-1.3B		DScoder-6.7B		StarCoder2-3B		StarCoder2-7B		CodeGemma-2B		CodeGemma-7B	
	EM	ES	EM	ES	EM	ES	EM	ES	EM	ES	EM	ES
P-Level: 0	6.18	46.19	33.94	70.98	28.32	66.87	31.45	69.09	26.93	62.13	36.69	74.42
P-Level: 1	6.55	46.58	36.20	71.90	30.73	67.97	34.43	70.65	29.30	63.46	39.55	75.70
P-Level: 2	9.83	49.63	34.73	70.89	30.02	66.41	31.26	68.24	27.34	61.13	38.31	74.32
+ <i>D-level:1</i>	9.45	49.44	36.87	72.14	29.91	66.96	32.62	69.11	28.93	62.03	39.17	75.16
+ <i>D-level:2</i>	8.70	48.61	36.38	71.66	29.64	66.99	32.96	69.13	28.44	61.76	39.06	74.91

Table 3: The results of completion using cross-file information with different pruning levels. + *D-level:x* denotes the model uses the cross-file information with dependency level x.

structured prompt, $\text{Length}(\mathcal{P})$ represents the length of the constructed prompt, and L represents context window size of the model.

4.2 Completion Error Analysis

To further investigate the issues of repository-level pre-trained Code LLMs in real-world completion tasks, we sampled 200 error examples from each model’s *Random-All* evaluation results for error analysis. Ultimately, we categorized the issues present in these models into eight classes: *Parameter Value Error*, *Non-existent Method Call*, *Improper Method Invocation*, *Missing Method Invocation*, *Redundant Content Generation*, *Partial Content Missing*, *Incorrect Content Generation*, and *Exact Match Error*. Figure 1 shows the error distribution statistics for six Repo-Code LLMs. In the appendix A, we provide examples of each type of error along with corresponding error analysis.

4.3 Topological Dependency Analysis

Definition 1. (Dependency Level) Let F denote a set of files in a code repository, and let $f \in F$ represent a specific file. We define the dependency levels as follows:

$$\begin{aligned}
 I(f) &= \{g \mid g \text{ is imported by } f\} \\
 D_0(f) &= \{f\} \\
 D_{i+1}(f) &= D_i(f) \cup I(D_i(f))
 \end{aligned} \tag{2}$$

We first identified the file requiring completion, then extracted all the import statements from the file with *Tree-Sitter*², and used a breadth-first search (BFS) method to progressively add dependent files. Algorithm 1 in appendix shows our specific dependency modeling process.

Figure 13 illustrates the growth in the number of dependent files (calculated by the length of the

tokenized prompt) as the number of dependency layers increases. We used median and average as statistical measures and found that in the vast majority of cases, the number of dependent files for a single file increases slowly after reaching four layers of dependencies. This suggests that using four layers of dependencies is sufficient to cover most scenarios. We further define:

$$D_\infty(f) = D_4(f) \cup \{F \setminus D_4(f)\} \tag{3}$$

to represent the prompt including all files in the repository.

In Table 2, the D-level rows show the results of completion using cross-file information with different dependency levels. The results indicate that although the maximum dependency depth of most files reaches 4 levels, only the information provided by $D_1(f)$ files is the most useful. Furthermore, the effectiveness of using $D_\infty(f)$ surpasses that of *Random-All*, indicating that besides $D_1(f)$ files, there are many other useful files within the repository.

4.4 Cross-File Content Analysis

Definition 2. (Pruning Level) We define the pruning levels into three categories:

- **P-Level 0:** No pruning is applied to the file content.
- **P-Level 1:** All global context content is removed from the file.
- **P-Level 2:** All global context content, function bodies and class method bodies are removed from the file.

Table 3 presents the results of completion using cross-file information with different pruning levels. We can see that the results of *P-level:1*

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

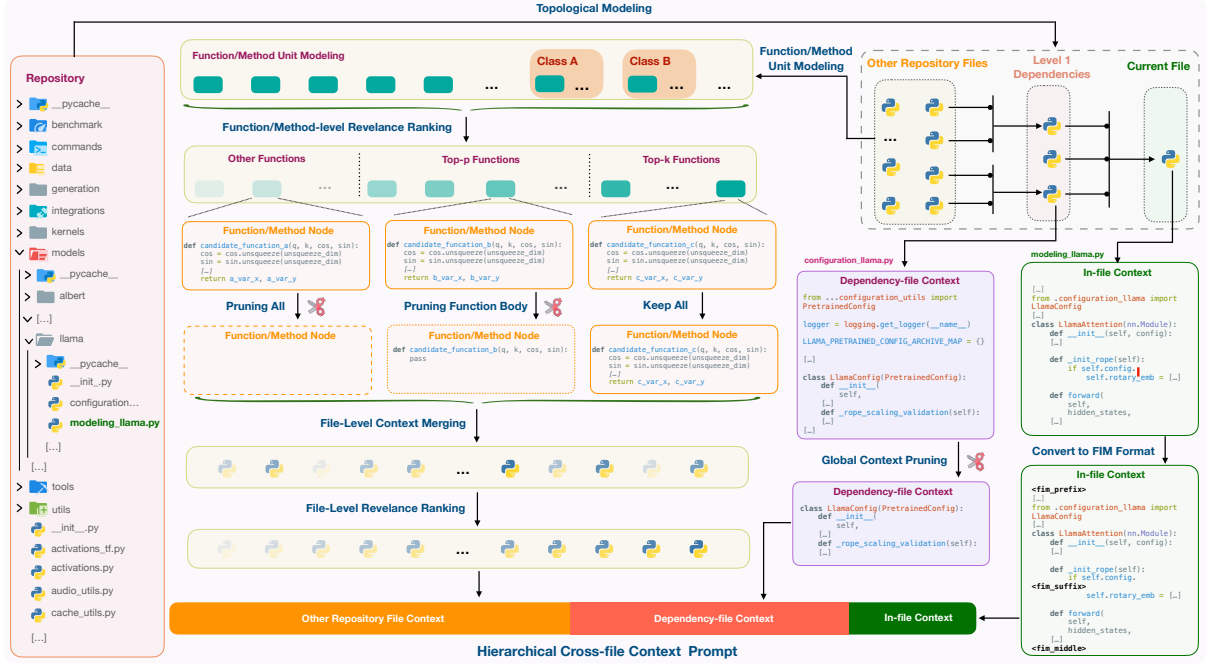


Figure 3: The framework of hierarchical context pruning for improving the performance of code large language models in real-world code completion tasks.

outperform those of $P\text{-level}:0$, indicating that the Global Context information from cross-file content has minimal impact on the completion of the current file. Additionally, the results of $P\text{-level}:2$ are only slightly worse than those of $D_\infty(f)$, and when combined with the information from $D_1(f)$, they are almost equivalent to the results of $D_\infty(f)$. This suggests that the specific implementations of most cross-file functions have minimal impact on the completion of the current file, and retaining only the function header information is sufficient.

5 Hierarchical Context Pruning

Based on the analysis results concerning the dependencies and content of the files, we attempt to construct a hierarchical context prompt based on the importance and relevance of the repository content. This approach aims to enhance the accuracy of code completion models while effectively reducing the length of the context. Figure 3 shows the specific process for constructing a hierarchical context prompt.

5.1 Fine-grained Repository Modeling

In order to precisely control the content within the code repository, we employ *Tree-Sitter* to parse the files within the repository. We model the content using three types of nodes:

- **Function Node:** Represents a function or a class method within a code file.
- **Class Node:** Represents a class in a code file, consisting of the class’s name, attributes, and Function Nodes.
- **File Node:** Represents a code file, comprising Nodes that represent the functions and classes within the file, along with global context information.

5.2 Hierarchical Context

As shown in the top right of Figure 3, following the settings in Section 4.3, we conduct a dependency analysis on the files in the repository. We perform a topological sort based on the dependency relationships, centering around the file currently being completed. According to the experimental results in Section 4.3, only files at dependency level 1 significantly enhance completion accuracy. Therefore, we select files designated as $D_1(f)$ to serve as dependency files. Ultimately, the files in the repository are categorized into three types: *current file*, *dependency files*, and *other files*. We will apply different strategies to optimize each type of file.

Current File. For the current file, any content within the file may be needed during completion,

XF-Context	Hierarchical Context Pruning (Top-p: 1.0)											
	DScoder-1.3B		DScoder-6.7B		StarCoder2-3B		StarCoder2-7B		CodeGemma-2B		CodeGemma-7B	
	EM	ES	EM	ES	EM	ES	EM	ES	EM	ES	EM	ES
Random-All	6.18	46.19	33.94	70.98	28.32	66.87	31.45	69.09	26.93	62.13	36.69	74.42
Top-k: 0	9.45	49.44	36.87	72.14	29.91	66.96	32.62	69.11	28.93	62.03	39.17	75.16
Top-k: 5	9.64	49.78	39.74	73.90	32.68	69.05	35.76	71.41	31.26	63.74	42.44	76.95
Top-k: 10	9.91	49.85	40.30	74.56	34.15	69.37	36.47	71.50	31.82	64.34	42.63	77.35

Table 4: The results of completion using hierarchical context pruning with different top-k values.

XF-Context	Hierarchical Context Pruning (Top-k: 5)											
	DScoder-1.3B		DScoder-6.7B		StarCoder2-3B		StarCoder2-7B		CodeGemma-2B		CodeGemma-7B	
	EM	ES	EM	ES	EM	ES	EM	ES	EM	ES	EM	ES
Random-All	6.18	46.19	33.94	70.98	28.32	66.87	31.45	69.09	26.93	62.13	36.69	74.42
Top-p: 0.1	14.27	53.94	37.85	73.11	32.99	68.75	34.16	70.43	29.19	62.09	40.98	76.26
Top-p: 0.2	13.52	53.20	38.04	73.13	33.15	68.59	34.84	70.40	29.72	62.32	40.94	76.25
Top-p: 0.3	12.88	52.60	38.49	73.19	32.84	68.31	35.22	70.64	30.13	62.77	41.21	76.20

Table 5: The results of completion using hierarchical context pruning with different top-p values.

so we retain all content of the file and convert it into the Fill-in-the-middle (FIM) format.

Dependency Files. According to the experimental results in Section 4.4, removing the global context across files does not affect the accuracy of completions. Therefore, for dependency files, we remove all global context from these files.

Other Files. We refer to files other than the current file and its direct dependency files, namely $\{F \setminus D_1(f)\} \setminus f$, collectively as other files. For the content in *other files*, we remove all global context, and then we employ **function-level** sampling and pruning methods to optimize the content of these files.

5.3 Function-level Sampling

In this study, we used OpenAI’s text-embedding API³ to embed each function (or class method) and query code snippet in the repository. We then used the pre-computed similarity of embeddings between the query and candidate functions (or class methods) as an indicator of relevance. We select the code from the current line of completion and the 10 lines before and after it as a query to find functions and class methods most relevant to the current completion content.

We implemented two sampling strategies (**top-k** and **top-p**) and designed distinct content pruning

strategies for the functions (or class methods) sampled under each strategy, see Section 5.4.

5.4 Function-level Pruning

According to the experimental results in Section 4.4, the global context from all non-current files and most of the function bodies (or class method bodies) within the code repository can be pruned. Appropriately pruning low-relevance content can significantly reduce the length of the prompt input to the model.

Let F denote the set of all functions and class methods in the repository, F_k represent the functions sampled using the top-k strategy, and F_p represent the functions sampled using the top-p strategy:

$$\begin{aligned} F_k &= \{f \mid f \in \text{Top}_k(F)\} \\ F_p &= \{f \mid f \in \text{Top}_p(F)\} \end{aligned} \quad (4)$$

where $F_k \subseteq F_p$. Content from functions and class methods not within the set $F_k \cup F_p$ was completely pruned.

Top-k Context Pruning. For functions (or class methods) within the set F_k , we retained their entire content.

Top-p Context Pruning. For functions (or class methods) in the set F_p but not in F_k , we prune their implementations and retained only their function headers (or class method headers).

³openai-text-embedding-ada-002

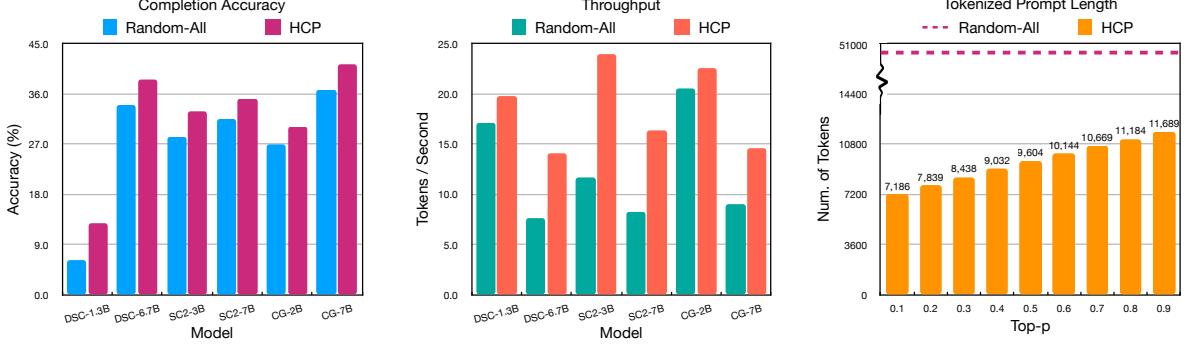


Figure 4: **left:** Comparison of completion results using random-all and the hierarchical context pruning across six models. **middle:** Comparison of throughput using random-all and the hierarchical context pruning across six models. **right:** Comparison of prompt length using random-all and the hierarchical context pruning of different top-p values (top-k=5).

5.5 File-level Relevance Ranking

Each function or class method in the repository has a similarity score. We assign different relevance weights to functions sampled using different sampling strategies.

$$W(f) = \begin{cases} 1.0, & \forall f \in F_k \\ 0.5, & \forall f \in F_p \setminus F_k \\ 0.0, & \forall f \in F \setminus (F_k \cup F_p) \end{cases} \quad (5)$$

where $\text{Top}_k(F)$ and $\text{Top}_p(F)$ represent the functions with the highest relevance scores sampled using the top-k and top-p strategies, respectively.

The similarity of a class is defined as the weighted sum of its class methods:

$$S(c) = \sum_{m \in c} W(m) * S(m) \quad (6)$$

where, c represents the class, and m represents the class method.

The similarity of a file is defined as the weighted sum of its functions and classes:

$$S(f) = \sum_{x \in \mathcal{F}} W(x) * S(x) + \sum_{c \in \mathcal{C}} S(c) \quad (7)$$

where, \mathcal{F} and \mathcal{C} represent the set of functions and classes in the file, respectively.

Finally, we sort the files at the file-level according to the relevance score to determine their relative positions in the prompt.

5.6 Experimental Results

We initially fixed top-p at 1.0 and tested the impact of different top-k values on completion accuracy. Table 4 presents some of the experimental results, while Table 11 in the Appendix E provides a more

comprehensive results. We observed that increasing the top-k value beyond 5 did not result in significant improvements in accuracy. Therefore, we conclude that a top-k value of 5 is sufficient.

We further fixed the top-k value at 5 and tested the impact of varying top-p values (ranging from 0.1 to 0.9) on completion accuracy. Partial experimental results are presented in Table 5, with more comprehensive results available in Table 12 in Appendix E. Our observations indicate that increasing the top-p value enhances completion accuracy; however, beyond a top-p value of 0.3, the improvement in accuracy slows considerably. Thus, we consider 0.3 to be a reasonable value.

Figure 4 visually compares the Hierarchical Context Pruning (HCP) strategy (top-k=5, top-p=0.3) with the method of randomly concatenating all repository code files across three dimensions: completion accuracy, throughput rate, and input length. The visualization shows that, compared to random concatenation, **HCP** significantly reduces input length (enhancing throughput) while improving the model’s completion accuracy.

6 Conclusion

In this study, we evaluated six Code LLMs pre-trained with repository-level code data. We conducted a detailed error analysis on these Code LLMs, performed topological dependency analysis on files within the code repositories, and analyzed the content of these files. Based on the results of these experiments, we proposed a strategy named Hierarchical Context Pruning to construct high-quality prompt inputs. Finally, we conducted experiments on six Repo-Code LLMs to verify the effectiveness of the proposed method.

Limitations

Benchmark. In this study, we utilized the Cross-CodeEval benchmark for evaluation. However, as demonstrated in the error analysis presented in Sections 4.2 and Appendix A, while the evaluation method based on exact matches is convenient and quick, it does not provide comprehensive results. Therefore, there may be a discrepancy between the evaluation outcomes and the actual capabilities of the model.

Function-level Sampling. In this study, sampling functions and class methods based on relevance required the use of a text embedding model. When the number of code files in the repository is excessive, this may reduce the sampling rate, leading to increased completion latency.

Ethical Statements

This study does not involve human participants, personal data, or hazardous materials, and primarily focuses on computational model performance. All resources used are open-source or properly licensed, ensuring compliance with relevant standards.

References

- Loubna Ben Allal, Raymond Li, Denis Kocetkov, Chenghao Mou, Christopher Akiki, Carlos Munoz Ferrandis, Niklas Muennighoff, Mayank Mishra, and Alex Gu et al. 2023. [Santacoder: don't reach for the stars!](#)
- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, and Charles Sutton. 2021. [Program synthesis with large language models.](#)
- Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao Liu, Aohan Zeng, Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. 2023. [Longbench: A bilingual, multi-task benchmark for long context understanding.](#)
- Mohammad Bavarian, Heewoo Jun, Nikolas Tezak, John Schulman, Christine McLeavey, Jerry Tworek, and Mark Chen. 2022. [Efficient training of language models to fill in the middle.](#)
- Zhangqian Bi, Yao Wan, Zheng Wang, Hongyu Zhang, Batu Guan, Fangxin Lu, Zili Zhang, Yulei Sui, Xuanhua Shi, and Hai Jin. 2024. [Iterative refinement of project-level code context for precise code generation with compiler feedback.](#)
- Federico Cassano, John Gouwar, Daniel Nguyen, Sydney Nguyen, Luna Phipps-Costin, Donald Pinckney, Ming-Ho Yee, Yangtian Zi, Carolyn Jane Anderson, Molly Q Feldman, Arjun Guha, Michael Greenberg, and Abhinav Jangda. 2022. [Multipl-e: A scalable and extensible approach to benchmarking neural code generation.](#)
- Longze Chen, Ziqiang Liu, Wanwei He, Yunshui Li, Run Luo, and Min Yang. 2024a. [Long context is not long at all: A prospector of long-dependency data for large language models.](#)
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, and et al. 2021. [Evaluating large language models trained on code.](#)
- Shouyuan Chen, Sherman Wong, Liangjian Chen, and Yuandong Tian. 2023. [Extending context window of large language models via positional interpolation.](#)
- Yukang Chen, Shengju Qian, Haotian Tang, Xin Lai, Zhijian Liu, Song Han, and Jiaya Jia. 2024b. [Longlora: Efficient fine-tuning of long-context large language models.](#)
- Yangruibo Ding, Zijian Wang, Wasi Uddin Ahmad, Hantian Ding, Ming Tan, Nihal Jain, Murali Krishna Ramanathan, Ramesh Nallapati, Parminder Bhatia, Dan Roth, and Bing Xiang. 2023. [Crosscodeeval: A diverse and multilingual benchmark for cross-file code completion.](#)
- Yiran Ding, Li Lyna Zhang, Chengruidong Zhang, Yuanxuan Xu, Ning Shang, Jiahang Xu, Fan Yang, and Mao Yang. 2024. [Longrope: Extending llm context window beyond 2 million tokens.](#)
- Richard Fang, Rohan Bindu, Akul Gupta, Qiusi Zhan, and Daniel Kang. 2024. [Llm agents can autonomously hack websites.](#)
- Daniel Fried, Armen Aghajanyan, Jessy Lin, Sida Wang, Eric Wallace, Freda Shi, Ruiqi Zhong, Wen tau Yih, Luke Zettlemoyer, and Mike Lewis. 2023. [InCoder: A generative model for code infilling and synthesis.](#)
- Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, Y. Wu, Y. K. Li, Fuli Luo, Yingfei Xiong, and Wenfeng Liang. 2024. [Deepseek-coder: When the large language model meets programming – the rise of code intelligence.](#)
- Dan Hendrycks, Steven Basart, Saurav Kadavath, Mantas Mazeika, Akul Arora, Ethan Guo, Collin Burns, Samir Puranik, Horace He, Dawn Song, and Jacob Steinhardt. 2021. [Measuring coding challenge competence with apps.](#)
- Samuel Holt, Max Ruiz Luyten, and Mihaela van der Schaar. 2024. [L2mac: Large language model automatic computer for extensive code generation.](#)

- Carlos E. Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. 2024. [Swe-bench: Can language models resolve real-world github issues?](#)
- Yuhang Lai, Chengxi Li, Yiming Wang, Tianyi Zhang, Ruiqi Zhong, Luke Zettlemoyer, Scott Wen tau Yih, Daniel Fried, Sida Wang, and Tao Yu. 2022. [Ds-1000: A natural and reliable benchmark for data science code generation.](#)
- Jia Li, Ge Li, Xuanming Zhang, Yihong Dong, and Zhi Jin. 2024a. [Evocodebench: An evolving code generation benchmark aligned with real-world code repositories.](#)
- Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou, Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, and et al. 2023. [Starcoder: may the source be with you!](#)
- Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, Thomas Hubert, Peter Choy, Cyprien de Masson d’Autume, Igor Babuschkin, Xinyun Chen, Po-Sen Huang, Johannes Welbl, Sven Gowal, Alexey Cherepanov, James Molloy, Daniel J. Mankowitz, Esme Sutherland Robson, Pushmeet Kohli, Nando de Freitas, Koray Kavukcuoglu, and Oriol Vinyals. 2022. [Competition-level code generation with alpha-code.](#) *Science*, 378(6624):1092–1097.
- Yunshui Li, Binyuan Hui, Xiaobo Xia, Jiayi Yang, Min Yang, Lei Zhang, Shuzheng Si, Ling-Hao Chen, Junhao Liu, Tongliang Liu, Fei Huang, and Yongbin Li. 2024b. [One-shot learning as instruction data prospector for large language models.](#)
- Ming Liang, Xiaoheng Xie, Gehao Zhang, Xunjin Zheng, Peng Di, wei jiang, Hongwei Chen, Chengpeng Wang, and Gang Fan. 2024. [Repofuse: Repository-level code completion with fused dual context.](#)
- Tianyang Liu, Canwen Xu, and Julian McAuley. 2023. [Repobench: Benchmarking repository-level code auto-completion systems.](#)
- Anton Lozhkov, Raymond Li, Loubna Ben Allal, Federico Cassano, Joel Lamy-Poirier, Nouamane Tazi, Ao Tang, Dmytro Pykhtar, Jiawei Liu, Yuxiang Wei, Tianyang Liu, Max Tian, Denis Kocetkov, and Arthur Zucker et al. 2024. [Starcoder 2 and the stack v2: The next generation.](#)
- Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xubo Geng, Wenxiang Hu, Chongyang Tao, Jing Ma, Qingwei Lin, and Daxin Jiang. 2023. [Wizardcoder: Empowering code large language models with evol-instruct.](#)
- Niklas Muennighoff, Qian Liu, Armel Zebaze, Qinkai Zheng, Binyuan Hui, Terry Yue Zhuo, Swayam Singh, Xiangru Tang, Leandro von Werra, and Shayne Longpre. 2024. [Octopack: Instruction tuning code large language models.](#)
- Erik Nijkamp, Hiroaki Hayashi, Caiming Xiong, Silvio Savarese, and Yingbo Zhou. 2023a. [Codegen2: Lessons for training llms on programming and natural languages.](#)
- Erik Nijkamp, Bo Pang, Hiroaki Hayashi, Lifu Tu, Huan Wang, Yingbo Zhou, Silvio Savarese, and Caiming Xiong. 2023b. [Codegen: An open large language model for code with multi-turn program synthesis.](#)
- Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. [Training language models to follow instructions with human feedback.](#)
- Huy N. Phan, Hoang N. Phan, Tien N. Nguyen, and Nghi D. Q. Bui. 2024. [Repohyper: Better context retrieval is all you need for repository-level code completion.](#)
- Nikhil Pinnaparaju, Reshith Adithyan, Duy Phung, Jonathan Tow, James Baicoianu, Ashish Datta, Maksym Zhuravinskyi, Dakota Mahan, Marco Bellagente, Carlos Riquelme, and Nathan Cooper. 2024. [Stable code technical report.](#)
- Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, Jérémy Rapin, and Artyom Kozhevnikov et al. 2024. [Code llama: Open foundation models for code.](#)
- Wenqi Shi, Ran Xu, Yuchen Zhuang, Yue Yu, Jieyu Zhang, Hang Wu, Yuanda Zhu, Joyce Ho, Carl Yang, and May D. Wang. 2024. [Ehrgent: Code empowers large language models for few-shot complex tabular reasoning on electronic health records.](#)
- Disha Shrivastava, Hugo Larochelle, and Daniel Tarlow. 2023. [Repository-level prompt generation for large language models of code.](#)
- Qiushi Sun, Zhirui Chen, Fangzhi Xu, Kanzhi Cheng, Chang Ma, Zhangyue Yin, Jianing Wang, Chengcheng Han, Renyu Zhu, Shuai Yuan, Qipeng Guo, Xipeng Qiu, Pengcheng Yin, Xiaoli Li, Fei Yuan, Lingpeng Kong, Xiang Li, and Zhiyong Wu. 2024. [A survey of neural code intelligence: Paradigms, advances and beyond.](#)
- CodeGemma Team, Ale Jakse Hartman, Andrea Hu, Christopher A. Choquette-Choo, Heri Zhao, Jane Fine, Jeffrey Hui, Jingyue Shen, Joe Kelley, Joshua Howland, Kshitij Bansal, Luke Vilnis, Mateo Wirth, Nam Nguyen, Paul Michel, Peter Choy, Pratik Joshi, Ravin Kumar, Sarmad Hashmi, Shubham Agrawal, Siqi Zuo, Tris Warkentin, and Zhitao Gong. 2024. [Codegemma: Open code models based on gemma.](#)
- Amitayush Thakur, George Tsoukalas, Yeming Wen, Jimmy Xin, and Swarat Chaudhuri. 2024. [An in-context learning agent for formal theorem-proving.](#)

- Xingyao Wang, Yangyi Chen, Lifan Yuan, Yizhe Zhang, Yunzhu Li, Hao Peng, and Heng Ji. 2024a. [Executable code actions elicit better llm agents](#).
- Yejie Wang, Keqing He, Guanting Dong, Pei Wang, Weihao Zeng, Muxi Diao, Yutao Mou, Mengdi Zhang, Jingang Wang, Xunliang Cai, and Weiran Xu. 2024b. [Dolphocoder: Echo-locating code large language models with diverse and multi-objective instruction tuning](#).
- Yizhong Wang, Hamish Ivison, Pradeep Dasigi, Jack Hessel, Tushar Khot, Khyathi Raghavi Chandu, David Wadden, Kelsey MacMillan, Noah A. Smith, Iz Beltagy, and Hannaneh Hajishirzi. 2023a. [How far can camels go? exploring the state of instruction tuning on open resources](#).
- Zhiruo Wang, Shuyan Zhou, Daniel Fried, and Graham Neubig. 2023b. [Execution-based evaluation for open-domain code generation](#).
- Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. 2024. [Efficient streaming language models with attention sinks](#).
- Yiheng Xu, Hongjin Su, Chen Xing, Boyu Mi, Qian Liu, Weijia Shi, Binyuan Hui, Fan Zhou, Yitao Liu, Tianbao Xie, Zhoujun Cheng, Siheng Zhao, Lingpeng Kong, Bailin Wang, Caiming Xiong, and Tao Yu. 2023. [Lemur: Harmonizing natural language and code for language agents](#).
- John Yang, Carlos E. Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik Narasimhan, and Ofir Press. 2024. [Swe-agent: Agent-computer interfaces enable automated software engineering](#).
- John Yang, Akshara Prabhakar, Karthik Narasimhan, and Shunyu Yao. 2023. [Intercode: Standardizing and benchmarking interactive coding with execution feedback](#).
- Pengcheng Yin, Wen-Ding Li, Kefan Xiao, Abhishek Rao, Yeming Wen, Kensen Shi, Joshua Howland, Paige Bailey, Michele Catasta, Henryk Michalewski, Alex Polozov, and Charles Sutton. 2022. [Natural language to code generation in interactive data science notebooks](#).
- Eric Zelikman, Eliana Lorch, Lester Mackey, and Adam Tauman Kalai. 2024. [Self-taught optimizer \(stop\): Recursively self-improving code generation](#).
- Fengji Zhang, Bei Chen, Yue Zhang, Jacky Keung, Jin Liu, Daoguang Zan, Yi Mao, Jian-Guang Lou, and Weizhu Chen. 2023a. [Repocoder: Repository-level code completion through iterative retrieval and generation](#).
- Lei Zhang, Yunshui Li, Ziqiang Liu, Jiayi Yang, Junhao Liu, and Min Yang. 2023b. [Marathon: A race through the realm of long context with large language models](#).
- Xinrong Zhang, Yingfa Chen, Shengding Hu, Zihang Xu, Junhao Chen, Moo Khai Hao, Xu Han, Zhen Leng Thai, Shuo Wang, Zhiyuan Liu, and Maosong Sun. 2024a. [\$\infty\$ bench: Extending long context evaluation beyond 100k tokens](#).
- Yakun Zhang, Wenjie Zhang, Dezhi Ran, Qihao Zhu, Chengfeng Dou, Dan Hao, Tao Xie, and Lu Zhang. 2024b. [Learning-based widget matching for migrating gui test cases](#). In *Proceedings of the 46th IEEE/ACM International Conference on Software Engineering, ICSE '24*. ACM.
- Tianyu Zheng, Ge Zhang, Tianhao Shen, Xueling Liu, Bill Yuchen Lin, Jie Fu, Wenhui Chen, and Xiang Yue. 2024. [Opencodeinterpreter: Integrating code generation with execution and refinement](#).

A Error Description and Analysis

In this section, we present detailed instances of various error types in model completions, accompanied by in-depth explanations and analyses of these errors. Figures 5-12 depict representatives for each error category. Each figure is bifurcated, with the left panel showing the output generated by the code model and the right panel presenting the corresponding ground truth. Errors in model completions are emphasized in *red italic* text, whereas the ground truth is denoted in *green italic*.

A.1 Redundant Content Generation

Redundant Content Generation means that the method is correctly called, but unnecessary additional content is generated. Figure 5 illustrates an example of a Redundant Content Generation error. The ground truth specifies `active is False`, yet the model's completion includes not only `active is False` but also additional irrelevant content.

Redundant Content Generation	
Completion	Groundtruth
<pre>[...] def test_client_agent_inactive(): client = Client(active=True, name="MyApp") assert client._config.options["active"] is True client.start() assert os.environ.get("_APPSIGNAL_ACTIVE") == "true" assert agent.active is False def test_client_agent_active_without_push_api_key(): [...]</pre>	<pre>[...] def test_client_agent_inactive(): client = Client(active=True, name="MyApp") assert client._config.options["active"] is True client.start() assert os.environ.get("_APPSIGNAL_ACTIVE") == "true" assert agent.active is False [...]</pre>

Figure 5: An example of redundant content generation error.

A.2 Partial Content Missing

Partial Content Missing indicates that the right method is called, but the generated content is incomplete, although this might still be acceptable to the user. Figure 6 presents an example of a Partial Content Missing error. The ground truth is `MiniGrid` and not `game_name.startswith('MiniGrid-')`, but the code completion model only managed to replicate a portion of this ground truth.

Partial Content Missing	
Completion	Groundtruth
<pre>[...] class TrainingConfig(): [...] def init_env_name(self, game_name, project_name): env_name = game_name self.env_source = EnvSrc.get_enum_env_src(self.env_source) if self.env_source == EnvSrc.MiniGrid: env_name = f'MiniGrid-{game_name}' [...]</pre>	<pre>[...] class TrainingConfig(): [...] def init_env_name(self, game_name, project_name): env_name = game_name self.env_source = EnvSrc.get_enum_env_src(self.env_source) if self.env_source == EnvSrc.MiniGrid and not game_name.startswith('MiniGrid-'): env_name = f'MiniGrid-{game_name}' [...]</pre>

Figure 6: An example of partial content missing error.

A.3 Parameter Value Error

The Parameter Value Error reflects the situation where the function call is correct, but the passed parameter values are incorrect. Figure 7 displays an instance of a Parameter Value Error. The code completion model correctly invokes the class method, but the parameters it employs differ from those specified in the ground truth.

A.4 Exact Match Error

Exact Match Error is a misjudgment due to the limitations of the exact match metric, such as using default values or specific strings when calling a function. Figure 8 illustrates an example of an Exact Match Error.

Parameter Value Error

Completion	Groundtruth
<pre>[...] llp.LLaMAConfig.set_model_path(input("Path to GGML LLaMA model weights: ")) model = Infilling(["Well, you see, every", " he", " to", " another", "!"]) # Run SMC for i,p in enumerate(llp.smc_steel(model, 8, 3)): print(f"Particle {i}: {p} (weight {p.weight})") [...]</pre>	<pre>[...] llp.LLaMAConfig.set_model_path(input("Path to GGML LLaMA model weights: ")) model = Infilling(["Well, you see, every", " he", " to", " another", "!"]) # Run SMC for i,p in enumerate(llp.smc_steel(model, 4,4)): print(f"Particle {i}: {p} (weight {p.weight})") [...]</pre>

Figure 7: An example of parameter value error.

The content completed by the code model is syntactically correct and semantically accurate, differing only slightly in textual terms from the ground truth. To avoid such misjudgments, a more reasonable evaluation method is necessary to assess the completion results.

Exact Match Error

Completion	Groundtruth
<pre>[...] def close_session(self): if self.session_handle is None: return LOGGER.debug(f"Closing session {self.session_name!r}") rc = self.CloseTrace(self.session_handle) if rc not in [winerror.ERROR_SUCCESS, winerror.ERROR_CTX_CLOSE_PENDING]: raise EtwConsumerException(f"CloseTrace failed [...]") self.session_handle = None [...]</pre>	<pre>[...] def close_session(self): if self.session_handle is None: return LOGGER.debug(f"closing session {self.session_name!r} ...") rc = self.CloseTrace(self.session_handle) if rc not in [winerror.ERROR_SUCCESS, winerror.ERROR_CTX_CLOSE_PENDING]: raise EtwConsumerException(f"CloseTrace failed [...]") self.session_handle = None [...]</pre>

Figure 8: An example of exact match error.

A.5 Non-existent Method Call

Non-existent Method Call indicates a call to a function, method, or property that does not exist. Figure 9 presents an example of a Non-existent Method Call error. The ground truth refers to a class method within the session class; however, the content generated by the code completion model erroneously calls a method that does not exist in the session class. This error can be regarded as a form of hallucination in the context of code completion.

Non-existent Method Call

Completion	Groundtruth
<pre>[...] @app.route("/api/userinput", methods=['POST']) def api_userinput(): data = request.get_json() user_input = data["user_input"] with generate_lock: result = Response(stream_with_context(session.generate(user_input, max_new_tokens = 1000)) return result [...]</pre>	<pre>[...] @app.route("/api/userinput", methods=['POST']) def api_userinput(): data = request.get_json() user_input = data["user_input"] with generate_lock: result = Response(stream_with_context(session.respond_multi(user_input)), mimetype = 'application/json') return result [...]</pre>

Figure 9: An example of non-existent method call error.

A.6 Improper Method Invocation

Improper Method Invocation represents the situation where the call is made to an existing method, but a different, more appropriate method should have been used. Figure 10 showcases an example of an Improper Method Invocation error. The code completion model generated a call to the class method

Transformer within the llp class, whereas the correct content should have invoked the class method Geometric within the same class.

Improper Method Invocation	
Completion	Groundtruth
<pre>[...] class Infilling(llp.Model): [...]</pre> <pre>def step(self): # Generate a token n = self.sample(llp.Transformer(self.ctx)) for _ in range(n): self.s += self.sample(llp.Transformer(self.ctx)) [...]</pre>	<pre>[...] class Infilling(llp.Model): [...]</pre> <pre>def step(self): # Generate a token n = self.sample(llp.Geometric(0.5)) + 1 for _ in range(n): self.s += self.sample(llp.Transformer(self.ctx)) [...]</pre>

Figure 10: An example of improper method invocation error.

A.7 Missing Method Invocation

Missing Method Invocation indicates that a function or method should have been called to achieve functionality, but the model failed to make this call. Figure 11 illustrates an example of a Missing Method Invocation error. The ground truth involves calling the class method paginate from the query class to obtain the queried variable. However, the code completion model failed to complete this method invocation and instead achieved the same functionality through multiple alternative class methods.

Missing Method Invocation	
Completion	Groundtruth
<pre>[...] class Page: [...]</pre> <pre>@staticmethod def set_iterator(client, set_query, map_lambda=None, mapper=None, page_size=None): def get_page(**kwargs): queried = query.set_(set_query) if kwargs: queried = query.filter_(queried, **kwargs) if map_lambda is not None:</pre> <pre>[...]</pre>	<pre>[...] class Page: [...]</pre> <pre>@staticmethod def set_iterator(client, set_query, map_lambda=None, mapper=None, page_size=None): def get_page(**kwargs): queried = query.paginate(set_query, **kwargs) if map_lambda is not None:</pre> <pre>[...]</pre>

Figure 11: An example of missing method invocation error.

A.8 Incorrect Content Generation

Incorrect Content Generation represents the situation where the generated content is illogical, irrelevant to the current code context, or completely incorrect. Figure 12 depicts an example of an Incorrect Content Generation error. The content produced by the code completion model is entirely unrelated to the ground truth and also lacks relevance to the current code context.

Incorrect Content Generation	
Completion	Groundtruth
<pre>[...] def test_envron_source(): [...]</pre> <pre>assert config.sources["environment"] == env_options final_options = Options() final_options.log_file_path = cwdir final_options.update(config.sources["system"]) final_options.update(env_options) assert config.options == final_options [...]</pre>	<pre>[...] def test_envron_source(): [...]</pre> <pre>assert config.sources["environment"] == env_options final_options = Options() final_options.update(config.sources["default"]) final_options.update(config.sources["system"]) final_options.update(env_options) assert config.options == final_options [...]</pre>

Figure 12: An example of incorrect content generation error.

B Special Tokens & Prompt Templates

Table 6 shows the special tokens used by DeepseekCoder, Starcoder2, and CodeGemma for fill-in-the-middle code completion. The prompt templates for DeepseekCoder, Starcoder2, and CodeGemma are shown in Table 7. Both Starcoder2 and CodeGemma utilize special tokens for segmenting code files, whereas DeepseekCoder does not employ such tokens, despite being trained on repository-level code data.

Model	Special Tokens
DeepseekCoder	< fim_begin >,< fim_hole >,< fim_end >
Starcoder2	<repo_name>,<file_sep>,<fim_pad>,<fim_prefix>,<fim_suffix>,<fim_middle>
CodeGemma	< file_separator >,< fim_prefix >,< fim_suffix >,< fim_middle >

Table 6: Special tokens used by DeepseekCoder, Starcoder2 and CodeGemma for fill-in-the-middle code completion.

Model	Fill-in-the-Middle Prompt Template
DeepseekCoder	#file_path0\ncode0\n#file_path1\ncode1\n#file_path2\ncode2\n#file_path3\n< fim_begin >prefix_code< fim_hole >suffix_code< fim_end >
Starcoder2	<repo_name>reponame<file_sep>file_path0\ncode0<file_sep>file_path1<fim_prefix>prefix_code<fim_suffix>suffix_code<fim_middle>
CodeGemma	< file_separator >file_path0\ncode0<file_separator>file_path1\n< fim_prefix >prefix_code< fim_suffix >suffix_code< fim_middle >

Table 7: Prompt templates for DeepseekCoder, Starcoder2 and CodeGemma.

C Prompt Length Distribution

Table 8 presents the average and median lengths of input sequences for three code completion models when utilizing contexts of varying dependency levels. Notably, Level ∞ , which incorporates the entire repository code into the input, results in an average input sequence length exceeding 50,000, far surpassing the context window supported by these models. To more visually observe the changes in input sequence length with respect to dependency levels, Figure 13 was created. It is evident that the median input sequence length begins to converge once the dependency level reaches 2, and the average input sequence length also starts to stabilize after reaching a dependency level of 3.

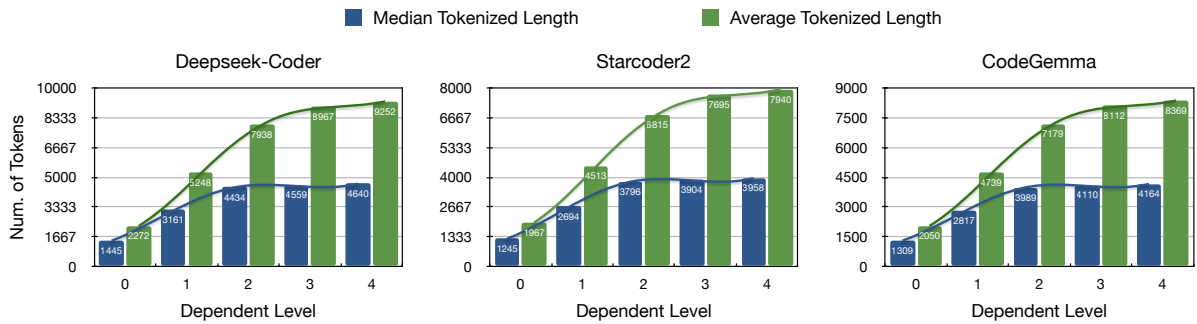


Figure 13: The distribution of tokenized prompt lengths in the CrossCodeEval benchmark. The x-axis represents the dependent level, and the y-axis represents the number of tokens. ■ denotes the median value of the tokenized prompt length. ■ denotes the average value of the tokenized prompt length.

Model	CrossCodeEval Benchmark: Python											
	Level 0		Level 1		Level 2		Level 3		Level 4		Level ∞	
	Median	Average	Median	Average	Median	Average	Median	Average	Median	Average	Median	Average
DeepseekCoder	1,445	2,272	3,161	5,248	4,434	7,938	4,559	8,967	4,640	9,252	44,475	58,217
StarCoder2	1,245	1,967	2,694	4,513	3,796	6,815	3,904	7,695	3,958	7,940	38,174	50,632
CodeGemma	1,309	2,050	2,817	4,739	3,989	7,179	4,110	8,112	4,164	8,369	39,647	52,875

Table 8: The median and average tokenized prompt lengths of the DeepseekCoder, StarCoder2 and CodeGemma models on the CrossCodeEval: Python benchmark.

D Dependency Level Analysis

D.1 Complete Experimental Results

Table 8 documents the comprehensive experimental results of repository file dependency analyses across six code completion models. It is observed that when the length of the input sequence exceeds the model’s context window, there is a significant decrease in completion accuracy. However, truncating the input sequence from the left to fit within the model’s context window size reveals that greater amounts of code repository content can enhance completion accuracy. Additionally, it was found that the DeepseekCoder-1.3B model exhibits a severe performance degradation in completion accuracy as the number of repository files increases.

Dependency	Topological Dependency Analysis											
	DSCoder-1.3B		DSCoder-6.7B		StarCoder2-3B		StarCoder2-7B		CodeGemma-2B		CodeGemma-7B	
	EM	ES	EM	ES	EM	ES	EM	ES	EM	ES	EM	ES
Dep-Level: 0	16.72	56.60	28.14	68.40	21.92	61.45	23.16	63.62	20.60	55.97	30.40	69.76
+ left truncate	16.72	56.58	28.14	68.36	21.92	61.49	22.98	63.58	20.64	56.26	30.58	70.36
Dep-Level: 1	14.99	54.33	32.20	68.57	26.33	64.54	28.66	67.00	23.16	55.00	32.17	65.77
+ left truncate	15.44	55.03	33.03	70.77	26.18	64.15	28.51	66.91	24.37	58.79	34.65	73.01
Dep-Level: 2	12.73	51.72	30.21	65.46	26.63	64.50	29.83	67.03	21.24	49.62	28.36	57.76
+ left truncate	13.63	53.45	33.56	70.74	26.70	64.58	29.45	67.03	25.31	59.27	35.67	73.26
Dep-Level: 3	12.28	50.90	28.93	63.67	26.74	64.52	29.42	66.58	20.30	47.64	27.16	55.66
+ left truncate	13.26	53.17	33.07	70.51	26.82	64.56	29.23	67.01	25.35	59.30	35.93	73.34
Dep-Level: 4	12.13	50.69	28.44	63.15	26.48	64.30	29.68	66.84	20.08	47.29	26.93	55.16
+ left truncate	13.37	53.20	33.22	70.57	26.59	64.46	29.53	67.07	25.54	59.42	36.12	73.54
Dep-Level: ∞	1.32	28.04	7.08	17.53	18.19	51.92	24.52	54.73	1.54	6.17	1.85	3.88
+ left truncate	5.76	46.22	35.29	71.51	30.43	67.34	33.03	69.57	29.08	62.91	39.32	75.35

Table 9: Comparison of completion results using different context dependency levels across 6 models. **EM** denotes Exact Match, and **ES** denotes Edit Similarity. ∞ denotes the prompt including all files in the repository. *+left truncate* denotes the prompt is truncated to the max context window of LLMs from the left.

D.2 Hit Count Changes

Table 10 collates the variations in correct and incorrect completions across six code completion models when input contexts of different dependency levels are used. It is evident that as the dependency level increases, the variations in the model’s completion results become more stable. This stability arises because the changes in the model’s input context diminish as the dependency level is elevated. This also indicates that augmenting the model’s input with additional content can enhance completion accuracy, albeit at the risk of turning some originally correct completions into incorrect ones.

We also observed that the DeepseekCoder series of models lack special tokens for delineating repository files; however, this deficiency does not result in more pronounced fluctuations in the outcomes. This suggests that the DeepseekCoder models are capable of effectively distinguishing between different files

in the repository, even without the aid of special tokens.

XF-Context	Hit Count Changes					
	DScoder-1.3B	DScoder-6.7B	Starcoder2-3B	Starcoder2-7B	CodeGemma-2B	CodeGemma-7B
Infile-Only	+444	+747	+582	+610	+548	+812
0 → 1	-108 +74	-47 +177	-44 +157	-37 +184	-31 +130	-68 +176
Level: 1	+408	+877	+695	+755	+647	+920
1 → 2	-61 +13	-33 +47	-41 +55	-33 +58	-30 +55	-44 +71
Level: 2	+362	+891	+709	+782	+672	+947
2 → 3	-15 +5	-20 +7	-13 +16	-19 +13	-10 +11	-11 +18
Level: 3	+352	+878	+712	+776	+673	+954
3 → 4	-3 +6	-1 +5	-10 +4	-3 +11	-3 +8	-5 +10
Level: 4	+355	+882	+706	+784	+678	+959
4 → ∞	-238 +36	-135 +190	-55 +157	-68 +161	-45 +139	-72 +157
Level: ∞	+153	+937	+808	+877	+772	+1044

Table 10: The changes in the hit counts of correct and incorrect completions across six code completion models when using different context dependency levels. The green values denote the number of test samples that were originally correct but became incorrect as the dependency level of the input context increased. The red values represent the number of test samples that were initially incorrect but became correct with the elevation of the input context’s dependency level.

XF-Context	Hierarchical Context Pruning (Top-p: 1.0)											
	DScoder-1.3B		DScoder-6.7B		Starcoder2-3B		Starcoder2-7B		CodeGemma-2B		CodeGemma-7B	
	EM	ES	EM	ES	EM	ES	EM	ES	EM	ES	EM	ES
Random-All	6.18	46.19	33.94	70.98	28.32	66.87	31.45	69.09	26.93	62.13	36.69	74.42
Top-k: 5	9.64	49.78	39.74	73.90	32.68	69.05	35.76	71.41	31.26	63.74	42.44	76.95
Top-k: 10	9.91	49.85	40.30	74.56	34.15	69.37	36.47	71.50	31.82	64.34	42.63	77.35
Top-k: 15	9.23	49.23	40.75	74.59	33.96	69.41	36.62	71.59	31.86	64.53	42.55	77.06
Top-k: 20	9.01	48.95	41.24	74.57	34.37	69.81	36.66	71.56	31.93	64.67	42.85	77.52
Top-k: 25	8.93	48.82	40.34	74.47	33.73	69.62	37.00	71.91	32.46	64.81	43.11	77.47
Top-k: 30	8.44	48.48	39.74	74.17	33.28	69.42	36.14	71.29	32.46	64.81	42.44	77.20

Table 11: The results of completion using hierarchical context pruning with different top-k values.

E Complete Function-Level Sampling Experiment Results

Due to space constraints, we report only a subset of the results from the function-level sampling experiments in the main body. Tables 11 and 12 provide a comprehensive statistical overview of the complete sampling experiments.

E.1 Top-k Sampling

Table 11 details the results of top-k sampling, where top-p is fixed at 1.0. It is observed that increasing the value of k does not significantly enhance the accuracy of completions; improvements become negligible when k exceeds 5.

E.2 Top-p Sampling

Table 12, on the other hand, presents the outcomes of top-p sampling with top-k fixed at 5. Here, increasing the value of p does not yield significant improvements, particularly when p exceeds 0.3.

XF-Context	Hierarchical Context Pruning (Top-k: 5)											
	DScoder-1.3B		DScoder-6.7B		StarCoder2-3B		StarCoder2-7B		CodeGemma-2B		CodeGemma-7B	
	EM	ES	EM	ES	EM	ES	EM	ES	EM	ES	EM	ES
Random-All	6.18	46.19	33.94	70.98	28.32	66.87	31.45	69.09	26.93	62.13	36.69	74.42
Top-p: 0.1	14.27	53.94	37.85	73.11	32.99	68.75	34.16	70.43	29.19	62.09	40.98	76.26
Top-p: 0.2	13.52	53.20	38.04	73.13	33.15	68.59	34.84	70.40	29.72	62.32	40.94	76.25
Top-p: 0.3	12.88	52.60	38.49	73.19	32.84	68.31	35.22	70.64	30.13	62.77	41.21	76.20
Top-p: 0.4	11.60	51.58	38.42	72.92	32.81	68.51	35.07	70.55	30.40	62.97	40.98	76.30
Top-p: 0.5	11.49	51.50	38.95	73.31	32.88	68.67	34.73	70.33	30.17	62.73	41.32	76.15
Top-p: 0.6	11.00	51.14	38.87	73.33	32.32	68.56	34.73	70.38	30.02	63.15	41.58	76.20
Top-p: 0.7	11.11	51.21	38.83	73.52	31.94	68.14	34.92	70.74	30.47	63.13	41.77	76.40
Top-p: 0.8	10.40	50.39	38.95	73.56	31.79	68.34	34.80	70.59	30.17	63.05	41.81	76.41
Top-p: 0.9	10.40	50.08	38.61	73.13	31.94	68.26	34.54	70.22	30.43	63.18	41.81	76.51

Table 12: The results of completion using hierarchical context pruning with different top-p values.

F Dependency Search Algorithm

Algorithm 1 delineates the specific process we employed for dependency modeling within code repositories. For more detailed implementation specifics, please visit our code repository.

G Additional Related Work

G.1 Long Context Code Large Language Models

Research on optimizing large language models for long contexts has been underway for some time, with many innovative long-context optimization techniques (Chen et al., 2023, 2024b; Xiao et al., 2024; Ding et al., 2024; Chen et al., 2024a) and evaluation sets (Bai et al., 2023; Zhang et al., 2023b, 2024a) being proposed and widely applied. Some Code LLMs (Guo et al., 2024; Lozhkov et al., 2024) utilize these techniques for fine-tuning to extend their context windows. Larger context windows allow Code LLMs to receive and process more complex code content, such as repository-level code completion and repository-level code repair (Ding et al., 2023; Liu et al., 2023; Zhang et al., 2024b; Li et al., 2024a).

G.2 Code Agent

Research on Code Agents (Yang et al., 2024; Fang et al., 2024; Thakur et al., 2024; Shi et al., 2024) focuses on developing intelligent systems that assist in software development by automating tasks like code generation and debugging (Holt et al., 2024; Wang et al., 2023b; Yin et al., 2022; Zelikman et al., 2024). The use of Code LLMs has proven effective in understanding complex code structures and semantics (Wang et al., 2024a; Yang et al., 2023). These models have been further refined to handle specific software development tasks, including repository-level code analysis and automated error correction (Jimenez et al., 2024).

Algorithm 1 Dependency Search Algorithm for Python Files

```
1: Input: file - initial Python file, maxDepth - maximum search depth
2: Output: List of dependent files
3: function FINDDEPENDENCIES(file, maxDepth)
4:     queue = [(file, 0)]
5:     visited = set()
6:     while queue do
7:         currentFile, currentDepth = queue.pop(0)
8:         if currentDepth > maxDepth then
9:             break
10:        end if
11:        imports = extractImports(currentFile)
12:        for imp in imports do
13:            if imp is local and imp not in visited then
14:                visited.add(imp)
15:                queue.append((imp, currentDepth + 1))
16:            end if
17:        end for
18:    end while
19:    return visited
20: end function
21: function EXTRACTIMPORTS(file)
22:     Use Tree-Sitter to parse file and extract all import statements
23:     return list of imports
24: end function
```
