

EXECREPOBENCH: Multi-level Executable Code Completion Evaluation

Anonymous ACL submission

Abstract

Code completion has become an essential tool for daily software development. Existing evaluation benchmarks often employ static methods that do not fully capture the dynamic nature of real-world coding environments and face significant challenges, including limited context length, reliance on superficial evaluation metrics, and potential overfitting to training datasets. In this work, we introduce a novel framework for enhancing code completion in software development through the creation of a repository-level benchmark EXECREPOBENCH and the instruction corpora REPO-INSTRUCT, aim at improving the functionality of open-source large language models (LLMs) in real-world coding scenarios that involve complex interdependencies across multiple files. EXECREPOBENCH include 1.2K samples from active Python repositories. Plus, we present a multi-level grammar-based completion methodology conditioned on the abstract syntax tree to mask code fragments at various logical units (e.g. statements, expressions, and functions). Then, we fine-tune the open-source LLM with 7B parameters on REPO-INSTRUCT to produce a strong code completion baseline model Qwen2.5-Coder-Instruct-C based on the open-source model. Qwen2.5-Coder-Instruct-C is rigorously evaluated against benchmarks, including MultiPL-E and EXECREPOBENCH, which consistently outperforms prior baselines across all programming languages. The deployment of Qwen2.5-Coder-Instruct-C can be used as a high-performance, local service for programming development¹.

1 Introduction

In the field of software engineering, the emergence of large language models (LLMs) designed specifically for code-related tasks has represented a significant advancement. These code LLMs (Li

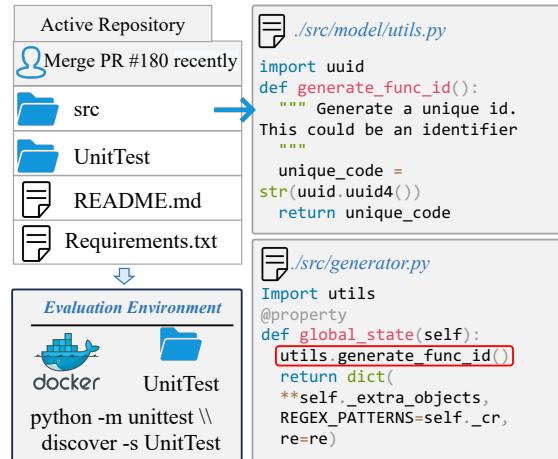


Figure 1: Executable Repository-level code evaluation with the given test cases.

et al., 2022; Allal et al., 2023), such as DeepSeekCoder (Guo et al., 2024a) and Qwen-Coder (Hui et al., 2024), have been pre-trained on extensive datasets comprising billions of code-related data. The advent of code LLMs has revolutionized the automation of software development tasks, providing contextually relevant code suggestions and facilitating code generation.

The code completion task holds paramount importance in modern software development, acting as a cornerstone for enhancing coding efficiency and accuracy. By analyzing the context of the ongoing work and using sophisticated algorithms to predict and suggest the next segments of code, code completion tools drastically reduce the time and effort programmers spend on writing boilerplate code, navigating large codebases, or recalling complex APIs and frameworks, which both accelerates the software development cycle and significantly diminishes the likelihood of syntax errors and bugs, leading to cleaner, more maintainable code. The recent code LLMs (Bavarian et al., 2022; Zheng et al., 2023) complete the middle code based on the prefix and suffix code through prefix-suffix-middle

¹The evaluation code and dataset will be released

(PSM) and suffix-prefix-middle (SPM) pre-training paradigm. To correctly evaluate the code completion capability, the HumanEval benchmark (Allal et al., 2023; Zheng et al., 2023) is extended to the infilling task by randomly masking some code spans and lines and prompting LLMs to predict the middle code. The recent works (Ding et al., 2023b, 2022, 2023a) propose to use the cross-file context to complete the current file and then score the results with n -gram string match. *However, the community still lacks an executable evaluation repository-level benchmark from live repositories and the corresponding instruction corpora.*

In this work, we benchmark, elicit, and enhance code repository-level completion tasks of open-source large language models (LLMs) by creating the repository-level instruction corpora **REPO-INSTRUCT** and the corresponding benchmark **EXECREPOBENCH** for utilization and evaluation for code completion in real-world software development scenarios, where projects frequently involve complex dependencies across multiple files. Unlike previous benchmarks with text-matching metrics (e.g. exact match (EM) and edit similarity (ES)), we EXECREPOBENCH is constructed with repository-level unit tests to verify the correctness of the completion code, which contains 1.2K samples from 50 active Python repositories. To facilitate the attention of the community for the code completion task, we propose the multi-level grammar-based completion to create REPO-INSTRUCT, where the code fragments under the different levels of logical units are masked for completion using the parsed abstract syntax tree (AST). During supervised finetuning (SFT), the code snippet of the repository is packed into the instruction data for the code completion LLMs Qwen2.5-Coder-Instruct-C, where the query gives the prefix code of the current file, suffix code of the current file, and code snippets of other files.

Qwen2.5-Coder-Instruct-C is evaluated on the code generation benchmark (Cassano et al., 2023) and our created code completion benchmark EXECREPOBENCH. The results demonstrate that Qwen2.5-Coder-Instruct-C consistently achieves state-of-the-art performance across all languages, notably surpassing the previous baselines. The contributions are summarized as follows:

- We introduce executable repository-level benchmark EXECREPOBENCH for code completion evaluation, which collects the active

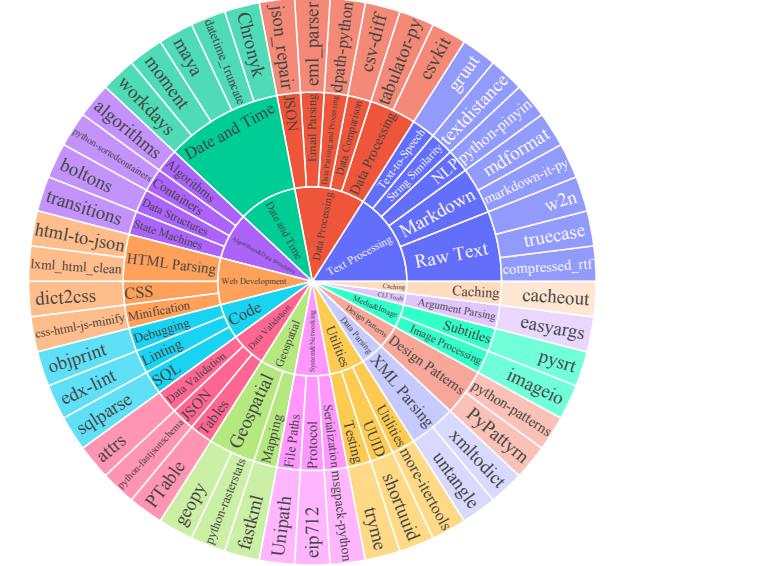


Figure 2: Classification of collected repositories.

repositories from GitHub and modify them into executable formats with test cases.

- We propose the multi-level grammar-based completion conditioned on the abstract syntax tree, where the statement-level, expression-level, function-level, and class-level code snippets are extracted for multi-level completion instruction corpora REPO-INSTRUCT
- Based on the open-source LLMs and the instruction corpora REPO-INSTRUCT, we fine-tune base LLMs with 7B parameters Qwen2.5-Coder-Instruct-C with a mixture of code completion data and standard instruction corpora, which can be used as a local service for programming developer.

2 EXECREPOBENCH Construction

Data Collection and Annotation The collected and refined repositories should follow the following guidelines: (1) Search Github code repositories of the Python language that have been continuously updated. (2) Given the collected repositories, the annotator should collect or create the test cases for evaluation. (3) All collected repositories should pass the test cases in a limited time for fast evaluation (< 2 minutes). In Table 2, we collect diverse repositories for comprehensive code completion evaluation. We feed the prefix tokens and suffix tokens of the current file with the context tokens into the LLM to predict the middle code tokens.

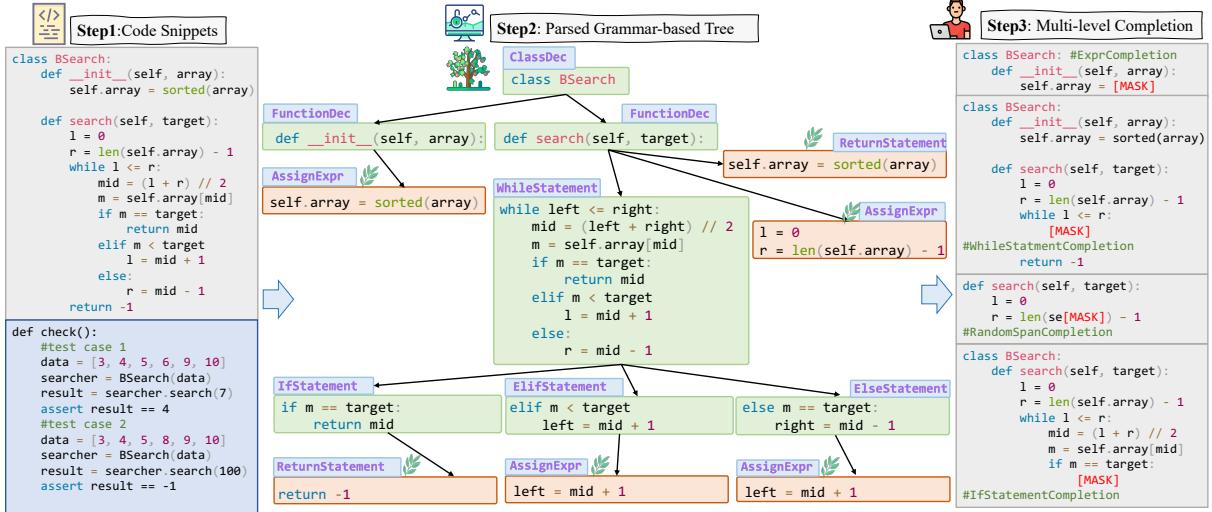


Figure 3: Multi-level Completion based on the parsed abstract syntax tree from the code snippet.

Decontainmation. To avoid data leakage, we remove exact matches (20-gram word overlap) from CrossCodeEval (Ding et al., 2023b) and the pre-training corpus stack V2 (Lozhkov et al., 2024).

Data Statistics To create the benchmark EXECREPOBENCH, we first construct the random span completion, random single-line completion, and random multi-line completion task by masking contiguous spans and lines of the chosen file of the whole repository. For the grammar-based completion, we first parse the code into an abstract syntax tree (AST) tree and randomly mask the node to match the input habits of programming developers habits. Besides, we sort the context files using the relevance between the current masked file and truncate the tokens exceeding the maximum supported length of the code LLM. The data statistic of EXECREPOBENCH is listed in Table 1.

3 Qwen2.5-Coder-Instruct-C

3.1 Problem Defintion

In-file Completion Given the code x_{L_k} of the current file of programming language L_k ($L_k \in L_{all} = \{L_{i=1}\}_{k=1}^K$), the LLM infills the middle code x_m conditioned on the prefix code x_p and the suffix code x_s as follow:

$$P(x_m^{L_k}) = P(x_m^{L_k} | x_p^{L_k}, x_s^{L_k}; \mathcal{M}) \quad (1)$$

where $x_p^{L_k}$, $x_s^{L_k}$, and $x_m^{L_k}$ are concatenated as the complete code to be executed with the given test cases to verify the correctness of the generated code.

Repository-level Completion Another more important completion scenario is the repository-level completion. Given the code snippets $z = \{z_{i=1}^{L_k}\}_{i=1}^N$ of N other files, the LLMs try to fill the part code of the current file. Based on the prefix code of the current file $x_p^{L_k}$, suffix code of the current file $x_s^{L_k}$, and the code snippets $z = \{z_{i=1}^{L_k}\}_{i=1}^N$ in other files in the same repository, the LLMs aim at producing the middle code x_m as:

$$P(x_m^{L_k}) = P(x_m^{L_k} | x_p^{L_k}, x_s^{L_k}, \{z_{i=1}^{L_k}\}_{i=1}^N; \mathcal{M}) \quad (2)$$

where the concatenation of $x_p^{L_k}$, $x_s^{L_k}$, and $x_m^{L_k}$ are used for repository-level execution for evaluation.

3.2 Multi-level Grammar-based Completion

Inspired by programming language syntax rules and user habits in practical scenarios, we leverage the tree-sitter-languages² to parse the code snippets and extract the basic logic blocks as the middle code to infill. For example, the abstract syntax tree (AST) represents the structure of Python code in a tree format, where each node in the tree represents a construct occurring in the source code. The tree's hierarchical nature reflects the syntactic nesting of constructs in the code, and includes various elements such as expressions, statements, and functions. By traversing and manipulating the AST, we can randomly extract the nodes of multiple levels and use the code context of the same file to uncover the masked node.

²<https://pypi.org/project/tree-sitter-languages/>

	Span	Random Completion		Grammar-based Completion		
		Single Line	Multiple Line	Expression	Statement	Function
Samples	42	34	38	407	266	377
Context Tokens	0/277.7K/27.7K	333/276.7K/22.6K	0/1484.1K/55.5K	0/1484.4K/36.2K	0/1484.6K/93.1K	0/1484.5K/65.1K
Prefix Tokens	0/32.2K/1.4K	0/38.3K/1.9K	0/7.2K/978.0	0/35.8K/1.5K	0/12.8K/786.0	0/38.3K/2.4K
Middle Tokens	3/22/7.0	4/48/15.0	4/156/40.0	2/150/13.0	2/74/9.0	7/123/33.0
Suffix Tokens	0/7.9K/924.0	0/2.0K/562.0	0/5.8K/935.0	1/39.0K/1.3K	1/40.1K/2.6K	1/37.8K/2.0K
Repository Overview	Repositories 50	Directories 2/115/15	Stars 1/39K/2.6K	Files 15/790/113	Python Files 4/411/38	Other Files 10/379/75

Table 1: Data statistics of EXECREPOBENCH.

Expression-level Completion At the expression level, we focus on completing sub-expressions within larger expressions or simple standalone expressions. This might involve filling in operand or operator gaps in binary operations or providing appropriate function arguments.

Statement-level Completion This level targets the completion of individual statements, such as variable assignments, control flow structures (if statements, for loops), and others. The goal is to maintain the logical flow and ensure syntactic correctness.

Function-level Completion At the function level, our approach involves completing entire function bodies or signature infillings. This includes parameter lists, return types, and the internal logic of the functions.

4 Heuristic Completion Techniques

To enhance the performance of our AST-based code infilling, we implement heuristic completion techniques to mimic the complementary habits of human users.

Random Line Completion We randomly select lines from the same file or similar files in the dataset to serve as candidates for completion. This process requires additional context-aware filtering to maintain relevance and accuracy.

Random Span Completion Instead of single lines, we randomly select code spans - sequences of lines that represent cohesive logical units. This approach suits larger blocks of code, needing a finer grasp of context and structure for effective completion.

4.1 Hybrid Instruction Tuning

Different from the base model trained with the FIM objective, we fine-tune the LLM with a mixture of the code completion data $(x_p^{L_k}, x_m^{L_k}, x_s^{L_k}) \in D_x = \{D_x^{L_k}\}_{k=1}^K$. The code completion training objective

is described as:

$$\mathcal{L}_c = -\frac{1}{K} \sum_{k=1}^K \mathbb{E}_{x_p^{L_k}, x_m^{L_k}, x_s^{L_k} \in D_x^{L_k}} [P(x_m^k | x_p, x_s, z; \mathcal{M})] \quad (3)$$

where the concatenation of $x_p^{L_k}$, $x_s^{L_k}$, and $x_m^{L_k}$ are used for repository-level execution for evaluation. $z = \{z_i^{L_k}\}_{i=1}^N$ are context code snippets of N files in the same repository.

We also adopt the standard instruction data $(q^{L_k}, a^{L_k}) \in D_{q,a} = \{D_{q,a}^{L_k}\}_{k=1}^K$. The question-answer instruction tuning on $D_{q,a}$ is calculated by:

$$\mathcal{L}_{qa} = -\frac{1}{K} \sum_{k=1}^K \mathbb{E}_{a^{L_k}, q^{L_k} \in D_{q,a}^{L_k}} \log P(q^{L_k} | a^{L_k}; \mathcal{M}) \quad (4)$$

where (q^{L_k}, a^{L_k}) are query and the corresponding response from the dataset $D_{x,y}$, including code generation, code summarization other code-related tasks.

We unify the capability of the code completion and question-answer in a single instruction model. The training objective of the hybrid instruction tuning is described as:

$$L_{all} = \mathcal{L}_c + \mathcal{L}_{qa} \quad (5)$$

where \mathcal{L}_c is the code completion objective and \mathcal{L}_{qa} is the question-answering objective.

5 Experiments

5.1 Code LLMs

We evaluate 30+ models with sizes ranging from 0.5B to 30B+ parameters for open-source code large language models and closed-source general LLMs. For general models, we evaluate GPTs (Brown et al., 2020; OpenAI, 2023) (GPT-3.5-Turbo, GPT4-o) and Claude series (Anthropic, 2023). For code models, we test CodeLlama (Rozière et al., 2023), StarCoder/StarCoder2 (Li et al., 2023; Lozhkov et al., 2024), CodeGeeX (Zheng et al., 2023), OpenCoder (Huang et al., 2024), Qwen-Coder (Hui et al., 2024), DeepSeekCoder (Guo et al., 2024a), CodeStral (MistralAI,

276 2024), Yi-Coder³, CodeGemma (Zhao et al., 2024),
277 and Granite-Coder (Mishra et al., 2024).

278 5.2 Implementation Details

279 We extract the repository-level code snippets from
280 the-stack-V2⁴ and filter the data with heuristic
281 rules (e.g. GitHub stars and file length). We keep
282 the mainstream programming language (Python, C-
283 sharp, Cpp, Java, Javascript, Typescript, Php) and
284 drop other long-tailed languages to obtain nearly
285 1.5M repositories. Finally, we obtain the instruc-
286 tion dataset REPO-INSTRUCT contains nearly 3M
287 completion samples. We fine-tune the open-source
288 base foundation LLM Qwen2.5-Coder on nearly
289 3M instruction samples used in Qwen2.5-Coder
290 (Hui et al., 2024) and code completion data (in-file
291 and cross-file completion data). Qwen2.5-Coder-
292 Instruct-C is fine-tuned on Megatron-LM⁵ with 64
293 NVIDIA H100 GPUs. The learning rate first in-
294 creases into 3×10^{-4} with 100 warmup steps and
295 then adopts a cosine decay scheduler. We adopt
296 the Adam optimizer (Kingma and Ba, 2015) with
297 a global batch size of 2048 samples, truncating
298 sentences to 32K tokens.

299 5.3 Evaluation Metrics

300 **Edit Similarity** We compare the generated code
301 and the ground-truth code using edit similarity (ES)
302 to report string-based scores.

303 **Pass@k** Similar to the in-file benchmark Hu-
304 manEval/MBPP, we employ the Pass@k met-
305 ric (Chen et al., 2021) based on the executable
306 results to get the reliability evaluation results. In
307 this work, we report the greedy Pass@1 score of all
308 LLMs with greedy inference for a fair comparison.

309 5.4 Evaluation Benchmarks

310 **EXECREPOBENCH** EXECREPOBENCH is cre-
311 ated with the repository-level unit tests to verify the
312 correctness of the completion code, comprised of
313 1.2K samples from 50 active Python repositories.
314 We separately report the ES score and Pass@1 in
315 the table.

316 **MultiPL-E** Since the mixture training of in-
317 struction samples and code completion samples,
318 Qwen2.5-Coder-Instruct-C also supports answer-
319 ing the code-related queries. We adopt MultiPL-E

(Cassano et al., 2023) for multilingual evaluation,
320 including 8 popular programming languages.
321

322 5.5 Main Results

323 **EXECREPOBENCH** Table 2 presents a compara-
324 tive analysis of various code completion mod-
325 els, highlighting their performance across differ-
326 ent metrics and parameter sizes. Code LLMs
327 (e.g. CodeLlama and StarCoder) are evaluated
328 across several completion tasks: random comple-
329 tion (span, single-line, multi-line), and grammar-
330 based completion (expression, statement, function).
331 Our proposed model Qwen2.5-Coder-Instruct-C,
332 significantly outperforms competing models in
333 all categories despite having only 7B parameters.
334 Compared to the base foundation model Qwen2.5-
335 Coder and DS-Coder, Qwen2.5-Coder-Instruct-C
336 enhanced by the multi-level grammar-based fine-
337 tuning achieves an impressive average score of
338 44.2, marking a substantial advancement in the
339 field of code completion technologies. From the
340 table, we can see that there exists a mismatch be-
341 tween the n-gram-based metric ES and execution-
342 based metric pass@1. Granite-Coder-8B gets a
343 good pass@1 score but a bad ES score, which em-
344 phasizes the importance of execution-based metric
345 pass@k for correctly evaluating the code comple-
346 tion capability of code LLMs. ES metric has its
347 own inherent flaws, where the score is calculated
348 by the comparison between the generated code and
349 ground-truth code.

350 **MultiPL-E** Table 3 showcases the evaluation re-
351 sults in terms of Pass@1 performance (%) across
352 various models on the MultiPL-E benchmark, fo-
353 cusing on different programming languages. The
354 comparison is categorically divided between propri-
355 etary models, like GPT-3.5 and GPT-4, and open-
356 source models, which include DS-Coder, Yi-Coder,
357 and Qwen2.5-Coder variants, among others. o1-
358 preview, a proprietary model, leads with an aver-
359 age of 85.3%, showcasing the difference in per-
360 formance capability between proprietary and open-
361 source models. The results highlight the effective-
362 ness of our method, particularly in optimizing per-
363 formance within the constraints of parameter size.
364 Notably, our method, Qwen2.5-Coder-Instruct-C,
365 with 7 billion parameters, outperforms other mod-
366 els in this parameter range across all listed pro-
367 gramming languages, achieving an average Pass@1
368 performance of 76.4%.

³<https://huggingface.co/01-ai/Yi-Coder-9B>

⁴<https://huggingface.co/datasets/bigcode/the-stack-v2>

⁵<https://github.com/NVIDIA/Megatron-LM>

Models	Params	Random Completion						Grammar-based Completion						Avg.	
		Span		Single-line		Multi-line		Expression		Statement		Function			
		ES	Pass@1	ES	Pass@1	ES	Pass@1	ES	Pass@1	ES	Pass@1	ES	Pass@1	ES	Pass@1
Code-Llama	7B	3.7	11.9	6.8	35.3	17.2	26.3	5.8	28.5	5.9	23.3	17.3	15.6	9.9	22.7
Code-Llama	13B	3.4	19.0	6.5	35.3	16.7	26.3	5.7	29.5	6.3	25.6	17.4	17.5	9.9	24.4
Code-Llama	34B	4.5	9.5	6.5	32.4	16.9	18.4	6.7	28.7	6.5	22.9	17.8	16.7	10.5	22.6
Code-Llama	70B	3.9	16.7	6.8	38.2	17.7	26.3	5.8	28.7	6.1	25.6	17.6	19.9	10.0	24.9
Codestral	22B	3.5	16.7	9.0	41.2	18.1	28.9	5.7	27.0	6.1	24.8	17.4	16.7	10.0	23.3
StarCoder	1B	31.9	23.8	23.5	41.2	34.8	21.1	19.8	36.6	27.5	25.6	22.8	19.6	23.6	27.7
StarCoder	3B	16.6	11.9	11.7	41.2	24.6	26.3	12.2	28.3	18.0	26.7	19.6	15.1	16.5	23.4
StarCoder	7B	41.6	31.0	30.5	47.1	40.2	28.9	27.6	43.5	31.9	35.3	29.8	22.0	30.3	33.8
StarCoder2	3B	2.9	14.3	5.8	38.2	15.1	21.1	4.8	26.3	5.3	21.4	15.0	15.1	8.5	21.3
StarCoder2	7B	3.1	21.4	5.6	35.3	15.1	23.7	4.9	25.8	5.2	22.9	14.9	15.1	8.5	21.7
StarCoder2	15B	2.7	19.0	5.6	41.2	14.8	23.7	4.9	27.5	5.2	23.7	15.1	15.6	8.5	22.8
DS-Coder	1.3B	35.0	21.4	25.7	47.1	32.1	28.9	27.2	32.9	14.9	27.1	27.6	14.1	24.9	25.3
DS-Coder	6.7B	40.2	28.6	41.0	50.0	47.5	39.5	45.7	37.3	36.1	33.8	45.9	15.1	43.3	29.5
DS-Coder	33B	47.1	33.3	52.0	64.7	50.1	44.7	46.3	40.8	37.8	37.2	49.5	17.0	45.7	32.8
DS-Coder-V2-Lite	2.4/16B	33.0	31.0	44.1	52.9	42.5	39.5	42.4	37.6	30.0	32.3	42.7	16.2	39.4	29.7
Granite-Coder	3B	35.0	21.4	25.7	47.1	32.1	28.9	27.2	32.9	14.9	27.1	27.6	14.1	24.9	25.3
Granite-Coder	8B	0.0	19.0	0.0	58.8	2.6	28.9	5.2	36.6	0.0	26.3	0.0	21.5	1.9	29.1
Granite-Coder	20B	3.2	16.7	7.8	35.3	14.9	23.7	5.2	26.3	5.7	21.4	15.8	15.1	9.1	21.4
Granite-Coder	34B	3.1	16.7	8.0	35.3	15.2	26.3	5.3	26.3	6.2	24.1	15.4	15.6	9.1	22.3
CodeQwen1.5	7B	13.5	16.7	13.8	41.2	17.6	26.3	9.8	26.0	11.3	24.4	14.9	13.0	12.3	21.6
Qwen2.5-Coder	0.5B	11.3	16.7	10.4	47.1	13.0	26.3	10.7	26.0	14.6	24.1	16.2	14.1	13.5	22.0
Qwen2.5-Coder	1.5B	3.5	14.3	3.2	29.4	7.9	15.8	4.0	21.9	3.5	16.9	9.1	11.7	5.6	17.2
Qwen2.5-Coder	3B	13.8	19.0	14.9	44.1	12.5	21.1	13.9	28.0	11.2	23.7	18.5	13.5	14.8	22.3
Qwen2.5-Coder	7B	7.8	16.7	10.2	35.3	12.4	23.7	5.3	24.3	8.1	21.1	11.0	12.5	8.3	19.8
Qwen2.5-Coder	14B	7.0	16.7	12.7	35.3	12.1	18.4	11.4	27.5	15.1	27.8	18.5	13.5	14.4	22.6
Qwen2.5-Coder	32B	3.9	16.7	32.5	47.1	20.3	23.7	21.2	29.5	26.1	33.5	33.0	15.4	25.8	25.7
OpenCoder	1.5B	1.7	11.9	3.4	38.2	5.4	23.7	3.2	26.3	3.2	27.4	6.6	15.4	4.3	22.8
OpenCoder	8B	2.7	14.3	4.4	32.4	10.8	21.1	4.0	29.5	3.7	24.1	7.8	16.7	5.3	23.4
Yi-Coder	1.5B	3.9	16.7	6.6	32.4	16.4	28.9	6.1	25.3	6.4	26.3	17.2	14.1	10.0	21.9
Yi-Coder	9B	3.4	16.7	6.8	29.4	17.7	26.3	5.8	28.0	6.3	26.3	17.6	17.8	10.1	23.9
CodeGemma	2B	21.3	21.4	23.5	38.2	19.8	18.4	24.1	23.6	28.3	21.1	23.4	12.5	24.6	19.6
CodeGemma	7B	12.4	19.0	14.3	35.3	28.2	26.3	15.4	29.2	18.7	36.5	25.8	18.3	19.8	27.1
Qwen2.5-Coder-Instruct-C	7B	75.8	38.1	68.0	41.2	60.2	28.9	76.4	58.7	78.7	45.5	63.9	30.2	72.1	44.2

Table 2: Completion evaluation Results of base foundation model on EXECREPOBENCH.

6 Analysis

Ablation Study Figure 4 emphasizes the essence of each component in our method by conducting the ablation study. CrossCodeEval (Ding et al., 2023b) is developed using a variety of real-world, openly available repositories with permissive licenses, covering four widely used programming languages: Python, Java, TypeScript, and C-sharp. Figure 4(a) shows the model results of the code completion task CrossCodeEval and Figure 4(b) plots the results on the instruction-following code benchmark MultiPL-E. By unifying the code generation and completion in the same model, Qwen2.5-Coder-Instruct-C can support multiple scenarios.

Case Study Figure 5 showcases a part of a Python module named BankOperation which focuses on simulating basic bank account operations. The module, assumed to be spread across files, includes the BankAccount class definition housed within the given code snippet. Within this class, methods are defined for initializing an account (`__init__`), depositing money (`deposit`), and displaying the account balance (`display_balance`). The

core segment provided adds a `withdraw` method to this class, which allows for deducting a specified amount from the account’s balance if the amount is positive and does not exceed the available balance. Each transaction (`deposit` and `withdrawal`) is followed up with a call to `A.sync()`, hinting at an operation to synchronize the current state of the account with a database, potentially managed by code within `A.py`. Error handling is incorporated within the `deposit` and `withdrawal` operations to ensure amounts are valid. The description wraps up this modular approach to implementing a banking system in Python, emphasizing object-oriented programming principles, error management, and database integration. This example shows that Qwen2.5-Coder-Instruct-C can successfully find the dependency from the context file.

7 Related Work

Code Large Language Models In software engineering, the advent of large language models (LLMs) tailored for code-centric tasks has proven to be transformative. Models (Feng et al., 2020; Chen et al., 2021; Scao et al., 2022; Li et al.,

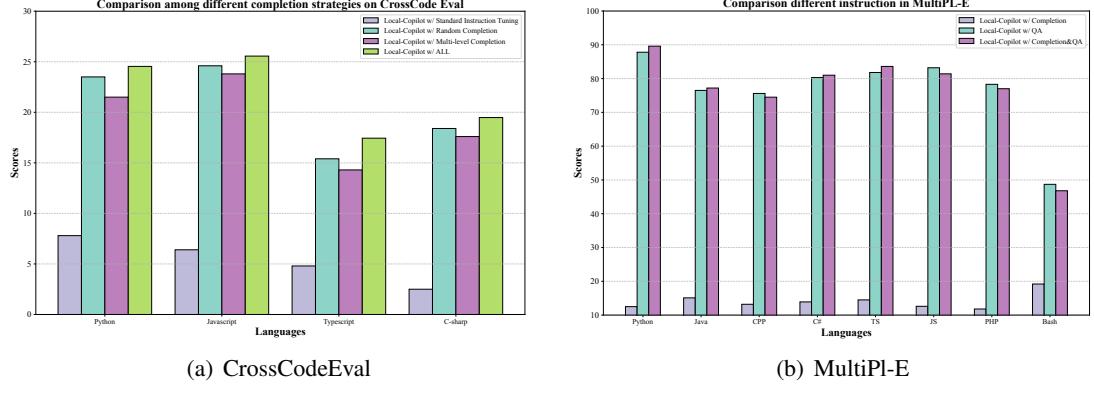


Figure 4: Evaluation results based on standard QA pairs and code completion.

Model	Size	HE	HE+	MBPP	MBPP+	Python	Java	C++	C#	TS	JS	PHP	Bash	Avg.
Closed-APIs														
Claude-3.5-Sonnet-20240620	🔒	89.0	81.1	87.6	72.0	89.6	86.1	82.6	85.4	84.3	84.5	80.7	48.1	80.2
Claude-3.5-Sonnet-20241022	🔒	92.1	86.0	91.0	74.6	93.9	86.7	88.2	87.3	88.1	91.3	82.6	52.5	83.8
GPT-4o-mini-2024-07-18	🔒	87.8	84.8	86.0	72.2	87.2	75.9	77.6	79.7	79.2	81.4	75.2	43.7	75.0
GPT-4o-2024-08-06	🔒	92.1	86.0	86.8	72.5	90.9	83.5	76.4	81.0	83.6	90.1	78.9	48.1	79.1
o1-mini	🔒	97.6	90.2	93.9	78.3	95.7	90.5	93.8	77.2	91.2	92.5	84.5	55.1	85.1
o1-preview	🔒	95.1	88.4	93.4	77.8	96.3	88.0	91.9	84.2	90.6	93.8	90.1	47.5	85.3
0.5B+ Models														
Qwen2.5-Coder-0.5B-Instruct	0.5B	61.6	57.3	52.4	43.7	61.6	57.3	52.4	43.7	50.3	50.3	52.8	27.8	49.6
1B+ Models														
DS-Coder-1.3B-Instruct	1.3B	65.9	60.4	65.3	54.8	65.2	51.9	45.3	55.1	59.7	52.2	45.3	12.7	48.4
Yi-Coder-1.5B-Chat	1.5B	69.5	64.0	65.9	57.7	67.7	51.9	49.1	57.6	57.9	59.6	52.2	19.0	51.9
Qwen2.5-Coder-1.5B-Instruct	1.5B	70.7	66.5	69.2	59.4	71.2	55.7	50.9	64.6	61.0	62.1	59.0	29.1	56.7
3B+ Models														
Qwen2.5-Coder-3B-Instruct	3B	84.1	80.5	73.6	62.4	83.5	74.7	68.3	78.5	79.9	75.2	73.3	43.0	72.1
6B+ Models														
CodeLlama-7B-Instruct	7B	40.9	33.5	54.0	44.4	34.8	30.4	31.1	21.6	32.7	-	28.6	10.1	-
DS-Coder-6.7B-Instruct	6.7B	74.4	71.3	74.9	65.6	78.6	68.4	63.4	72.8	67.2	72.7	68.9	36.7	66.1
CodeQwen1.5-7B-Chat	7B	83.5	78.7	77.7	67.2	84.1	73.4	74.5	77.8	71.7	75.2	70.8	39.2	70.8
Yi-Coder-9B-Chat	9B	82.3	74.4	82.0	69.0	85.4	76.0	67.7	76.6	72.3	78.9	72.1	45.6	71.8
DS-Coder-V2-Lite-Instruct	2.4/16B	81.1	75.6	82.8	70.4	81.1	76.6	75.8	76.6	80.5	77.6	74.5	43.0	73.2
Qwen2.5-Coder-7B-Instruct	7B	88.4	84.1	83.5	71.7	87.8	76.5	75.6	80.3	81.8	83.2	78.3	48.7	76.5
OpenCoder-8B-Instruct	8B	83.5	78.7	79.1	69.0	83.5	72.2	61.5	75.9	78.0	79.5	73.3	44.3	71.0
13B+ Models														
CodeLlama-13B-Instruct	13B	40.2	32.3	60.3	51.1	42.7	40.5	42.2	24.0	39.0	-	32.3	13.9	-
StarCoder2-15B-Instruct-v0.1	15B	67.7	60.4	78.0	65.1	68.9	53.8	50.9	62.7	57.9	59.6	53.4	24.7	54.0
Qwen2.5-Coder-14B-Instruct	14B	89.6	87.2	86.2	72.8	89.0	79.7	85.1	84.2	86.8	84.5	80.1	47.5	79.6
20B+ Models														
CodeLlama-34B-Instruct	34B	48.2	40.2	61.1	50.5	41.5	43.7	45.3	31.0	40.3	-	36.6	19.6	-
CodeStral-22B-v0.1	22B	81.1	73.2	78.2	62.2	81.1	63.3	65.2	43.7	68.6	-	68.9	42.4	-
DS-Coder-33B-Instruct	33B	81.1	75.0	80.4	70.1	79.3	73.4	68.9	74.1	67.9	73.9	72.7	43.0	69.2
CodeLlama-70B-Instruct	70B	72.0	65.9	77.8	64.6	67.8	58.2	53.4	36.7	39.0	-	58.4	29.7	-
DS-Coder-V2-Instruct	21/236B	85.4	82.3	89.4	75.1	90.2	82.3	84.8	82.3	83.0	84.5	79.5	52.5	79.9
Qwen2.5-Coder-32B-Instruct	32B	92.7	87.2	90.2	75.1	92.7	80.4	79.5	82.9	86.8	85.7	78.9	48.1	79.4
Qwen2.5-32B-Instruct	32B	87.8	82.9	86.8	70.9	88.4	80.4	81.0	74.5	83.5	82.4	78.3	46.8	76.9
Qwen2.5-72B-Instruct	32B	85.4	79.3	90.5	77.0	82.9	81.0	80.7	81.6	81.1	82.0	77.0	48.7	75.1
Qwen2.5-SynCoder	32B	92.7	87.8	86.2	74.7	92.1	80.4	80.7	81.6	83.0	85.7	77.6	49.4	78.8
Qwen2.5-Coder-Instruct-C	7B	87.2	81.1	81.7	68.5	89.6	77.2	74.5	81.0	83.6	81.4	77.0	46.8	76.4

Table 3: The performance of different instruction LLMs on EvalPlus and MultiPL-E. ‘‘HE’’ denotes the HumanEval, ‘‘HE+’’ denotes the plus version with more test cases, and ‘‘MBPP+’’ denotes the plus version with more test cases.

2022; Allal et al., 2023; Fried et al., 2022; Wang et al., 2021; Zheng et al., 2024; Jiang et al., 2024; Nijkamp et al., 2023; Wei et al., 2023;

Zhao et al., 2024) like CodeLlama (Rozière et al., 2023), DeepSeek-Coder (Guo et al., 2024a), OpenCoder (Huang et al., 2024) and Qwen2.5-Coder

```

[System Prompt]
You are a helpful code completion assistant.

[Question]
##Repo Name##:
BankOperation
##Repo-level Context Code##:
{A.py}
##Prefix Code##:
Import A
class BankAccount:
# Initialization method
    def __init__(self, owner, balance=0):
        """Initializes the account owner."""
        self.owner = owner
        self.balance = balance
    def deposit(self, amount):
        """Increases the account balance."""
        if amount > 0:
            self.balance += amount
            A.sync() #sync database
            return f"current balance is: {self.balance}"
        else:
            return "must be greater than zero."
##Suffix Code##:
def display_balance(self):
    """Displays the current account balance."""
    return f"current account balance is: {self.balance}"

[Answer]
##Middle Code##:
def withdraw(self, amount):
    """Decreases the account balance."""
    if 0 < amount <= self.balance:
        self.balance -= amount
        A.sync() #sync database
        return f"current balance is: {self.balance}"
    else:
        return "The withdrawal amount is invalid."

```

Figure 5: Example of Qwen2.5-Coder-Instruct-C.

Figure 6: A completion example in instruction format of Qwen2.5-Coder-Instruct-C.

(Hui et al., 2024) — all trained on vast corpuses comprising billions of code snippets — have fundamentally augmented the development process. These Code LLMs are instrumental in automating repetitive software tasks, proposing code improvements, and facilitating the conversion of natural language into executable code. Notable among these are Starcoder (Li et al., 2023; Lozhkov et al., 2024) and Qwen2.5-Coder (Hui et al., 2024), each bringing unique contributions to the enhancement of coding assistance tools. Inspired by the success of the grammar-based parsed tree in many fields, we adopt the abstract syntax tree to augment the code completion training.

Repo-level Code Evaluation In the domain of code evaluation, a rich tapestry of benchmarks (Zheng et al., 2023; Yu et al., 2024; Yin et al., 2023; Peng et al., 2024; Khan et al., 2023; Guo et al., 2024b; Lai et al., 2023) has been woven to address the challenges of accurately assessing code quality and functionality, such as HumanEval/MBPP (Chen et al., 2021; Austin et al.,

2021), their upgraded version EvalPlus (Liu et al., 2023a), and the multilingual benchmark MultiPLE (Cassano et al., 2023), McEval (Chai et al., 2024), and MdEval (Liu et al., 2024b). Realistic scenarios (Liu et al., 2024c,a), such as BigCodeBench (Zhuo et al., 2024), CodeArena (Yang et al., 2024) and SAFIM (Gong et al., 2024), separately evaluate code LLMs for more diverse scenarios and code preferences. Studies have explored a variety of approaches, ranging from static analysis techniques (e.g. exact match (EM) and edit similarity (ES)), which examine code without executing it, to dynamic methods that involve code execution in controlled environments (e.g. Pass@k). The current benchmarks support code models to evaluate a series of different types of tasks, such as code repair (Lin et al., 2017; Tian et al., 2024; Jimenez et al., 2023) and multilingual scenarios (Wang et al., 2023; Athiwaratkun et al., 2023; Zheng et al., 2023; Peng et al., 2024; Orlanski et al., 2023). An important task FIM (Fried et al., 2022; Bavarian et al., 2022; Ding et al., 2024) is to fill the middle code, given the prefix and suffix code, which provides substantial assistance for software development. Repo-level completion, such as RepoEval (Zhang et al., 2023), CrossCodeEval (Ding et al., 2023b; Wu et al., 2024) and RepoBench (Liu et al., 2023b) only using EM and ES without code execution can not accurately reflect the model performance.

8 Conclusion

In this work, we represent a significant leap forward in the realm of code completion, driven by the advancements in large language models (LLMs) tailored for coding tasks. By introducing an executable repository-level benchmark EXECREPOBENCH and a multi-level grammar-based instruction corpora REPO-INSTRUCT, we both tackles the limitations of existing benchmarks and set a new standard for evaluating code completion tools in real-world software development scenarios, where the EXECREPOBENCH is collected from real-world repositories. The fine-tuning of base LLMs with 7B parameters, culminating in Qwen2.5-Coder-Instruct-C, demonstrates a remarkable improvement in code completion accuracy and efficiency across various programming languages, outperforming previous LLMs. We pave the way for more accurate and context-aware code completion, promising to enhance developer productivity and make software development more efficient.

493 9 Limitations

494 We acknowledge the following limitations of this
495 study: (1) The evaluation in repository-level mul-
496 tilingual scenarios are not fully explored. (2) The
497 code completion model Qwen2.5-Coder-Instruct-C
498 is mainly supervised fine-tuned on the 7B open-
499 source base LLMs. In the future, we will try the
500 (3) The fine-tuned model can be further improved
501 using RLHF for better user experience, such as
502 DPO.

503 Ethics Statement

504 This research adheres to ethical guidelines for AI
505 development. We aim to enhance the capabilities
506 of large language models (LLMs) while acknowl-
507 edging potential risks such as bias, misuse, and
508 privacy concerns. To mitigate these, we advocate
509 for transparency, rigorous bias testing, robust se-
510 curity measures, and human oversight in AI appli-
511 cations. Our goal is to contribute positively to the
512 field and to encourage responsible AI development
513 and deployment.

514 References

515 Loubna Ben Allal, Raymond Li, Denis Kocetkov,
516 Chenghao Mou, Christopher Akiki, Carlos Munoz
517 Ferrandis, Niklas Muennighoff, Mayank Mishra,
518 Alex Gu, Manan Dey, et al. 2023. **SantaCoder: Don’t**
519 **reach for the stars!** *arXiv preprint arXiv:2301.03988*.

520 Anthropic. 2023. **Introducing Claude**.

521 Ben Athiwaratkun, Sanjay Krishna Gouda, Zijian Wang,
522 Xiaopeng Li, Yuchen Tian, Ming Tan, Wasi Uddin
523 Ahmad, Shiqi Wang, Qing Sun, Mingyue Shang, Su-
524 jan Kumar Gonugondla, Hantian Ding, Varun Ku-
525 mar, Nathan Fulton, Arash Farahani, Siddhartha Jain,
526 Robert Giaquinto, Haifeng Qian, Murali Krishna
527 Ramanathan, and Ramesh Nallapati. 2023. **Muli-**
528 **lingual evaluation of code generation models**. In *The*
529 *Eleventh International Conference on Learning Re-*
530 *presentations, ICLR 2023, Kigali, Rwanda, May 1-5,*
531 *2023*. OpenReview.net.

532 Jacob Austin, Augustus Odena, Maxwell Nye, Maarten
533 Bosma, Henryk Michalewski, David Dohan, Ellen
534 Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. 2021.
535 **Program synthesis with large language models**. *arXiv*
536 *preprint arXiv:2108.07732*.

537 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang,
538 Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei
539 Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin,
540 Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu,
541 Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren,
542 Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong

543 Tu, Peng Wang, Shijie Wang, Wei Wang, Sheng-
544 guang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang,
545 Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu,
546 Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingx-
547 uan Zhang, Yichang Zhang, Zhenru Zhang, Chang
548 Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang
549 Zhu. 2023. **Qwen technical report**. *arXiv preprint*
550 *arXiv:2309.16609*, abs/2309.16609.

551 Mohammad Bavarian, Heewoo Jun, Nikolas Tezak,
552 John Schulman, Christine McLeavey, Jerry Tworek,
553 and Mark Chen. 2022. Efficient training of lan-
554 guage models to fill in the middle. *arXiv preprint*
555 *arXiv:2207.14255*.

556 Tom Brown, Benjamin Mann, Nick Ryder, Melanie
557 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
558 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
559 Askell, et al. 2020. Language models are few-shot
560 learners. *Advances in neural information processing*
561 *systems*, 33:1877–1901.

562 Federico Cassano, John Gouwar, Daniel Nguyen, Syd-
563 ney Nguyen, Luna Phipps-Costin, Donald Pinckney,
564 Ming-Ho Yee, Yangtian Zi, Carolyn Jane Anderson,
565 Molly Q Feldman, et al. 2023. Multipl-e: a scal-
566 able and polyglot approach to benchmarking neural
567 code generation. *IEEE Transactions on Software*
568 *Engineering*.

569 Linzheng Chai, Shukai Liu, Jian Yang, Yuwei Yin,
570 Ke Jin, Jiaheng Liu, Tao Sun, Ge Zhang, Changyu
571 Ren, Hongcheng Guo, et al. 2024. Mceval: Mas-
572 sively multilingual code evaluation. *arXiv preprint*
573 *arXiv:2406.07436*.

574 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan,
575 Henrique Pondé de Oliveira Pinto, Jared Kaplan,
576 Harrison Edwards, Yuri Burda, Nicholas Joseph,
577 Greg Brockman, Alex Ray, Raul Puri, Gretchen
578 Krueger, Michael Petrov, Heidy Khlaaf, Girish Sas-
579 try, Pamela Mishkin, Brooke Chan, Scott Gray,
580 Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz
581 Kaiser, Mohammad Bavarian, Clemens Winter,
582 Philippe Tillet, Felipe Petroski Such, Dave Cum-
583 mings, Matthias Plappert, Fotios Chantzis, Eliza-
584 beth Barnes, Ariel Herbert-Voss, William Heben
585 Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie
586 Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain,
587 William Saunders, Christopher Hesse, Andrew N.
588 Carr, Jan Leike, Joshua Achiam, Vedant Misra, Evan
589 Morikawa, Alec Radford, Matthew Knight, Miles
590 Brundage, Mira Murati, Katie Mayer, Peter Welinder,
591 Bob McGrew, Dario Amodei, Sam McCandlish, Ilya
592 Sutskever, and Wojciech Zaremba. 2021. **Evaluat-**
593 **ing large language models trained on code**. *arXiv*
594 *preprint arXiv:2107.03374*, abs/2107.03374.

595 Hantian Ding, Varun Kumar, Yuchen Tian, Zijian Wang,
596 Rob Kwiatkowski, Xiaopeng Li, Murali Krishna Ra-
597 manathan, Baishakhi Ray, Parminder Bhatia, Sudipta
598 Sengupta, et al. 2023a. A static evaluation of code
599 completion by large language models. *arXiv preprint*
600 *arXiv:2306.03203*.

601	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. 2023b. Crosscodeeval: A diverse and multilingual benchmark for cross-file code completion. In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023.</i>	657
602		658
603		659
604		660
605		
606		
607		
608		
609		
610	Yangruibo Ding, Zijian Wang, Wasi Uddin Ahmad, Murali Krishna Ramanathan, Ramesh Nallapati, Parminder Bhatia, Dan Roth, and Bing Xiang. 2022. Cocomic: Code completion by jointly modeling in-file and cross-file context. <i>arXiv preprint arXiv:2212.10007.</i>	661
611		662
612		663
613		664
614		665
615		666
616	Yifeng Ding, Hantian Ding, Shiqi Wang, Qing Sun, Varun Kumar, and Zijian Wang. 2024. Horizon-length prediction: Advancing fill-in-the-middle capabilities for code generation with lookahead planning. <i>arXiv preprint arXiv:2410.03103.</i>	667
617		668
618		669
619		670
620		671
621	Zhangyin Feng, Daya Guo, Duyu Tang, Nan Duan, Xiaocheng Feng, Ming Gong, Linjun Shou, Bing Qin, Ting Liu, Dixin Jiang, and Ming Zhou. 2020. Codebert: A pre-trained model for programming and natural languages. In <i>Findings of the Association for Computational Linguistics: EMNLP 2020, Online Event, 16-20 November 2020</i> , volume EMNLP 2020 of <i>Findings of ACL</i> , pages 1536–1547. Association for Computational Linguistics.	672
622		673
623		674
624		675
625		676
626		677
627		
628		
629		
630	Daniel Fried, Armen Aghajanyan, Jessy Lin, Sida I. Wang, Eric Wallace, Freda Shi, Ruiqi Zhong, Wentao Yih, Luke Zettlemoyer, and Mike Lewis. 2022. Incoder: A generative model for code infilling and synthesis. <i>arXiv preprint arXiv:2204.05999, abs/2204.05999.</i>	678
631		679
632		680
633		681
634		682
635		
636	Linyuan Gong, Sida Wang, Mostafa Elhoushi, and Alvin Cheung. 2024. Evaluation of llms on syntax-aware code fill-in-the-middle tasks. In <i>Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024</i> . OpenReview.net.	683
637		684
638		685
639		686
640		687
641	Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, Y Wu, YK Li, et al. 2024a. Deepseek-coder: When the large language model meets programming—the rise of code intelligence. <i>arXiv preprint arXiv:2401.14196.</i>	688
642		689
643		690
644		691
645		
646		
647	Jiawei Guo, Ziming Li, Xueling Liu, Kaijing Ma, Tianyu Zheng, Zhouliang Yu, Ding Pan, Yizhi Li, Ruibo Liu, Yue Wang, et al. 2024b. Codeeditor-bench: Evaluating code editing capability of large language models. <i>arXiv preprint arXiv:2404.03543.</i>	692
648		693
649		694
650		695
651		696
652	Siming Huang, Tianhao Cheng, Jason Klein Liu, Jiaran Hao, Liuyihan Song, Yang Xu, J Yang, JH Liu, Chenchen Zhang, Linzheng Chai, et al. 2024. Opencoder: The open cookbook for top-tier code large language models. <i>arXiv preprint arXiv:2411.04905.</i>	711
653		712
654		713
655		714
656		715
		716

717	Yujia Li, David H. Choi, Junyoung Chung, Nate Kushman, Julian Schrittweiser, Rémi Leblond, Tom Ecclles, 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 AlphaCode. <i>arXiv preprint arXiv:2203.07814</i> , abs/2203.07814.	773 774 775 776 777 778 779 780 781 782
729	Derrick Lin, James Koppel, Angela Chen, and Armando Solar-Lezama. 2017. Quixbugs: A multi-lingual program repair benchmark set based on the quixey challenge. In <i>Proceedings Companion of the 2017 ACM SIGPLAN international conference on systems, programming, languages, and applications: software for humanity</i> , pages 55–56.	773 774 775 776 777 778 779 780 781 782
736	Jiawei Liu, Thanh Nguyen, Mingyue Shang, Hantian Ding, Xiaopeng Li, Yu Yu, Varun Kumar, and Zijian Wang. 2024a. Learning code preference via synthetic evolution. <i>arXiv preprint arXiv:2410.03837</i> .	773 774 775 776 777 778 779 780 781 782
740	Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. 2023a. Is your code generated by chatgpt really correct? rigorous evaluation of large language models for code generation. <i>arXiv preprint arXiv:2305.01210</i> , abs/2305.01210.	773 774 775 776 777 778 779 780 781 782
745	Shukai Liu, Linzheng Chai, Jian Yang, Jiajun Shi, He Zhu, Liran Wang, Ke Jin, Wei Zhang, Hualei Zhu, Shuyue Guo, et al. 2024b. Mdeval: Massively multilingual code debugging. <i>arXiv preprint arXiv:2411.02310</i> .	773 774 775 776 777 778 779 780 781 782
750	Siyao Liu, He Zhu, Jerry Liu, Shulin Xin, Aoyan Li, Rui Long, Li Chen, Jack Yang, Jinxiang Xia, ZY Peng, et al. 2024c. Fullstack bench: Evaluating llms as full stack coder. <i>arXiv preprint arXiv:2412.00535</i> .	773 774 775 776 777 778 779 780 781 782
754	Tianyang Liu, Canwen Xu, and Julian McAuley. 2023b. Repobench: Benchmarking repository-level code auto-completion systems. <i>arXiv preprint arXiv:2306.03091</i> .	773 774 775 776 777 778 779 780 781 782
758	Anton Lozhkov, Raymond Li, Loubna Ben Allal, Federico Cassano, Joel Lamy-Poirier, Nouamane Tazi, Ao Tang, Dmytro Pykhtar, Jiawei Liu, Yuxiang Wei, et al. 2024. Starcoder 2 and the stack v2: The next generation. <i>arXiv preprint arXiv:2402.19173</i> .	773 774 775 776 777 778 779 780 781 782
763	Mayank Mishra, Matt Stallone, Gaoyuan Zhang, Yikang Shen, Aditya Prasad, Adriana Meza Soria, Michele Merler, Parameswaran Selvam, Saptha Surendran, Shivdeep Singh, et al. 2024. Granite code models: A family of open foundation models for code intelligence. <i>arXiv preprint arXiv:2405.04324</i> .	773 774 775 776 777 778 779 780 781 782
769	MistralAI. 2024. Codestral. https://mistral.ai/news/codestratal . 2024.05.29.	773 774 775 776 777 778 779 780 781 782
771	Erik Nijkamp, Hiroaki Hayashi, Caiming Xiong, Silvio Savarese, and Yingbo Zhou. 2023. Codegen2: Lessons for training llms on programming and natural languages. <i>CoRR</i> , abs/2305.02309.	773 774 775 776 777 778 779 780 781 782
779	OpenAI. 2023. Gpt-4 technical report. <i>arXiv preprint arXiv:2303.08774</i> .	773 774 775 776 777 778 779 780 781 782
789	Gabriel Orlanski, Kefan Xiao, Xavier Garcia, Jeffrey Hui, Joshua Howland, Jonathan Malmaud, Jacob Austin, Rishabh Singh, and Michele Catasta. 2023. Measuring the impact of programming language distribution. In <i>International Conference on Machine Learning</i> , pages 26619–26645. PMLR.	773 774 775 776 777 778 779 780 781 782
799	Qiwei Peng, Yekun Chai, and Xuhong Li. 2024. Humaneval-xl: A multilingual code generation benchmark for cross-lingual natural language generalization. <i>arXiv preprint arXiv:2402.16694</i> .	773 774 775 776 777 778 779 780 781 782
809	Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémie Rapin, et al. 2023. Code Llama: Open foundation models for code. <i>arXiv preprint arXiv:2308.12950</i> .	773 774 775 776 777 778 779 780 781 782
819	Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. 2022. BLOOM: A 176B-parameter open-access multilingual language model. <i>arXiv preprint arXiv:2211.05100</i> .	773 774 775 776 777 778 779 780 781 782
829	Runchu Tian, Yining Ye, Yujia Qin, Xin Cong, Yankai Lin, Zhiyuan Liu, and Maosong Sun. 2024. Debugbench: Evaluating debugging capability of large language models. <i>arXiv preprint arXiv:2401.04621</i> .	773 774 775 776 777 778 779 780 781 782
839	Yue Wang, Weishi Wang, Shafiq Joty, and Steven CH Hoi. 2021. CodeT5: Identifier-aware unified pre-trained encoder-decoder models for code understanding and generation. <i>arXiv preprint arXiv:2109.00859</i> .	773 774 775 776 777 778 779 780 781 782
849	Zhiruo Wang, Shuyan Zhou, Daniel Fried, and Graham Neubig. 2023. Execution-based evaluation for open-domain code generation. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 1271–1290.	773 774 775 776 777 778 779 780 781 782
859	Yuxiang Wei, Zhe Wang, Jiawei Liu, Yifeng Ding, and Lingming Zhang. 2023. Magicoder: Source code is all you need. <i>arXiv preprint arXiv:2312.02120</i> , abs/2312.02120.	773 774 775 776 777 778 779 780 781 782
869	Di Wu, Wasi Uddin Ahmad, Dejiao Zhang, Murali Krishna Ramanathan, and Xiaofei Ma. 2024. Reporformer: Selective retrieval for repository-level code completion. In <i>Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024</i> . OpenReview.net.	773 774 775 776 777 778 779 780 781 782
879	Weixiang Yan, Yuchen Tian, Yunzhe Li, Qian Chen, and Wen Wang. 2023. Codetransocean: A comprehensive multilingual benchmark for code translation. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 5067–5089. Association for Computational Linguistics.	773 774 775 776 777 778 779 780 781 782

829 Jian Yang, Jiaxi Yang, Ke Jin, Yibo Miao, Lei Zhang,
830 Liqun Yang, Zeyu Cui, Yichang Zhang, Binyuan
831 Hui, and Junyang Lin. 2024. Evaluating and align-
832 ing codellms on human preference. *arXiv preprint*
833 *arXiv:2412.05210*.

834 Pengcheng Yin, Wen-Ding Li, Kefan Xiao, Abhishek
835 Rao, Yeming Wen, Kensen Shi, Joshua Howland,
836 Paige Bailey, Michele Catasta, Henryk Michalewski,
837 et al. 2023. Natural language to code generation in
838 interactive data science notebooks. In *Proceedings*
839 of the 61st Annual Meeting of the Association for
840 Computational Linguistics (Volume 1: Long Papers),
841 pages 126–173.

842 Hao Yu, Bo Shen, Dezhi Ran, Jiaxin Zhang, Qi Zhang,
843 Yuchi Ma, Guangtai Liang, Ying Li, Qianxiang Wang,
844 and Tao Xie. 2024. Codereval: A benchmark of prag-
845 matic code generation with generative pre-trained
846 models. In *Proceedings of the 46th IEEE/ACM Inter-
847 national Conference on Software Engineering*, pages
848 1–12.

849 Fengji Zhang, Bei Chen, Yue Zhang, Jacky Keung, Jin
850 Liu, Daoguang Zan, Yi Mao, Jian-Guang Lou, and
851 Weizhu Chen. 2023. Repocoder: Repository-level
852 code completion through iterative retrieval and gen-
853 eration. In *Proceedings of the 2023 Conference on*
854 *Empirical Methods in Natural Language Process-*
855 *ing, EMNLP 2023, Singapore, December 6-10, 2023*,
856 pages 2471–2484. Association for Computational
857 Linguistics.

858 Heri Zhao, Jeffrey Hui, Joshua Howland, Nam Nguyen,
859 Siqi Zuo, Andrea Hu, Christopher A. Choquette-
860 Choo, Jingyue Shen, Joe Kelley, Kshitij Bansal,
861 Luke Vilnis, Mateo Wirth, Paul Michel, Peter
862 Choy, Pratik Joshi, Ravin Kumar, Sarmad Hashmi,
863 Shubham Agrawal, Zhitao Gong, Jane Fine, Tris
864 Warkentin, Ale Jakse Hartman, Bin Ni, Kathy Ko-
865 revec, Kelly Schaefer, and Scott Huffman. 2024.
866 **Codegemma: Open code models based on gemma.**
867 *CoRR*, abs/2406.11409.

868 Qinkai Zheng, Xiao Xia, Xu Zou, Yuxiao Dong, Shan
869 Wang, Yufei Xue, Zihan Wang, Lei Shen, Andi Wang,
870 Yang Li, Teng Su, Zhilin Yang, and Jie Tang. 2023.
871 **Codegeex: A pre-trained model for code generation**
872 **with multilingual evaluations on humaneval-x.** *arXiv*
873 *preprint arXiv:2303.17568*, abs/2303.17568.

874 Tianyu Zheng, Ge Zhang, Tianhao Shen, Xueling Liu,
875 Bill Yuchen Lin, Jie Fu, Wenhui Chen, and Xiang
876 Yue. 2024. Opencodeinterpreter: Integrating code
877 generation with execution and refinement. *arXiv*
878 *preprint arXiv:2402.14658*.

879 Terry Yue Zhuo, Minh Chien Vu, Jenny Chim, Han Hu,
880 Wenhao Yu, Ratnadira Widayarsi, Imam Nur Bani
881 Yusuf, Haolan Zhan, Junda He, Indraneil Paul, et al.
882 2024. Bigcodebench: Benchmarking code genera-
883 tion with diverse function calls and complex instruc-
884 tions. *arXiv preprint arXiv:2406.15877*.

885 A Related Work

886 **Code Large Language Models** In software engineering, the advent of large language models
887 (LLMs) tailored for code-centric tasks has proven
888 to be transformative. Models (Feng et al., 2020;
889 Chen et al., 2021; Scao et al., 2022; Li et al.,
890 2022; Allal et al., 2023; Fried et al., 2022; Wang
891 et al., 2021; Zheng et al., 2024; Jiang et al.,
892 2024; Nijkamp et al., 2023; Wei et al., 2023;
893 Zhao et al., 2024) like CodeLlama (Rozière et al.,
894 2023), DeepSeek-Coder (Guo et al., 2024a), Open-
895 Coder (Huang et al., 2024) and Qwen2.5-Coder
896 (Hui et al., 2024) — all trained on vast corpuses
897 comprising billions of code snippets — have funda-
898 mentally augmented the development process.
899 These Code LLMs are instrumental in automating
900 repetitive software tasks, proposing code improve-
901 ments, and facilitating the conversion of natural
902 language into executable code. Notable among
903 these are Starcoder (Li et al., 2023; Lozhkov et al.,
904 2024), CodeLlama (Rozière et al., 2023), and Code-
905 Qwen (Bai et al., 2023), each bringing unique con-
906 tributions to the enhancement of coding assistance
907 tools. With these advancements, Code LLMs show-
908 case a promising trajectory for further revolutioniz-
909 ing how developers interact with code, promising
910 ever-greater efficiency and intuitiveness in software
911 creation. Inspired by the success of the grammar-
912 based parsed tree in many fields, we adopt the ab-
913 stract syntax tree to augment the code completion
914 training.

916 **Repo-level Code Evaluation** In the domain of
917 code evaluation, a rich tapestry of benchmarks
918 (Zheng et al., 2023; Yu et al., 2024; Yin et al.,
919 2023; Peng et al., 2024; Khan et al., 2023; Guo
920 et al., 2024b; Lai et al., 2023) has been woven
921 to address the challenges of accurately assessing
922 code quality, functionality, and efficiency, such
923 as HumanEval/MBPP (Chen et al., 2021; Austin
924 et al., 2021), their upgraded version EvalPlus (Liu
925 et al., 2023a), and the multilingual benchmark
926 MultiPL-E (Cassano et al., 2023), McEval (Chai
927 et al., 2024), and MdEval (Liu et al., 2024b). Big-
928 CodeBench (Zhuo et al., 2024), fullstack (Liu
929 et al., 2024c), CodeFavor (Liu et al., 2024a),
930 CodeArena (Yang et al., 2024) and SAFIM (Gong
931 et al., 2024) separately evaluate code LLMs for
932 more diverse scenarios and code preferences. Stud-
933 ies have explored a variety of approaches, ranging
934 from static analysis techniques (e.g. exact match
935 (EM) and edit similarity (ES)), which examine code

936 without executing it, to dynamic methods that in-
937 volve code execution in controlled environments
938 (e.g. Pass@k). The current benchmarks support
939 code models to evaluate a series of different types
940 of tasks, such as code understanding, code repair
941 (Lin et al., 2017; Tian et al., 2024; Jimenez et al.,
942 2023), code translation (Yan et al., 2023), and multi-
943 lingual scenarios (Wang et al., 2023; Athiwaratkun
944 et al., 2023; Zheng et al., 2023; Peng et al., 2024;
945 Orlanski et al., 2023). An important task FIM
946 (Fried et al., 2022; Bavarian et al., 2022; Ding
947 et al., 2024) is to fill the middle code, given the
948 prefix and suffix code, which provides substantial
949 assistance for software development. Repo-level
950 completion, such as RepoEval (Zhang et al., 2023),
951 CrossCodeEval (Ding et al., 2023b; Wu et al., 2024)
952 and RepoBench (Liu et al., 2023b) only using exact
953 match and edit similarity without code execution
954 can not accurately reflect the model performance
955 and Humaneval-Fim (Zheng et al., 2023) focus in-
956 file evaluation.