

# 000 001 002 003 004 005 006 007 008 009 010 011 012 013 014 015 016 017 018 019 020 021 022 023 024 025 026 027 028 029 030 031 032 033 034 035 036 037 038 039 040 041 042 043 044 045 046 047 048 049 050 051 052 053 MCEVAL: MASSIVELY MULTILINGUAL CODE EVALUA- TION

Anonymous authors

Paper under double-blind review

## ABSTRACT

Code large language models (LLMs) have shown remarkable advances in code understanding, completion, and generation tasks. Programming benchmarks, comprised of a selection of code challenges and corresponding test cases, serve as a standard to evaluate the capability of different LLMs in such tasks. However, most existing benchmarks primarily focus on Python and are still restricted to a limited number of languages, where other languages are translated from the Python samples (e.g. MultiPL-E) degrading the data diversity. To further facilitate the research of code LLMs, we propose a massively multilingual code benchmark covering 40 programming languages (MCEVAL) with 16K test samples, which substantially pushes the limits of code LLMs in multilingual scenarios. The benchmark contains challenging code completion, understanding, and generation evaluation tasks with finely curated massively multilingual instruction corpora MCEVAL-INSTRUCT. In addition, we introduce an effective multilingual coder MCODER trained on MCEVAL-INSTRUCT to support multilingual programming language generation. Extensive experimental results on MCEVAL show that there is still a difficult journey between open-source models and closed-source LLMs (e.g. GPT-series models) in numerous languages.

## 1 INTRODUCTION

Large language models (LLMs) designed for code, such as Codex (Chen et al., 2021), CodeGen (Nijkamp et al., 2023), Code Llama (Rozière et al., 2023), DeepSeekCoder (Guo et al., 2024), and CodeQwen (Hui et al., 2024) excel at code understanding, completion, and generation tasks.

Code LLMs with a large number of parameters (e.g. 7B, 13B, or larger) are pre-trained on large-scale code databases with self-supervised autoregressive objectives, followed by instruction tuning (Ouyang et al., 2022) for aligning to human preferences and downstream code-related tasks. Most code benchmarks (Chen et al., 2021; Austin et al., 2021; Athiwaratkun et al., 2023) are introduced to evaluate the performance of code LLMs by assessing their ability to generate executable code based on the problem descriptions. The assessments aim to gauge the capacity of the models to understand and generate code effectively, thereby contributing to facilitating and streamlining the programming process for developers. The execution-based method executes generated code against test cases to measure the success rate. Due to the difficulty of creating the problem and its corresponding solution (requiring specialized programming staff), the development of evaluation benchmarks is limited within Python, with a few other languages being translated from Python. *Therefore, the*

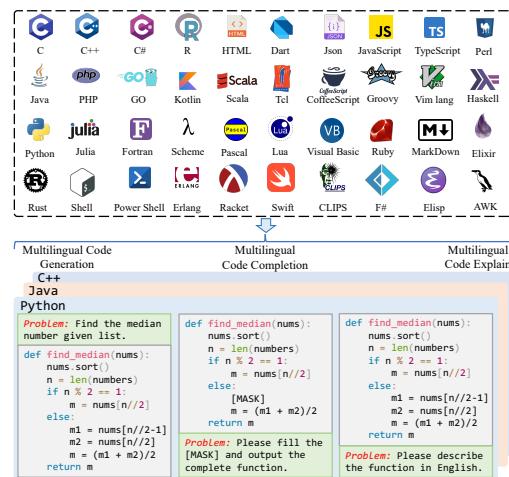


Figure 1: Massively multilingual evaluation task comprised of three tasks, including code generation, code completion, and code explanation.

054 *community desperately needs a massively multilingual programming benchmark (not from HumanEval or MBPP) comprised of instruction corpora and evaluation set to comprehensively facilitate*  
 055 *and evaluate the generation, completion, and understanding capability of LLMs.*

056  
 057 To facilitate the development of code LLMs, we introduce a complete framework that includes  
 058 the multilingual code instruction corpora, multilingual coder (MCODER), and multilingual code  
 059 evaluation benchmark. First, we propose MCEVAL, the first massively multilingual code evaluation  
 060 benchmark (from human handwriting) covering 40 languages (16K samples in total), encompassing  
 061 multilingual code generation, multilingual code explanation, and multilingual code completion  
 062 tasks. Then, we create a massively multilingual instruction corpora MCEVAL-INSTRUCT of 40  
 063 languages. We initially select and refine high-quality code snippets from various programming  
 064 languages (PLs) using an LLM. The LLM then generates clear and self-contained instructional  
 065 content, including problem descriptions and corresponding solutions, based on the refined snippets.  
 066 To ensure consistency and enhance learning across languages, we introduce cross-lingual code  
 067 transfer, adapting instructional content to different PLs while increasing sample complexity. Based  
 068 on open-source models and MCEVAL-INSTRUCT, MCODER is used as a strong baseline to explore  
 069 the transferability of LLMs among different PLs.

070 The contributions are summarized as follows: (1) We propose MCEVAL with enough test samples  
 071 (16K), a true massively multilingual multitask code evaluation benchmark (not from HumanEval  
 072 or MBPP) covering 40 languages, encompassing multilingual code generation, multilingual code  
 073 explanation, and multilingual code completion tasks. (2) We introduce MCEVAL-INSTRUCT, the  
 074 massively multilingual code instruction corpora covering from the multilingual code snippet from  
 075 40 languages. Based on MCEVAL-INSTRUCT, an effective multilingual coder MCODER is used as  
 076 a strong baseline for MCEVAL. (3) We systematically evaluate the understanding and generation  
 077 capabilities of 20+ models on our created MCEVAL and create a leaderboard to evaluate them on  
 078 40 programming languages dynamically. Notably, extensive experiments suggest that comprehensive  
 079 multilingual multitask evaluation can realistically measure the gap between open-source (e.g.  
 080 DeepSeekCoder and CodeQwen1.5) and closed-source models (e.g. GPT-3.5 and GPT-4).

## 082 2 MULTILINGUAL CODE EVALUATION: MCEVAL

### 084 2.1 DATASET STATISTICS

086 The created MCEVAL is comprised of three key code-related  
 087 tasks covering 40 programming languages, including multi-  
 088 lingual code generation, multilingual code explanation, and  
 089 multilingual code completion tasks. The multilingual code  
 090 generation and explanation tasks separately contain 2K sam-  
 091 ples, where each language has nearly 50 samples. The code  
 092 completion task can be decomposed into *multi-line comple-*  
 093 *tion* (3K samples), *single-line completion* (3K samples), *span*  
 094 *completion* (4K samples), and *span completion (light)* (2K  
 095 samples) (Bavarian et al., 2022).

096 In Table 1, we display the number of questions, test cases, and  
 097 difficulty levels corresponding to the three tasks in MCEVAL  
 098 and the number of questions in the four sub-tasks of the  
 099 completion task. Moreover, we counted the token length  
 100 of the prompt and solutions. (The tokens are calculated  
 101 based on the Llama-3 tokenizer.) Among these tasks, the  
 102 *span completion (light)* task is similar in form to the *span*  
 103 *completion* task. However, in the *span completion (light)*  
 104 task, each problem is paired with all the corresponding code,  
 105 making it a balanced version of the *span completion* task  
 106 (fewer samples for fast inference and the same test size of  
 107 each programming language). The results of *span completion*  
 108 (*light*) can better reflect the differences in model performance  
 109 across different languages.

Table 1: MCEVAL dataset statistics.

| Statistics              | Value        |
|-------------------------|--------------|
| <b>Questions</b>        |              |
| Code Generation         | 2,007        |
| Code Explanation        | 2,007        |
| Code Completion         | 12,017       |
| - Single-Line           | 2,998        |
| - Multi-Line            | 2,998        |
| - Span                  | 4,014        |
| - Span(light)           | 2,007        |
| Total Test Cases        | 10,086       |
| <b>Difficulty Level</b> |              |
| - Easy                  | 1,221        |
| - Medium                | 401          |
| - Hard                  | 385          |
| <b>Length</b>           |              |
| Prompt                  |              |
| - maximum length        | 793 tokens   |
| - minimum length        | 16 tokens    |
| - avg length            | 173.8 tokens |
| Solution(Output)        |              |
| - maximum length        | 666 tokens   |
| - minimum length        | 4 tokens     |
| - avg length            | 120.9 tokens |

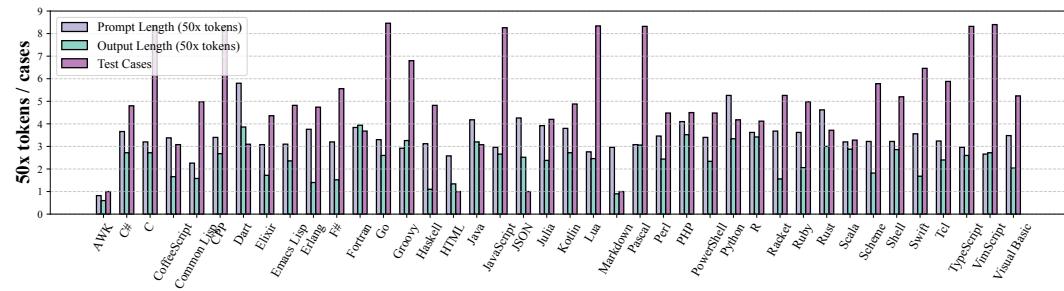


Figure 2: Data statistics of the MCEVAL benchmark involving 40 programming languages.

Figure 2 plots the length of the prompt, solution(output), and the number of test cases of each programming language.

In Table 2, We compared MCEVAL with other multilingual benchmarks. It is noteworthy that our benchmark provides a significant supplement to current benchmarks in terms of both the variety of programming languages and the number of questions.

Table 2: Comparison between MCEVAL and other multilingual code benchmarks.  $\diamond$ The number of each of the three tasks (Generation, Explanation, and Completion).

| Benchmark                         | Multi-Task | #Languages | Data source  | #Questions                          |
|-----------------------------------|------------|------------|--------------|-------------------------------------|
| MuliPL-E (Cassano et al., 2023)   | ✗          | 18         | Translate    | ~3,000                              |
| MBXP (Athiwaratkun et al., 2023)  | ✓          | 10         | Translate    | 12,425                              |
| HumanEval-X (Zheng et al., 2023b) | ✓          | 5          | Hand-Written | 820                                 |
| HumanEval-XL (Peng et al., 2024)  | ✗          | 12         | Hand-Written | 22,080                              |
| MCEVAL                            | ✓          | 40         | Hand-Written | 16,031 (2007/2007/12017) $\diamond$ |

## 2.2 HUMAN ANNOTATION & QUALITY CONTROL

To create the massively multilingual code evaluation benchmark MCEVAL, the annotation of multilingual code samples is conducted utilizing a comprehensive and systematic human annotation procedure, underpinned by rigorously defined guidelines to ensure accuracy and consistency. Initially, 10 software developers in computer science are recruited as multilingual programming annotators with proven proficiency in the respective programming languages. Following a detailed training session on the annotation protocol, which emphasizes the importance of context, syntactical correctness, and semantic fidelity across languages, annotators are tasked with creating problem definitions and the corresponding solution. The annotators should follow: (1) Provide a clear and self-contained problem definition, answer the question with any tools, and design the test cases to evaluate the correctness of the code. (2) Classify them into multiple difficulties (Easy/Middle/Hard), based on algorithmic complexity and functionality. Each sample is independently annotated by at least two annotators to minimize subjective bias and errors. Discrepancies between annotators are resolved through consensus or adjudication by a senior annotator. Finally, three volunteers are employed to evaluate the correctness of the benchmark (> 90% accuracy) and correct the errors. (See Appendix A.2 for more details).

## 2.3 EVALUATION TASKS

**Multilingual Code Generation.** Given the  $k$ -th programming language  $L_k \in \{L_i\}_{i=1}^K$ , where  $K = 40$  is the number of programming languages, we provide the problem description  $q^{L_k}$  and examples test cases  $e^{L_k}$  as the input for code LLMs  $\mathcal{M}$  to generate the corresponding code  $a^{L_k}$ . We obtain the sampled code result from the code generation distribution  $P(a^{L_k}|q^{L_k}, e^{L_k}; \mathcal{M})$  from code LLM  $\mathcal{M}$ , and then feed the test cases into the generated code, where the generated outputs by code should equal the expected outputs. The process can be described as:

$$r^{L_k} = \mathbb{I}(P(a^{L_k}|q^{L_k}, e^{L_k}; \mathcal{M}); u^{L_k}) \quad (1)$$

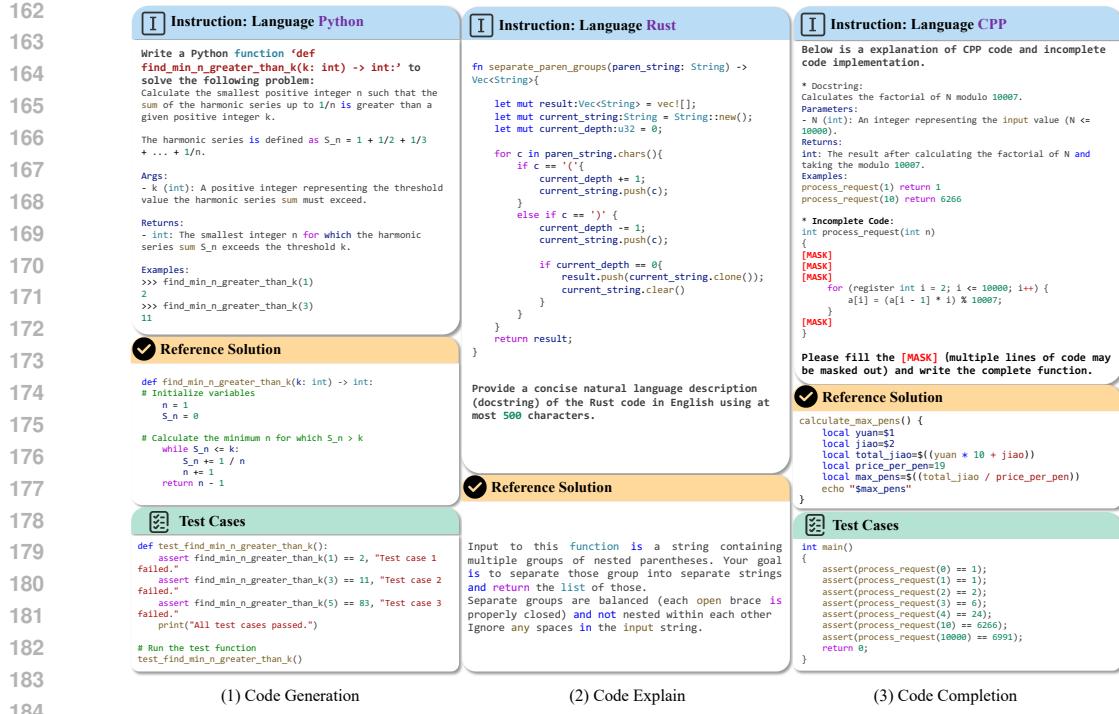


Figure 3: Examples of multilingual code generation, explanation, and completion.

where  $\mathbb{I}(\cdot)$  is the indicator function by executing the generated code with the given test cases  $u^{L_k}$ . when the generated code  $a^{L_k}$  passes all test cases, the evaluation result  $r = 1$ , else  $r = 0$ .

**Multilingual Code Explanation.** To evaluate the understanding capability of code LLMs, we adopt two-pass generation (Code-to-Natural-Language and Natural-Language-to-Code), since the text-similar metrics (e.g. BLEU (Papineni et al., 2002)) are hindered by the  $n$ -gram text matching and can not produce an accurate score. We first prompt the code LLMs to generate the natural language description  $t^{L_k}$  based on the code  $a^{L_k}$  and then we force the model to restore the original code based on  $t^{L_k}$ . The sampled code from  $P(a^{L_k}|t^{L_k}; \mathcal{M})$  is used to evaluate the understanding capability as:

$$r = \mathbb{I}(P(t^{L_k}|a^{L_k}; \mathcal{M})P(a^{L_k}|t^{L_k}; \mathcal{M}); u^{L_k}) \quad (2)$$

where  $\mathbb{I}(\cdot)$  is used to check the correctness of the generated code by running the code with test cases.

**Multilingual Code Completion.** Another important scenario is code completion, where the code LLM produces the middle code  $a_m^{L_k}$  based on the prefix code  $a_p^{L_k}$  and suffix code snippet  $a_s^{L_k}$ . Hence, we concatenate  $a_p^{L_k}$ ,  $a_m^{L_k}$ , and  $a_s^{L_k}$  as the complete code for evaluation as:

$$r = \mathbb{I}(P(a_m^{L_k}|a_p^{L_k}, a_s^{L_k}; \mathcal{M}); u^{L_k}) \quad (3)$$

where  $a_p^{L_k}$ ,  $a_q^{L_k}$ , and  $a_m^{L_k}$  are concatenated as the complete code to be executed with test cases  $u^{L_k}$ .

### 3 MCODER

#### 3.1 MCEVAL-INSTRUCT

**Collection from Code Snippet.** For a programming language  $L_k$  ( $L_k \in \{L_i\}_{i=1}^K$ ) and  $K$  is the number of programming languages), consider an existing code snippet  $c \in D_c^{L_k}$ , we prompt the LLM to select the high-quality code and refine the code to a self-contained code snippet by using the prompt “{Code Snippet}\nDetermine its educational value for a student whose goal is to learn basic coding concepts.\n\nIf the answer is ‘YES’. Please refine the code with clear variable definitions, comments, and docstring.”. Then, we can obtain the multilingual refined code snippets. (More details can be found in Appendix A.3)

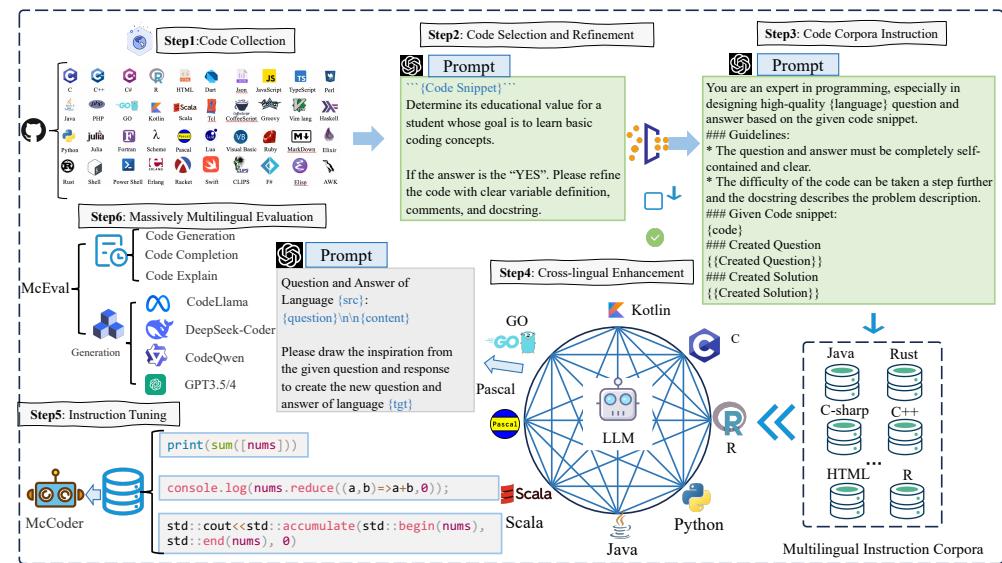


Figure 4: The framework of MCODER. We first create MCEVAL-INSTRUCT covering 40 languages from code snippets to fine-tune MCODER. 20+ existing LLMs and MCODER are then evaluated on MCEVAL comprised of multilingual code generation, explanation, and completion.

**Instruction Corpora Generation.** To construct a comprehensive massively multilingual code instruction corpora  $\{D^{L_i}\}_{i=1}^K$ , we prompt the LLMs (`gpt-4-1106-preview`) to create a problem description  $q^{L_k}$  and the corresponding solution  $a^{L_k}$  by drawing inspiration from the refined code snippet  $c^{L_k}$ . We use LLM to generate instruction dataset by using the prompt “You are an expert in programming, especially in designing high-quality language question and answer based on the given code snippet.\n\n### Guidelines: \* The question and answer must be completely self-contained and clear.\*\n\nThe difficulty of the code can be taken a step further and the docstring describes the problem description.\n\n### Given Code snippet: code\n### Created Question Created Question\n### Created Solution\n### Created Solution” in Figure 4.

**Cross-lingual Code Transfer.** Since the created instruction samples of different programming languages focus on different aspects of coding, we adopt the cross-lingual code transfer to minimize the gap among multiple languages. Given the instruction dataset  $D^{L_i}$  of language  $L_i$ , we randomly sample a pair  $(q^{L_i}, a^{L_i})$  and force the LLM to modify them to another language  $L_j$  with a more complex sample  $(q^{L_i \rightarrow L_j}, a^{L_i \rightarrow L_j})$ . In this way, we can get the derived instruction corpora  $\{D^{L_i \rightarrow L_j}\}_{i \neq j \wedge 1 \leq i, j \leq K}$ . Finally, we combine  $\{D^{L_k}\}_{k=1}^K$  and  $\{D^{L_i \rightarrow L_j}\}_{i \neq j \wedge 1 \leq i, j \leq K}$  as the multilingual instruction corpora MCEVAL-INSTRUCT  $\{D_{mc}^{L_k}\}_{k=1}^K$  covering 40 programming languages.

### 3.2 MULTILINGUAL CODE INSTRUCTION TUNING

The training objective  $\mathcal{L}_{all}$  of the multilingual instruction fine-tuning can be described as:

$$\mathcal{L}_{all} = - \sum_{k=1}^K \mathbb{E}_{q^{L_k}, a^{L_k} \sim \{D^{L_k}\}_{k=1}^K} [\log P(a^{L_k} | q^{L_k}; \mathcal{M})] \quad (4)$$

where  $q^{L_k}$  and  $a^{L_k}$  are the code-related question and answer from the dataset  $D^{L_k}$  of language  $L_k$ , respectively.  $K$  is the number of programming languages.

## 4 EXPERIMENTS

### 4.1 EXPERIMENT SETUP

**Code LLMs.** We evaluate 30 + models with sizes ranging from 7B to 236B parameters, including general/code LLMs, open/closed-source models, and base/instruction models. For general models, we evaluate GPT series (Brown et al., 2020; OpenAI, 2023), Qwen1.5 (Bai et al., 2023), Llama3 (AI,

Table 3: Pass@1 (%) scores of different code LLMs for multilingual code generation tasks on MCEVAL. “Avg<sub>all</sub>” represents the average scores of all code languages.

| Method                         | Size | AWK         | C           | C++         | C#          | Clip        | Coffee      | Dart        | Erlang      | Elixir      | Fortran     | F#          | Go          | Groovy      | Haskell     | Html        | Java        | JS          | Json        | Julia       |             |      |
|--------------------------------|------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|------|
| GPT-4o (240513)                | 54.0 | <b>60.0</b> | 58.0        | <b>72.0</b> | 60.0        | <b>82.0</b> | <b>54.9</b> | <b>64.0</b> | <b>66.0</b> | 44.0        | <b>66.0</b> | <b>78.0</b> | <b>62.0</b> | <b>80.0</b> | <b>90.0</b> | 32.0        | <b>81.1</b> | 62.0        | <b>74.0</b> | <b>72.0</b> |             |      |
| GPT-4-Turbo (231106)           | 70.0 | <b>60.0</b> | 60.0        | 44.0        | 68.0        | 54.0        | 76.0        | 41.2        | 26.0        | 30.0        | 46.0        | 40.0        | 68.0        | 54.0        | <b>86.0</b> | 50.0        | 18.0        | 71.7        | 60.0        | 76.0        | 54.0        |      |
| Yi-Large-Turbo                 | 14.0 | 58.0        | 58.0        | 44.0        | 68.0        | 54.0        | 76.0        | 41.2        | 26.0        | 30.0        | 46.0        | 40.0        | 68.0        | 54.0        | <b>86.0</b> | 50.0        | 18.0        | 71.7        | 60.0        | 76.0        | 54.0        |      |
| DeepSeekCoder-V2-Instruct      | 52.0 | 44.0        | 50.0        | 54.0        | 54.0        | 52.0        | 31.4        | 30.0        | 46.0        | 30.0        | 58.0        | 52.0        | 46.0        | 80.0        | 58.0        | 22.0        | 41.5        | 52.0        | <b>78.0</b> | 56.0        |             |      |
| DeepSeekCoder-V2-1.5-Instruct  | 236B | 62.0        | <b>60.0</b> | <b>66.0</b> | 72.0        | <b>74.0</b> | 80.0        | 43.1        | 52.0        | 62.0        | 32.0        | <b>80.0</b> | 76.0        | 56.0        | 84.0        | 72.0        | 30.0        | 77.4        | <b>64.0</b> | 74.0        | 70.0        |      |
| Qwen1.5-Chat                   | 72B  | 50.0        | 48.0        | 46.0        | 56.0        | 46.0        | 44.0        | 19.6        | 14.0        | 18.0        | 24.0        | 34.0        | 32.0        | 9.0         | 42.0        | 32.0        | 22.0        | 39.6        | 46.0        | 74.0        | 32.0        |      |
| CodeLlama-Instruct             | 34B  | 38.0        | 32.0        | 40.0        | 42.0        | 34.0        | 7.8         | 16.0        | 28.0        | 28.0        | 32.0        | 24.0        | 18.0        | 34.0        | 20.0        | 32.0        | 14.0        | 26.4        | 42.0        | 68.0        | 28.0        |      |
| WizardCoder-Python             | 34B  | 36.0        | 42.0        | 46.0        | 52.0        | 42.0        | 46.0        | 13.7        | 14.0        | 38.0        | 38.0        | 26.0        | 26.0        | 50.0        | 54.0        | 40.0        | 26.0        | 43.4        | 52.0        | 70.0        | 52.0        |      |
| DeepSeekCoder-Instruct         | 33B  | 50.0        | 55.0        | 50.0        | 55.0        | 50.0        | 55.0        | 25.5        | 50.0        | 40.0        | 46.0        | 40.0        | 50.0        | 50.0        | 50.0        | 50.0        | 30.0        | 73.4        | 62.0        | 42.0        | 56.0        |      |
| Codestar-0.1                   | 22B  | 54.0        | 48.0        | 56.0        | 66.0        | 60.0        | 62.0        | 43.1        | 28.0        | 34.0        | 24.0        | 56.0        | 58.0        | 38.0        | 76.0        | 52.0        | <b>34.0</b> | 73.6        | 80.0        | 72.0        | 52.0        |      |
| DeepSeekCoder-V2-Lite-Instruct | 16B  | 58.0        | 58.0        | 58.0        | 64.0        | 71.0        | 56.0        | 26.0        | 35.3        | 30.0        | 40.0        | 26.0        | 58.0        | 56.0        | 66.0        | 60.0        | 26.0        | 64.0        | 64.0        | 68.0        | 60.0        |      |
| OCTOCODER                      | 16B  | 28.0        | 28.0        | 28.0        | 38.0        | 40.0        | 18.0        | 5.9         | 14.0        | 28.0        | 16.0        | 22.0        | 4.0         | 34.0        | 30.0        | 20.0        | 8.0         | 34.0        | 30.0        | 58.0        | 20.0        |      |
| WizardCoder-V1.0               | 15B  | 18.0        | 38.0        | 28.0        | 36.0        | 36.0        | 50.0        | 17.6        | 0.0         | 24.0        | 24.0        | 28.0        | 30.0        | 38.0        | 62.0        | 40.0        | 6.0         | 30.2        | 52.0        | 14.0        | 38.0        |      |
| Granite-34B-code-instruct-8K   | 34B  | <b>48.0</b> | <b>38.0</b> | <b>54.0</b> | <b>50.0</b> | <b>60.0</b> | <b>64.0</b> | <b>19.6</b> | <b>16.0</b> | <b>36.0</b> | <b>36.0</b> | <b>50.0</b> | <b>40.0</b> | <b>52.0</b> | <b>64.0</b> | <b>44.0</b> | <b>28.0</b> | <b>37.7</b> | <b>54.0</b> | <b>66.0</b> | <b>48.0</b> |      |
| Granite-20B-code-instruct-8K   | 20B  | <b>38.0</b> | <b>42.0</b> | <b>40.0</b> | <b>56.0</b> | <b>40.0</b> | <b>56.0</b> | <b>23.5</b> | <b>18.0</b> | <b>40.0</b> | <b>16.0</b> | <b>32.0</b> | <b>42.0</b> | <b>46.0</b> | <b>48.0</b> | <b>40.0</b> | <b>20B</b>  | <b>35.8</b> | <b>50.0</b> | <b>74.0</b> | <b>48.0</b> |      |
| Granite-8B-code-instruct-4K    | 8B   | <b>56.0</b> | <b>40.0</b> | <b>50.0</b> | <b>50.0</b> | <b>50.0</b> | <b>56.0</b> | <b>11.0</b> | <b>14.0</b> | <b>36.0</b> | <b>36.0</b> | <b>56.0</b> | <b>36.0</b> | <b>40.0</b> | <b>50.0</b> | <b>40.0</b> | <b>20B</b>  | <b>35.8</b> | <b>50.0</b> | <b>64.0</b> | <b>48.0</b> |      |
| Granite-3B-code-instruct-128K  | 3B   | 16.0        | 14.0        | 4.0         | 38.0        | 30.0        | 30.0        | 9.8         | 12.0        | 36.0        | 16.0        | 26.0        | 26.0        | 30.0        | 32.0        | 34.0        | 22.0        | 26.4        | 50.0        | 62.0        | 28.0        |      |
| Phi-3-medium-4k-instruct       | 14B  | 52.0        | 46.0        | 40.0        | 52.0        | 34.0        | 42.0        | 13.7        | 14.0        | 16.0        | 8.0         | 30.0        | 32.0        | 42.0        | 42.0        | 42.0        | 20.0        | 39.6        | 52.0        | 68.0        | 48.0        |      |
| CodeLlama-Instruct             | 13B  | 36.0        | 38.0        | 38.0        | 40.0        | 46.0        | 30.0        | 7.8         | 16.0        | 32.0        | 32.0        | 16.0        | 26.0        | 34.0        | 22.0        | 38.0        | 18.0        | 22.6        | 34.0        | 56.0        | 32.0        |      |
| Llama-3-Model                  | 8B   | 32.0        | 46.0        | 50.0        | 54.0        | 38.0        | 48.0        | 15.7        | 14.0        | 32.0        | 30.0        | 12.0        | 26.0        | 48.0        | 52.0        | 38.0        | 16.0        | 45.3        | 54.0        | 70.0        | 40.0        |      |
| CodeQwen-2.5-Chat              | 7B   | 58.0        | 58.0        | 58.0        | 74.0        | 56.0        | <b>86.0</b> | 37.3        | 54.0        | 56.0        | 38.0        | 56.0        | 62.0        | 58.0        | 84.0        | 66.0        | 34.0        | 69.8        | 60.0        | 76.0        | 64.0        |      |
| Codegemm-it                    | 7B   | 26.0        | 40.0        | 42.0        | 48.0        | 20.0        | 23.5        | 10.0        | 4.0         | 4.0         | 8.0         | 22.0        | 46.0        | 58.0        | 32.0        | 24.0        | 56.6        | 48.0        | 70.0        | 38.0        |             |      |
| CodeLlama-Instruct             | 7B   | 32.0        | 34.0        | 24.0        | 40.0        | 42.0        | 32.0        | 5.9         | 14.0        | 18.0        | 22.0        | 20.0        | 40.0        | 14.0        | 30.0        | 32.0        | 14.0        | 18.8        | 28.0        | 64.0        | 24.0        |      |
| CodeShell-chat                 | 7B   | 24.0        | 30.0        | 26.0        | 36.0        | 28.0        | 34.0        | 5.9         | 4.0         | 14.0        | 6.0         | 8.0         | 28.0        | 30.0        | 22.0        | 24.0        | 22.6        | 44.0        | 66.0        | 24.0        |             |      |
| DeepSeekCoder-1.5-Instruct     | 7B   | 40.0        | 54.0        | 56.0        | 60.0        | 56.0        | 56.0        | 23.5        | 24.0        | 40.0        | 40.0        | 46.0        | 46.0        | 52.0        | 52.0        | 40.0        | 24.0        | 60.4        | 66.0        | 42.0        |             |      |
| Magicode-S-DS                  | 7B   | 44.0        | 50.0        | 50.0        | 60.0        | 58.0        | 72.0        | 19.6        | 32.0        | 34.0        | <b>62.0</b> | 62.0        | 54.0        | 50.0        | 80.0        | 54.0        | 16.0        | 66.0        | 60.0        | 56.0        | 40.0        |      |
| Ncode-CO-Orpo                  | 7B   | 38.0        | 52.0        | 58.0        | 50.0        | 46.0        | 62.0        | 23.5        | 22.0        | 38.0        | 36.0        | 52.0        | 46.0        | 52.0        | 72.0        | 46.0        | 28.0        | 56.6        | 50.0        | 64.0        | 62.0        |      |
| OpenCodeInterpreter-DS         | 7B   | 38.0        | 54.0        | 52.0        | 68.0        | 44.0        | 78.0        | 17.6        | 30.0        | 42.0        | 48.0        | 52.0        | 54.0        | 48.0        | 72.0        | 40.0        | 28.0        | 66.0        | 46.0        | 72.0        | 34.0        |      |
| CodeQwen-1.5-Chat              | 7B   | 40.0        | 52.0        | 56.0        | 62.0        | 48.0        | 62.0        | 29.4        | 22.0        | 38.0        | 38.0        | 50.0        | 50.0        | 70.0        | 44.0        | 30.0        | 58.5        | 54.0        | 64.0        | 62.0        |             |      |
| <b>MCODER (Our Method)</b>     | 7B   | 40.0        | 44.0        | 52.0        | 62.0        | 46.0        | <b>66.0</b> | 21.6        | 30.0        | 44.0        | 52.0        | 56.0        | 44.0        | 48.0        | <b>70.0</b> | 32.0        | 34.0        | <b>54.7</b> | <b>54.0</b> | <b>66.0</b> | <b>56.0</b> |      |
| <b>MCODER (Our Method)</b>     | 1    | 48.0        | 52.0        | 30.0        | 42.0        | 36.0        | 32.0        | 54.0        | 44.0        | 40.0        | 36.0        | 48.0        | 52.8        | 58.0        | 44.0        | 46.0        | 64.0        | 38.0        | 52.0        | 58.0        | 20.0        | 46.7 |

2024), Phi-3 (Abdin et al., 2024), and Yi (Young et al., 2024). For code models, we test CodeQwen (Hui et al., 2024), DeepSeekCoder (Guo et al., 2024), CodeLlama (Rozière et al., 2023), OCTOCODER (Muennighoff et al., 2023), CodeShell (Xie et al., 2024), MagiCoder (Wei et al., 2023), WizardCoder (Luo et al., 2023), Codegemma (Gemma Team, 2024) and Granite (Mishra et al., 2024). Furthermore, we further fine-tune MCODER based on CodeQwen1.5 and DeepSeekCoder to explore the language transfer capabilities of code LLMs.

304 **Evaluation Metrics.** We assess the models by executing the code and evaluating it using the  
305 Pass@1 metric. Pass@1 is a widely recognized measure in machine learning, particularly for code  
306 generation, as it gauges the model’s accuracy in producing correct solutions on the first attempt.

307 **Instruction Corpora.** The resulting dataset, MCEVAL-INSTRUCT (110K samples), is comprised  
308 of created question-answer pairs and open-source collection (Wei et al., 2023). We apply data  
309 decontamination before training our MCODER. Following Li et al. (2023); Wei et al. (2023), we adopt  
310 the N-gram exact match decontamination method with MCEVAL, HumanEval(Chen et al., 2021),  
311 MultiPL-E(Cassano et al., 2023), MBPP(Austin et al., 2021). For supervised fine-tuning (SFT), we  
312 utilize CodeQwen-1.5 as the foundational code LLMs. Specifically, we select all Python data from  
313 MCEVAL-INSTRUCT, comprising 50K training samples, for MCODER-Python training.

314 **Optimization & Evaluation.** Our MCODER based on CodeQwen1.5 are trained for 2 epochs with a  
315 cosine scheduler, starting at a learning rate of 2e-5 (3% warmup steps). We use AdamW (Loshchilov  
316 & Hutter, 2017) as the optimizer and a batch size of 512 (max length 4096). We adopt the greedy  
317 Pass@1 (%) metric (Kulal et al., 2019; Chen et al., 2021) for evaluations. For closed-source LLMs,  
318 the answers are generated by the official API. For code explanation, we prompt the LLM to describe  
319 the code and then restore the descriptions to the original code. (Details can be found in Appendix  
320 A.6).

## 321 4.2 MAIN RESULTS

322 **Multilingual Code Generation.** Table 3 shows the Pass@1 results of various models on MCEVAL  
323 for multilingual code generation task. The results reveal a significant disparity between closed-source

Table 4: Pass@1 (%) scores of different models for multilingual code completion tasks on MCEVAL.  
“Avgall” represents the average scores of all code languages.

|     |                          | Single-line Completion |             |             |             |             |             |             |             |             |             |             |             |             |             |             |             |             |             |             |             |             |
|-----|--------------------------|------------------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
|     | Method                   | Size                   | AWK         | C           | C++         | C#          | Clip        | Coffee      | Dart        | Erlang      | Elixir      | Erlang      | Fortran     | F#          | Go          | Groovy      | Haskell     | Html        | Java        | JS          | Json        | Julia       |
| 324 | GPT-4 Turbo (231106)     | ■                      | 90.0        | 74.4        | 75.6        | 89.3        | 91.0        | 97.5        | 76.5        | 82.4        | 82.4        | 82.4        | 82.4        | 82.4        | 92.6        | 76.2        | 57.1        | 93.9        | 80.0        | 93.4        | 93.3        |             |
| 325 | GPT-3.5 Turbo (240125)   | ●                      | 14.3        | 29.0        | 20.7        | 35.8        | 28.5        | 44.7        | 19.1        | 6.3         | 14.3        | 37.5        | 67.1        | 76.3        | 11.1        | 40.1        | 35.0        | 2.2         | 76.7        | 3.0         | 0.0         |             |
| 326 | DeepSeekCoder-Instruct   | 33B                    | 50.0        | 61.0        | 72.0        | 79.8        | 62.5        | 78.3        | 70.4        | 63.2        | 66.2        | 45.8        | 68.4        | 60.7        | 67.1        | 83.0        | 38.1        | 52.4        | 82.7        | 65.0        | 77.8        | 83.3        |
| 327 | OCTOCODER                | 16B                    | 42.9        | 47.6        | 52.4        | 82.1        | 35.7        | 56.3        | 60.2        | 34.2        | 8.8         | 0.0         | 48.2        | 58.5        | 73.4        | 0.1         | 2.4         | 81.6        | 56.3        | 6.9         | 0.0         |             |
| 328 | StarCoder2-instruct-V0.1 | 15B                    | 28.6        | 74.4        | 81.7        | 86.9        | 71.4        | 7.5         | 82.7        | 68.4        | 75.0        | 62.5        | 76.5        | 64.3        | 82.9        | 88.3        | 47.6        | 0.0         | 91.8        | 83.8        | 16.7        | 83.3        |
| 329 | WizardCoder-V1.0         | 15B                    | 21.4        | 37.8        | 34.9        | 36.1        | 51.3        | 37.8        | 6.6         | 28.6        | 27.9        | 14.7        | 9.2         | 10.7        | 59.8        | 66.0        | 14.3        | 16.7        | 38.8        | 41.3        | 0.0         | 53.3        |
| 330 | Qwen1.5-Chat             | 14B                    | 35.7        | 50.0        | 54.9        | 61.9        | 19.6        | 51.3        | 37.8        | 28.0        | 38.2        | 35.6        | 30.4        | 46.3        | 51.1        | 28.6        | 28.6        | 59.2        | 50.0        | 54.2        | 53.3        |             |
| 331 | CodeLlama-Instruct       | 7B                     | 75.8        | 77.0        | 78.3        | 79.8        | 82.5        | 85.6        | 83.7        | 82.5        | 41.7        | 41.7        | 41.7        | 41.7        | 46.3        | 68.3        | 68.3        | 26.2        | 21.4        | 52.4        | 61.3        | 66.7        |
| 332 | Yi-1.5-Chat              | 69.0                   | 35.7        | 63.4        | 65.9        | 39.2        | 61.3        | 66.3        | 38.2        | 48.5        | 27.1        | 68.4        | 35.7        | 59.8        | 85.1        | 38.1        | 21.4        | 77.6        | 58.8        | 69.0        | 75.6        |             |
| 333 | CodeQwen1.5-Chat         | 7B                     | 35.7        | 68.3        | 63.4        | 76.2        | 67.9        | 47.5        | 78.6        | 35.5        | 72.1        | 45.8        | 57.1        | 60.7        | 73.2        | 69.1        | 31.0        | 42.9        | 90.8        | 68.8        | 69.4        | 88.9        |
| 334 | Magicode-S-DS            | 7B                     | 71.4        | 67.1        | 72.0        | 84.5        | 71.4        | 77.5        | 80.6        | 69.7        | 79.4        | 62.5        | 85.7        | 85.7        | 75.6        | 96.8        | 52.4        | 40.5        | 94.9        | 72.5        | 73.6        | 74.4        |
| 335 | mCODER                   | 7B                     | 85.7        | 74.4        | 82.9        | 90.5        | 69.6        | 91.3        | 78.6        | 76.9        | 67.7        | 70.8        | 87.4        | 80.5        | 83.0        | 45.2        | 28.6        | 95.9        | 81.3        | 54.2        | 94.4        |             |
| 336 | Method                   | Kodin                  | Lua         | Md          | Pascal      | Perl        | PHP         | Power       | Python      | R           | Racket      | Ruby        | Rust        | Scala       | Scheme      | Shell       | Swift       | Tcl         | TS          | VB          | VimL        | Avgall      |
| 337 | GPT-4 Turbo (231106)     | 83.3                   | <b>80.0</b> | <b>28.6</b> | 70.7        | <b>91.4</b> | <b>86.2</b> | <b>92.3</b> | <b>86.7</b> | <b>80.0</b> | <b>88.7</b> | <b>89.5</b> | <b>81.5</b> | <b>87.8</b> | <b>85.3</b> | <b>87.8</b> | <b>86.1</b> | <b>85.9</b> | <b>77.5</b> | <b>76.9</b> | 63.2        | 82.7        |
| 338 | GPT-3.5 Turbo (240125)   | 64.3                   | 58.8        | 21.4        | 20.7        | 52.9        | 20.2        | 65.6        | 84.4        | 23.9        | 41.9        | 47.4        | 64.2        | 20.0        | 59.7        | 54.9        | 63.9        | 57.7        | 46.3        | 28.2        | 80.3        | 43.6        |
| 339 | DeepSeekCoder-Instruct   | 33B                    | 7.13        | 14.3        | 68.3        | 81.9        | 86.7        | 83.3        | 79.3        | 59.7        | 86.8        | 75.0        | 84.1        | 72.1        | 80.5        | 79.2        | 83.3        | 57.5        | 70.5        | 67.1        | 72.4        |             |
| 340 | OCTOCODER                | 70.2                   | 50.0        | 3.6         | 15.9        | 32.9        | 0.0         | 28.9        | 78.9        | 2.2         | 38.7        | 15.8        | 0.0         | 56.1        | 41.2        | 48.8        | 27.8        | 37.2        | 67.5        | 59.0        | 69.7        | 39.2        |
| 341 | StarCoder2-instruct-V0.1 | 15B                    | 85.7        | 61.0        | 0.0         | 55.6        | 75.6        | 80.0        | 86.0        | 80.0        | 87.8        | 86.0        | 86.0        | 86.0        | 87.8        | 86.0        | 87.8        | 86.0        | 87.8        | 86.0        | 87.8        |             |
| 342 | WizardCoder-V1.0         | 15B                    | 42.9        | 28.0        | 31.7        | 34.5        | 3.6         | 13.8        | 21.4        | 0.0         | 20.6        | 10.4        | 11.2        | 10.7        | 39.0        | 61.7        | 9.5         | 14.3        | 24.5        | 20.0        | 1.4         | 28.9        |
| 343 | Qwen1.5-Chat             | 14B                    | 7.1         | 40.2        | 41.5        | 42.9        | 7.1         | 26.3        | 27.6        | 5.3         | 14.7        | 4.2         | 23.5        | 12.5        | 36.6        | 47.9        | 11.9        | 14.3        | 49.0        | 42.5        | 41.7        | 45.6        |
| 344 | CodeLlama-Instruct       | 13B                    | 50.0        | 39.0        | 60.0        | 63.1        | 25.0        | 40.0        | 31.6        | 32.2        | 21.1        | 18.8        | 16.3        | 23.2        | 35.4        | 50.0        | 19.0        | 7.1         | 41.8        | 45.0        | 54.2        | 45.6        |
| 345 | CodeLlama-Instruct       | 9B                     | 35.7        | 47.6        | 48.8        | 64.3        | 14.3        | 52.5        | 42.9        | 18.4        | 20.6        | 16.7        | 51.0        | 25.0        | 39.0        | 73.4        | 23.8        | 16.7        | 61.2        | 55.0        | 75.0        | 60.0        |
| 346 | CodeQwen1.5-Chat         | 7B                     | 28.6        | 36.6        | 43.1        | 42.9        | 17.9        | 46.3        | 16.3        | 17.1        | 22.1        | 25.0        | 22.4        | 28.6        | 37.8        | 43.6        | 23.8        | 31.0        | 24.5        | 12.5        | 73.6        | 34.4        |
| 347 | CodeQwen1.5-Chat         | 7B                     | 0.0         | 32.4        | 32.7        | 37.8        | 9.6         | 25.0        | 31.4        | 32.0        | 34.2        | 34.2        | 34.8        | 31.0        | 42.9        | 58.7        | 31.0        | 35.7        | 78.6        | 71.3        | 72.2        | 68.9        |
| 348 | Magicode-S-DS            | 7B                     | 21.4        | 64.6        | 64.6        | 61.0        | 51.8        | 55.0        | 59.2        | 52.6        | 60.3        | 45.8        | 69.4        | 66.1        | 62.5        | 91.5        | 35.5        | 16.7        | 78.6        | 65.0        | 69.4        | 62.2        |
| 349 | mCODER                   | 7B                     | 21.4        | 57.3        | 53.7        | 78.6        | 42.9        | 71.3        | 56.1        | 50.0        | 57.4        | 47.9        | 45.9        | 51.8        | 56.1        | 80.9        | 23.8        | 23.8        | 80.6        | 75.0        | 62.5        | 72.2        |
| 350 | Method                   | Kodin                  | Lua         | Md          | Pascal      | Perl        | PHP         | Power       | Python      | R           | Racket      | Ruby        | Rust        | Scala       | Scheme      | Shell       | Swift       | Tcl         | TS          | VB          | VimL        | Avgall      |
| 351 | GPT-4 Turbo (231106)     | <b>84.5</b>            | <b>66.3</b> | <b>32.1</b> | <b>67.1</b> | <b>84.3</b> | <b>85.1</b> | <b>94.4</b> | <b>83.0</b> | <b>80.4</b> | <b>64.5</b> | <b>78.9</b> | <b>81.5</b> | <b>84.1</b> | <b>90.2</b> | <b>81.9</b> | <b>78.2</b> | <b>75.0</b> | <b>75.6</b> | <b>52.6</b> | <b>84.4</b> |             |
| 352 | GPT-3.5 Turbo (240125)   | 69.0                   | 58.8        | 21.4        | 50.0        | 60.0        | 39.1        | 52.2        | 43.3        | 27.4        | 2.6         | 0.0         | 34.1        | 23.5        | 32.9        | 33.3        | 25.6        | 45.0        | 33.3        | 35.5        | 27.1        |             |
| 353 | DeepSeekCoder-Instruct   | 53.6                   | 30.0        | 3.7         | 35.0        | 57.0        | 8.0         | 26.7        | 52.6        | 43.3        | 27.4        | 2.6         | 0.0         | 34.1        | 23.5        | 32.9        | 33.3        | 25.6        | 45.0        | 33.3        | 35.5        |             |
| 354 | OCTOCODER                | 16B                    | 43.4        | 35.7        | 32.0        | 33.0        | 21.2        | 23.2        | 30.6        | 7.9         | 5.0         | 0.0         | 26.2        | 26.0        | 37.0        | 37.0        | 30.4        | 40.0        | 19.0        | 1.1         |             |             |
| 355 | StarCoder2-instruct-V0.1 | 15B                    | 76.2        | 48.8        | 0.0         | 64.6        | 64.3        | 80.9        | 76.7        | 75.6        | 72.7        | 48.4        | 51.3        | 64.1        | 60.8        | 68.3        | 47.2        | 50.0        | 60.0        | 60.3        | 47.4        | 58.3        |
| 356 | WizardCoder-V1.0         | 46.4                   | 15.0        | 0.0         | 7.3         | 30.0        | 31.9        | 22.2        | 38.9        | 33.7        | 0.0         | 23.7        | 16.6        | 14.6        | 0.0         | 15.9        | 19.4        | 23.1        | 23.8        | 37.2        | 14.5        | 22.8        |
| 357 | Qwen1.5-Chat             | 36.9                   | 30.3        | 3.6         | 19.5        | 43.3        | 48.9        | 34.4        | 32.0        | 39.1        | 8.1         | 31.6        | 29.3        | 12.0        | 32.9        | 37.8        | 42.4        | 21.8        | 38.8        | 24.4        | 19.7        | 29.9        |
| 358 | CodeLlama-Instruct       | 47.6                   | 37.5        | 0.0         | 27.0        | 31.4        | 33.6        | 13.6        | 40.0        | 60.0        | 46.2        | 22.6        | 39.5        | 35.0        | 46.2        | 47.6        | 42.5        | 47.4        | 34.2        | 39.0        |             |             |
| 359 | CodeQwen1.5-Chat         | 44.0                   | 20.0        | 0.0         | 14.6        | 30.0        | 8.5         | 38.9        | 45.6        | 27.2        | 22.6        | 31.6        | 21.2        | 40.2        | 47.4        | 39.0        | 42.5        | 47.4        | 42.5        | 47.4        | 46.6        |             |
| 360 | Magicode-S-DS            | 69.0                   | 51.3        | 3.6         | 47.6        | 61.4        | 70.7        | 67.6        | 57.7        | 52.6        | 46.8        | 44.7        | 67.4        | 69.5        | 38.2        | 76.8        | 70.8        | 60.3        | 56.3        | 46.2        | 40.8        |             |
| 361 | mCODER                   | 66.7                   | 56.3        | 0.0         | 58.5        | 57.1        | 71.3        | 80.0        | 66.7        | 64.1        | 50.0        | 57.9        | 71.3        | 64.4        | 54.4        | 80.5        | 63.9        | 57.7        | 60.0        | 66.7        | 32.9        | 60.7        |
| 362 | Method                   | Kodin                  | Lua         | Md          | Pascal      | Perl        | PHP         | Power       | Python      | R           | Racket      | Ruby        | Rust        | Scala       | Scheme      | Shell       | Swift       | Tcl         | TS          | VB          | VimL        | Avgall      |
| 363 | GPT-4 Turbo (231106)     | <b>90.0</b>            | <b>66.0</b> | <b>34.0</b> | <b>59.0</b> | <b>92.0</b> | <b>80.0</b> | <b>87.0</b> | <b>87.0</b> | <b>80.0</b> | <b>80.0</b> | <b>80.0</b> | <b>72.0</b> | <b>72.0</b> | <b>66.0</b> | <b>9.0</b>  | <b>80.0</b> | <b>86.0</b> | <b>86.0</b> | <b>71.0</b> | <b>67.0</b> | <b>73.0</b> |
| 364 | GPT-3.5 Turbo (240125)   | 68.0                   | 61.0        | 25.0        | 45.0        | 68.0        | 8.0         | 62.0        | 80.0        | 13.0        | 56.0        | 24.0        | 26.4        | 39.0        | 47.0        | 72.0        | 70.0        | 57.0        | 60.0        | 68.0        | 58.0        | 48.1        |
| 365 | DeepSeekCoder-Instruct   | 69.0                   | 58.0        | 16.0        | 50.0        | 54.0        | 65.0        | 72.0        | 66.0        | 61.0        | 58.0        | 63.0        | 67.0        | 73.0        | 53.0        | 64.0        | 68.0        | 59.0        | 60.0        | 62.0        | 59.5        |             |
| 366 | OCTOCODER                | 47.0                   | 47.0        | 0.0         | 24.0        | 30.0        | 31.0        | 51.0        | 31.0        | 37.0        | 2.0         | 0.0         | 37.0        | 25.0        | 32.0        | 38.0        | 39.0        | 17.7        | 26.4        | 43.0        |             |             |
| 367 | StarCoder2-instruct-V0.1 | 15B                    | 34.0        | 0.0         | 58.0        | 69.0        | 69.0        | 70.0        | 67.0        | 67.0        | 58.0        | 58.0        | 67.0        | 67.0        | 67.0        | 67.0        | 67.0        | 22.0        | 17.0        | 23.0        | 32.0        |             |
| 368 | WizardCoder-V1.0         | 15B                    | 34.0        | 48.0        | 6.0         | 74.0        | 62.0        | 40.4        | 47.1        | 44.0        | 64.0        | 56.0        | 58.0        | 68.0        | 56.0        | 62.0        | 60.0        | 51.0        | 63.0        | 59.0        | 58.8        |             |
| 369 | Qwen1.5-Chat             | 33.0                   | 47.0        | 8.0         | 30.0        | 37.0        | 44.0        | 38.0        | 45.0        | 46.0        | 25.0        | 41.0        | 32.1        | 14.0        | 12.0        | 30.0        | 42.0        | 23.0        | 22.0        | 38.0        | 39.6        |             |
| 370 | CodeLlama-Instruct       | 45.0                   | 44.0        | 0.0         | 27.0        | 51.0        | 38.0        | 49.0        | 47.0        | 44.0        | 38.0        | 48.0        | 42.0        | 38.0        | 42.0        | 47.0        | 33.0        | 35.0        | 40.0        | 44.0        | 40.5        |             |
| 371 | CodeQwen1.5-Chat         | 56.0                   | 52.0        | 12.0        | 49.0        | 57.0        | 61.0        | 63.0        | 64.0        | 56.0        | 33.0        | 57.0        | 37.7        | 28.0        | 30.0        | 43.0        | 66.0        | 37.0        | 47.3        |             |             |             |
| 372 | Magicode-S-DS            | 7B                     | 30.0        | 58.0        | 36.0        | 44.0        | 30.0        | 32.0        | 19.6        | 24.0        | 34.0        | 38.0        | 32.0        | 42.0        | 50.0        | 34.0        | 20.0        | 32.0        | 38.0        | 50.0        | 38.0        |             |
| 373 | mCODER                   | 7B                     | 48.0        | 52.0        | 6.0         | 78.0        | 60.0        | 62.0        | 72.0        | 69.0        | 64.0        | 72.0        | 67.0        | 72.0        | 78.0        | 58.0        | 70.0        | 88.0        | 64.0        | 40.0        | 83.0        | 58.0        |
| 374 | Method                   | Kodin                  | Lua         | Md          | Pascal      | Perl        | PHP         | Power       | Python      | R           | Racket      | Ruby        | Rust        | Scala       | Scheme      | Shell       | Swift       | Tcl         | TS          | VB          | VimL        | Avgall      |
| 375 | GPT-4 Turbo (231106)     | <b>94.0</b>            | <b>44.0</b> | <b>70.0</b> | <b>90.0</b> | <b></b>     |             |             |             |             |             |             |             |             |             |             |             |             |             |             |             |             |

378 Table 5: Pass@1 (%) scores of different models for multilingual code explanation tasks on MCEVAL.  
379 “Avg<sub>all</sub>” represents the average scores of all code languages.

| Method                     | Size | AWK  | C    | C++  | C#   | Cisp | Coffee | Dart | Erlip | Elixir | Erlang | Fortran | F#   | Go   | Groovy | Haskell | Html | Java | JS   | Json | Julia |      |
|----------------------------|------|------|------|------|------|------|--------|------|-------|--------|--------|---------|------|------|--------|---------|------|------|------|------|-------|------|
| GPT-4o (240513)            | ■    | 74.0 | 68.0 | 72.0 | 72.0 | 78.0 | 76.0   | 54.9 | 60.0  | 66.0   | 38.0   | 62.0    | 78.0 | 74.0 | 88.0   | 86.0    | 16.0 | 79.2 | 68.0 | 24.0 | 76.0  |      |
| GPT-4 Turbo (231106)       | ■    | 80.0 | 76.0 | 72.0 | 52.0 | 60.0 | 54.0   | 54.9 | 50.0  | 64.0   | 42.0   | 28.0    | 56.0 | 54.0 | 84.0   | 54.0    | 10.0 | 81.1 | 74.0 | 18.0 | 70.0  |      |
| GPT-3.5 Turbo (240125)     | ■    | 8.0  | 68.0 | 64.0 | 62.0 | 56.0 | 66.0   | 29.2 | 26.0  | 42.0   | 34.0   | 42.0    | 42.0 | 60.0 | 54.0   | 74.0    | 68.0 | 14.0 | 47.2 | 48.0 | 26.0  | 62.0 |
| Yi-Large-Turbo             | ■    | 76.0 | 52.0 | 48.0 | 50.0 | 60.0 | 25.5   | 44.0 | 42.0  | 44.0   | 42.0   | 42.0    | 60.0 | 54.0 | 84.0   | 54.0    | 14.0 | 37.7 | 36.0 | 8.0  | 28.0  |      |
| DeepSpeech-Coder-Instruct  | 33B  | 70.0 | 62.0 | 66.0 | 78.0 | 56.0 | 68.0   | 45.1 | 44.0  | 58.0   | 40.0   | 42.0    | 64.0 | 64.0 | 56.0   | 88.0    | 52.0 | 8.0  | 67.9 | 52.0 | 32.0  | 70.0 |
| OCTOCODER                  | 16B  | 32.0 | 32.0 | 32.0 | 26.0 | 48.0 | 4.0    | 5.9  | 22.0  | 32.0   | 34.0   | 22.0    | 16.0 | 52.0 | 42.0   | 34.0    | 2.0  | 35.8 | 44.0 | 14.0 | 42.0  |      |
| Qwen1.5-Chat               | 14B  | 50.0 | 52.0 | 30.0 | 42.0 | 36.0 | 26.0   | 23.5 | 28.0  | 24.0   | 12.0   | 14.0    | 34.0 | 44.0 | 28.0   | 46.0    | 8.0  | 35.8 | 44.0 | 14.0 | 42.0  |      |
| CodeLlama-Instruct         | 13B  | 32.0 | 42.0 | 34.0 | 44.0 | 44.0 | 4.0    | 7.8  | 16.0  | 30.0   | 38.0   | 16.0    | 30.0 | 32.0 | 32.0   | 40.0    | 10.0 | 28.3 | 36.0 | 10.0 | 24.0  |      |
| Yi-1.5-Chat                | 9B   | 38.0 | 62.0 | 60.0 | 58.0 | 38.0 | 38.0   | 29.4 | 14.0  | 50.0   | 26.0   | 36.0    | 20.0 | 52.0 | 68.0   | 44.0    | 46.0 | 22.0 | 66.0 |      |       |      |
| CodeLlama-Instruct         | 7B   | 34.0 | 34.0 | 32.0 | 42.0 | 34.0 | 8.0    | 11.8 | 22.0  | 28.0   | 24.0   | 14.0    | 26.0 | 40.0 | 22.0   | 36.0    | 0.0  | 34.0 | 22.0 | 8.0  | 14.0  |      |
| CodeQwen1.5-Chat           | 7B   | 58.0 | 62.0 | 58.0 | 50.0 | 54.0 | 56.0   | 17.6 | 26.0  | 50.0   | 52.0   | 34.0    | 50.0 | 48.0 | 76.0   | 46.0    | 2.0  | 67.9 | 58.0 | 18.0 | 56.0  |      |
| MagicCoder-S-DS            | 7B   | 58.0 | 62.0 | 56.0 | 70.0 | 58.0 | 4.0    | 37.3 | 28.0  | 50.0   | 58.0   | 54.0    | 66.0 | 58.0 | 86.0   | 60.0    | 4.0  | 84.9 | 60.0 | 2.0  | 54.0  |      |
| <b>MCODER (Our Method)</b> | 7B   | 52.0 | 62.0 | 56.0 | 68.0 | 70.0 | 48.0   | 33.3 | 64.0  | 40.0   | 40.0   | 36.0    | 70.0 | 66.0 | 58.0   | 6.0     | 71.7 | 60.0 | 28.0 | 60.0 |       |      |
| <b>MCODER (Our Method)</b> |      | 62.0 | 52.0 | 0.0  | 48.0 | 46.0 | 62.0   | 58.0 | 38.0  | 50.0   | 74.0   | 67.9    | 32.0 | 46.0 | 60.0   | 50.0    | 58.0 | 58.0 | 50.0 | 51.4 |       |      |

393 exhibits clear improvement over the base model in nearly all the studied programming languages. It is  
394 noteworthy that MCODER, despite being trained with very limited multilingual data, still outperforms  
395 other large language models (LLMs) of similar or even larger sizes.

397 **Multilingual Code Explanation.** Table 5 displays the Pass@1 results for multilingual code explanation  
398 tasks. The results show that GPT models still significantly outperform open-source models in  
399 the code explanation task. For markup languages (Json and Markdown), the complexity of the code  
400 structure makes it difficult to describe accurately in natural language, resulting in generally poorer  
401 performance. Code LLMs need instruction-following capabilities for such complex structures.

403 **Multilingual Code Completion.** The completion tasks consist of *single-line completion*, *multi-line*  
404 *completion*, *span completion*, and *span completion (light)*. As shown in Table 4, the Pass@1 results  
405 for multilingual code completion tasks indicate that GPT-4 Turbo still achieves the best performance.  
406 Additionally, since this task is relatively easier compared to code generation, some open-source  
407 models perform comparably to GPT-4 Turbo in certain programming languages.

## 5 FURTHER ANALYSIS

411 **Programming Classification.** In Figure 5, we categorize the programming languages of MCEVAL  
412 into 5 programming paradigms and 11 application scenarios and summarize the performance of code  
413 LLMs on the code generation task in Figure 6. It can be observed that code LLMs generally perform  
414 better in object-oriented and multi-paradigm programming languages (high-resource languages),  
415 while perform worse in functional and procedural programming languages (low-resource Languages).  
416 In areas like web development and scientific computing, the gap between open-source and closed-  
417 source models is narrowing. However, for application scenarios, there is still a substantial gap  
418 between open-source models and the closed-source GPT-4 series in low-resource languages related to  
419 scripting, mobile development, and educational research. MCODER performs superior over multiple  
420 same-size models and even some larger open-source models.

421 **Unbalance on Different Languages.** We compare the results of several open-source models on  
422 the MultiPL-E multilingual benchmark with corresponding languages on MCEVAL. We obtained  
423 scores for 11 programming languages (including Python, Java, JavaScript, C++, PHP, Rust, Swift,  
424 R, Lua, Racket, Julia) from the BigCode leaderboard.<sup>1</sup> As shown in Figure 7(1), due to the sim-  
425 plicity of Python language tasks in this dataset, many models exhibit significant score discrepancies  
426 between the two benchmarks. Figure 7(2) highlights a majority of models within the blue circle,  
427 indicating that the current state-of-the-art performance of most models primarily lies in high-resource  
428 languages like Python, while their proficiency in low-resource languages awaits further exploration  
429 and enhancement. By examining Figure 7(2) and (3), it becomes evident that all LLMs demonstrate  
430 consistent multilingual capabilities between MultiPL-E and MCEVAL.

431 <sup>1</sup><https://huggingface.co/spaces/bigcode/bigcode-models-leaderboard>

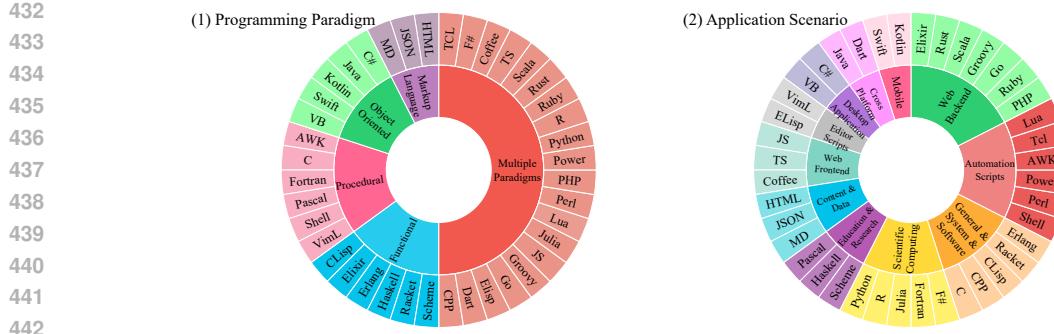


Figure 5: Classification of MCEVAL. The programming languages in MCEVAL can be categorized into 5 programming paradigms and 11 application scenarios.

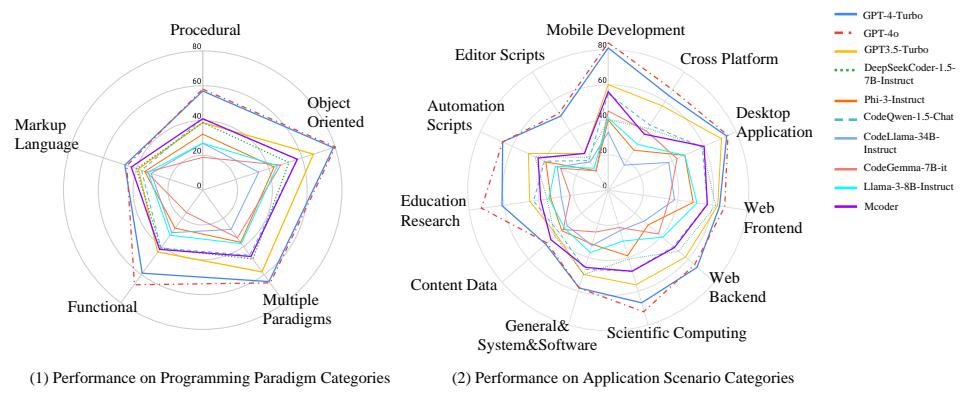


Figure 6: The performance of models in code completion tasks under different categories.

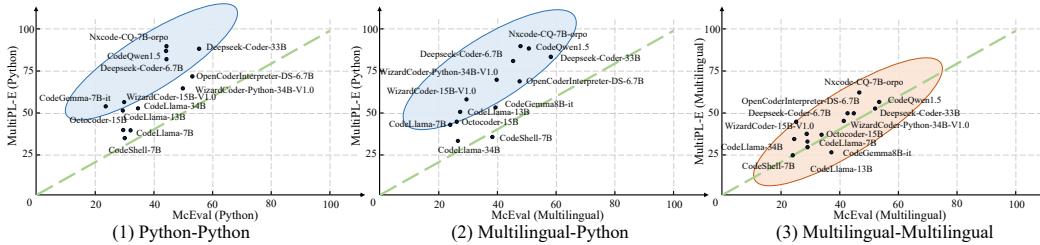


Figure 7: Unbalance performance on different languages across MultiPL-E and our MCEVAL.

**Cross-lingual Transfer.** We fine-tune the CodeQwen-1.5 model using Python-only data in MCEVAL-INSTRUCT and compare it with MCODER. In Figure 8, CodeQwen-1.5 performs well in most high-resource languages, but CodeQwen without alignment exhibits unsatisfactory results in some low-resource languages due to the inability to follow instructions. As such, with fine-tuning using only Python data, CodeQwen-1.5-Python improves significantly across most languages. It shows that the CodeQwen foundation model already possesses strong coding capabilities but lacks adequate instruction-following skills. Therefore, fine-tuning with Python-only data can still effectively transfer instruction-following abilities to other languages, resulting in superior multilingual performance.

**Difficulty of MCEVAL.** Based on algorithmic complexity, we classify MCEVAL into three levels (Easy/Medium/Hard). In Figure 9, we conduct a statistical analysis of CodeQwen-1.5-Chat’s performance on code generation tasks across various languages. For most languages, the code LLM can answer the majority of easy questions but struggles with medium and hard ones.

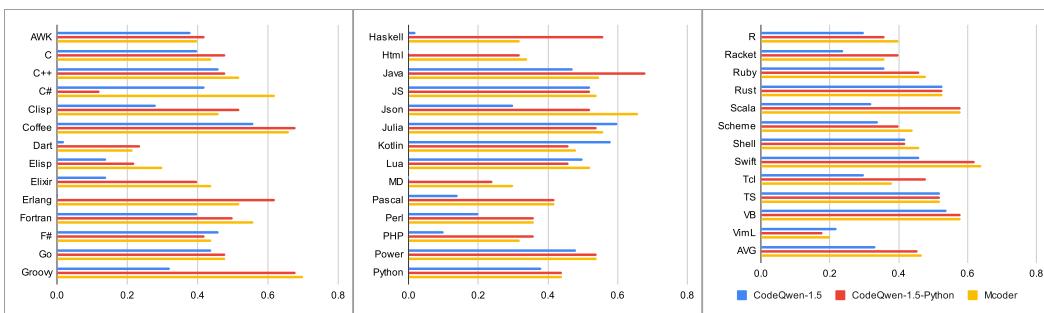


Figure 8: Cross-lingual transferability of LLMs among different languages. We fine-tune CodeQwen-1.5 using Python data in MCEVAL-INSTRUCT and OSS-Instruct to create CodeQwen-1.5-Python.

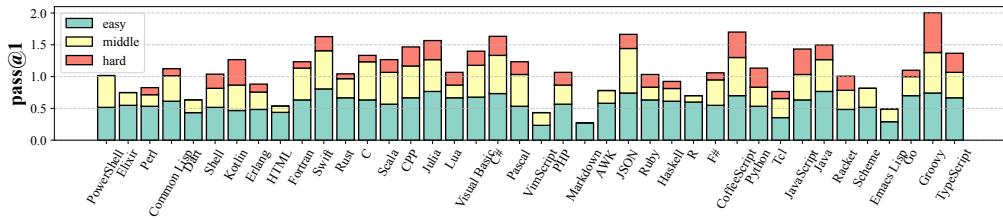


Figure 9: CodeQwen-1.5-Chat performance on MCEVAL for problems of different difficulty levels.

## 6 RELATED WORK

For the field of soft engineering, code LLMs (Feng et al., 2020; Chen et al., 2021; Scao et al., 2022; Li et al., 2022; Allal et al., 2023; Fried et al., 2022; Wang et al., 2021; Zheng et al., 2024; Guo et al., 2024) pre-trained on billions of code snippets, such as StarCoder (Li et al., 2023; Lozhkov et al., 2024), CodeLlama (Rozière et al., 2023), DeepSeekCoder (Guo et al., 2024), and Code-Qwen (Hui et al., 2024). The development and refinement of code LLMs have been pivotal in automating software development tasks, and supporting code generation/translation/summarization.

Many benchmarks (Yu et al., 2024; Yin et al., 2023; Khan et al., 2023; Orlanski et al., 2023; Jain et al., 2024) have been woven to accurately assess code quality, functionality, and efficiency, such as HumanEval (Chen et al., 2021), MBPP (Austin et al., 2021), their upgraded version EvalPlus (Liu et al., 2023b). Studies have explored a variety of approaches, ranging from static evaluation using text matching to dynamic methods that involve code execution under a controlled environment. The current benchmarks support code LLMs to evaluate a series of different types of tasks, such as code understanding, function calling (Zhuo et al., 2024), code repair (Lin et al., 2017; Tian et al., 2024; Jimenez et al., 2023; Zhang et al., 2023; Prener & Robbes, 2023; He et al., 2022), code translation (Yan et al., 2023). Some works focus on the multilingual scenarios (Wang et al., 2023; Athiwaratkun et al., 2023; Peng et al., 2024; Zheng et al., 2023b) by extending the Python-only HumanEval/MBPP benchmark (e.g. MultiPL-E (Cassano et al., 2023)), which is challenged by the number of the languages.

## 7 CONCLUSION

In this work, we push a significant advancement in the assessment of code LLMs by proposing the first massively multilingual code evaluation benchmark (MCEVAL) by involving an annotation and verification process conducted by professional developers, which spans 40 programming languages and helps comprehensively tackle various tasks, including code generation, explanation, and completion. The multilingual SFT on created instruction corpora MCEVAL-INSTRUCT further emphasizes the proficiency of LLMs in multiple coding languages. Systematic evaluations of existing code LLMs on MCEVAL illuminate the performance disparities among open-source and closed-source models. Extensive multilingual multitask assessment on MCEVAL provides a realistic and comprehensive measurement of code LLMs, marking a leap forward for developers utilizing AI techniques to understand and generate code effectively across a wide spectrum of programming languages.

540 REFERENCES  
541

- 542 Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany  
543 Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat Behl, et al. Phi-3 technical report:  
544 A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*, 2024.  
545 URL <https://arxiv.org/abs/2404.14219>.
- 546 Meta AI. Introducing meta llama 3: The most capable openly available llm to date. <https://ai.meta.com/blog/meta-llama-3/>, apr 2024.  
547  
548
- 549 Loubna Ben Allal, Raymond Li, Denis Kocetkov, Chenghao Mou, Christopher Akiki, Carlos Munoz  
550 Ferrandis, Niklas Muennighoff, Mayank Mishra, Alex Gu, Manan Dey, et al. SantaCoder: Don't  
551 reach for the stars! *arXiv preprint arXiv:2301.03988*, 2023. URL <https://arxiv.org/abs/2301.03988>.  
552
- 553 Ben Athiwaratkun, Sanjay Krishna Gouda, Zijian Wang, Xiaopeng Li, Yuchen Tian, Ming Tan,  
554 Wasi Uddin Ahmad, Shiqi Wang, Qing Sun, Mingyue Shang, Sujan Kumar Gonugondla, Hantian  
555 Ding, Varun Kumar, Nathan Fulton, Arash Farahani, Siddhartha Jain, Robert Giaquinto, Haifeng  
556 Qian, Murali Krishna Ramanathan, and Ramesh Nallapati. Multi-lingual evaluation of code  
557 generation models. In *The Eleventh International Conference on Learning Representations, ICLR*  
558 *2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. URL <https://openreview.net/forum?id=Bo7eeXm6An8>.  
559
- 560 Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan,  
561 Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language  
562 models. *arXiv preprint arXiv:2108.07732*, 2021. URL <https://arxiv.org/abs/2108.07732>.  
563
- 564 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge,  
565 Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu,  
566 Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan,  
567 Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin  
568 Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng  
569 Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou,  
570 Xiaohuan Zhou, and Tianhang Zhu. Qwen technical report. *arXiv preprint arXiv:2309.16609*,  
571 <https://arxiv.org/abs/2309.16609>, 2023. URL <https://arxiv.org/abs/2309.16609>.  
572
- 573 Mohammad Bavarian, Heewoo Jun, Nikolas Tezak, John Schulman, Christine McLeavey, Jerry  
574 Tworek, and Mark Chen. Efficient training of language models to fill in the middle. *arXiv preprint*  
575 *arXiv:2207.14255*, 2022. URL <https://arxiv.org/abs/2207.14255>.  
576
- 577 Luca Beurer-Kellner, Marc Fischer, and Martin T. Vechev. Prompting is programming: A query  
578 language for large language models. *Proc. ACM Program. Lang.*, 7(PLDI):1946–1969, 2023. doi:  
579 10.1145/3591300. URL <https://doi.org/10.1145/3591300>.
- 580 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhari-  
581 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language  
582 models are few-shot learners. *Advances in neural information processing systems*, 33:  
583 1877–1901, 2020. URL <https://proceedings.neurips.cc/paper/2020/hash/1457c0d6bfcb4967418fb8ac142f64a-Abstract.html>.  
584
- 585 Federico Cassano, John Gouwar, Daniel Nguyen, Sydney Nguyen, Luna Phipps-Costin, Donald  
586 Pinckney, Ming-Ho Yee, Yangtian Zi, Carolyn Jane Anderson, Molly Q Feldman, et al. Multipl-e:  
587 A scalable and polyglot approach to benchmarking neural code generation. *IEEE Transactions*  
588 *on Software Engineering*, 2023. URL <https://ieeexplore.ieee.org/abstract/document/10103177>.  
589
- 590 Bei Chen, Fengji Zhang, Anh Nguyen, Daoguang Zan, Zeqi Lin, Jian-Guang Lou, and Weizhu  
591 Chen. Codet: Code generation with generated tests. In *The Eleventh International Conference on*  
592 *Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023.  
593 URL <https://openreview.net/pdf?id=ktrw68Cmu9c>.

- 594 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Pondé de Oliveira Pinto, Jared  
 595 Kaplan, Harrison Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri,  
 596 Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan,  
 597 Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian,  
 598 Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios  
 599 Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Heben Guss, Alex Nichol, Alex Paino,  
 600 Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders,  
 601 Christopher Hesse, Andrew N. Carr, Jan Leike, Joshua Achiam, Vedant Misra, Evan Morikawa,  
 602 Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob  
 603 McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating  
 604 large language models trained on code. *arXiv preprint arXiv:2107.03374*, abs/2107.03374, 2021.  
 605 URL <https://arxiv.org/abs/2107.03374>.
- 606 Wenhui Chen, Xueguang Ma, Xinyi Wang, and William W. Cohen. Program of thoughts prompting:  
 607 Disentangling computation from reasoning for numerical reasoning tasks. *arXiv preprint  
 608 arXiv:2211.12588*, abs/2211.12588, 2022. doi: 10.48550/ARXIV.2211.12588. URL <https://doi.org/10.48550/arXiv.2211.12588>.
- 609 Chris Cummins, Volker Seeker, Dejan Grubisic, Mostafa Elhoushi, Youwei Liang, Baptiste Rozière,  
 610 Jonas Gehring, Fabian Gloeckle, Kim M. Hazelwood, Gabriel Synnaeve, and Hugh Leather. Large  
 611 language models for compiler optimization. *arXiv preprint arXiv:2309.07062*, abs/2309.07062,  
 612 2023. doi: 10.48550/ARXIV.2309.07062. URL <https://doi.org/10.48550/arXiv.2309.07062>.
- 613 Gautier Dagan, Frank Keller, and Alex Lascarides. Dynamic planning with a LLM. *arXiv preprint  
 614 arXiv:2308.06391*, abs/2308.06391, 2023. doi: 10.48550/ARXIV.2308.06391. URL <https://doi.org/10.48550/arXiv.2308.06391>.
- 615 Peng Di, Jianguo Li, Hang Yu, Wei Jiang, Wenting Cai, Yang Cao, Chaoyu Chen, Dajun Chen,  
 616 Hongwei Chen, Liang Chen, Gang Fan, Jie Gong, Zi Gong, Wen Hu, Tingting Guo, Zhichao Lei,  
 617 Ting Li, Zheng Li, Ming Liang, Cong Liao, Bingchang Liu, Jiachen Liu, Zhiwei Liu, Shaojun  
 618 Lu, Min Shen, Guangpei Wang, Huan Wang, Zhi Wang, Zhaogui Xu, Jiawei Yang, Qing Ye,  
 619 Gehao Zhang, Yu Zhang, Zelin Zhao, Xunjin Zheng, Hailian Zhou, Lifu Zhu, and Xianying  
 620 Zhu. Codefuse-13b: A pretrained multi-lingual code large language model. *arXiv preprint  
 621 arXiv:2310.06266*, abs/2310.06266, 2023. doi: 10.48550/ARXIV.2310.06266. URL <https://doi.org/10.48550/arXiv.2310.06266>.
- 622 Yihong Dong, Xue Jiang, Zhi Jin, and Ge Li. Self-collaboration code generation via chatgpt. *arXiv  
 623 preprint arXiv:2304.07590*, abs/2304.07590, 2023.
- 624 Zhangyin Feng, Daya Guo, Duyu Tang, Nan Duan, Xiaocheng Feng, Ming Gong, Linjun Shou, Bing  
 625 Qin, Ting Liu, Dixin Jiang, and Ming Zhou. Codebert: A pre-trained model for programming and  
 626 natural languages. In Trevor Cohn, Yulan He, and Yang Liu (eds.), *Findings of the Association for  
 627 Computational Linguistics: EMNLP 2020*, pp. 1536–1547, Online, November 2020. Association  
 628 for Computational Linguistics. doi: 10.18653/v1/2020.findings-emnlp.139. URL <https://aclanthology.org/2020.findings-emnlp.139>.
- 629 Daniel Fried, Armen Aghajanyan, Jessy Lin, Sida I. Wang, Eric Wallace, Freda Shi, Ruiqi Zhong,  
 630 Wen tau Yih, Luke Zettlemoyer, and Mike Lewis. Incoder: A generative model for code infilling  
 631 and synthesis. *arXiv preprint arXiv:2204.05999*, abs/2204.05999, 2022. URL <https://arxiv.org/abs/2204.05999>.
- 632 Google Gemma Team. Gemma: Open models based on gemini research and technology. *arXiv  
 633 preprint arXiv:2403.08295*, 2024. URL <https://arxiv.org/abs/2403.08295>.
- 634 Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao  
 635 Bi, Y Wu, YK Li, et al. Deepseek-coder: When the large language model meets programming –  
 636 the rise of code intelligence. *arXiv preprint arXiv:2401.14196*, 2024. URL <https://arxiv.org/abs/2401.14196>.
- 637 Jingxuan He, Luca Beurer-Kellner, and Martin Vechev. On distribution shift in learning-based bug  
 638 detectors. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and

- 648 Sivan Sabato (eds.), *Proceedings of the 39th International Conference on Machine Learning*,  
 649 volume 162 of *Proceedings of Machine Learning Research*, pp. 8559–8580. PMLR, 17–23 Jul  
 650 2022. URL <https://proceedings.mlr.press/v162/he22a.html>.
- 651  
 652 Binyuan Hui, Jian Yang, Zeyu Cui, Jiaxi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang,  
 653 Bowen Yu, Kai Dang, et al. Qwen2. 5-coder technical report. *arXiv preprint arXiv:2409.12186*,  
 654 2024.
- 655 Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando  
 656 Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free  
 657 evaluation of large language models for code. *arXiv preprint arXiv:2403.07974*, 2024. URL  
 658 <https://arxiv.org/abs/2403.07974>.
- 659  
 660 Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik  
 661 Narasimhan. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint*  
 662 *arXiv:2310.06770*, 2023. URL <https://arxiv.org/abs/2310.06770>.
- 663 Mohammad Abdullah Matin Khan, M Saiful Bari, Xuan Long Do, Weishi Wang, Md Rizwan  
 664 Parvez, and Shafiq Joty. xcodeeval: A large scale multilingual multitask benchmark for code  
 665 understanding, generation, translation and retrieval. *arXiv preprint arXiv:2303.03004*, 2023. URL  
 666 <https://arxiv.org/abs/2303.03004>.
- 667  
 668 Sumith Kulal, Panupong Pasupat, Kartik Chandra, Mina Lee, Oded Padon, Alex Aiken, and Percy S  
 669 Liang. Spoc: Search-based pseudocode to code. *Advances in Neural Information Processing*  
 670 *Systems*, 32, 2019. URL <https://proceedings.neurips.cc/paper/2019/hash/7298332f04ac004a0ca44cc69ecf6f6b-Abstract.html>.
- 671  
 672 Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E.  
 673 Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model  
 674 serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating*  
 675 *Systems Principles*, 2023. URL <https://dl.acm.org/doi/abs/10.1145/3600006.3613165>.
- 676  
 677 Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao  
 678 Mou, Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, Qian Liu, Evgenii Zheltonozhskii,  
 679 Terry Yue Zhuo, Thomas Wang, Olivier Dehaene, Mishig Davaadorj, Joel Lamy-Poirier, João  
 680 Monteiro, Oleh Shliazhko, Nicolas Gontier, Nicholas Meade, Armel Zebaze, Ming-Ho Yee,  
 681 Logesh Kumar Umapathi, Jian Zhu, Benjamin Lipkin, Muhtasham Oblokulov, Zhiruo Wang,  
 682 Rudra Murthy V, Jason Stillerman, Siva Sankalp Patel, Dmitry Abulkhanov, Marco Zocca, Manan  
 683 Dey, Zhihan Zhang, Nour Moustafa-Fahmy, Urvashi Bhattacharyya, Wenhao Yu, Swayam Singh,  
 684 Sasha Luccioni, Paulo Villegas, Maxim Kunakov, Fedor Zhdanov, Manuel Romero, Tony Lee,  
 685 Nadav Timor, Jennifer Ding, Claire Schlesinger, Hailey Schoelkopf, Jan Ebert, Tri Dao, Mayank  
 686 Mishra, Alex Gu, Jennifer Robinson, Carolyn Jane Anderson, Brendan Dolan-Gavitt, Danish  
 687 Contractor, Siva Reddy, Daniel Fried, Dzmitry Bahdanau, Yacine Jernite, Carlos Muñoz Ferrandis,  
 688 Sean Hughes, Thomas Wolf, Arjun Guha, Leandro von Werra, and Harm de Vries. Starcoder:  
 689 may the source be with you! *arXiv preprint arXiv:2305.06161*, abs/2305.06161, 2023. doi:  
 690 10.48550/arXiv.2305.06161. URL <https://arxiv.org/abs/2305.06161>.
- 691  
 692 Yujia Li, David H. Choi, Junyoung Chung, Nate Kushman, Julian Schrittweiser, Rémi Leblond,  
 693 Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, Thomas Hubert, Peter Choy,  
 694 Cyprien de Masson d’Autume, Igor Babuschkin, Xinyun Chen, Po-Sen Huang, Johannes Welbl,  
 695 Sven Gowal, Alexey Cherepanov, James Molloy, Daniel J. Mankowitz, Esme Sutherland Robson,  
 696 Pushmeet Kohli, Nando de Freitas, Koray Kavukcuoglu, and Oriol Vinyals. Competition-level  
 697 code generation with alphacode. *arXiv preprint arXiv:2203.07814*, abs/2203.07814, 2022. URL  
 698 <https://arxiv.org/abs/2203.07814>.
- 699  
 700 Derrick Lin, James Koppel, Angela Chen, and Armando Solar-Lezama. Quixbugs: a multi-lingual  
 701 program repair benchmark set based on the quixey challenge. In *Proceedings Companion of*  
*the 2017 ACM SIGPLAN international conference on systems, programming, languages, and*  
*applications: software for humanity*, pp. 55–56, 2017. URL <https://dl.acm.org/doi/abs/10.1145/3135932.3135941>.

- 702 Chao Liu, Xuanlin Bao, Hongyu Zhang, Neng Zhang, Haibo Hu, Xiaohong Zhang, and Meng Yan.  
 703 Improving chatgpt prompt for code generation. *arXiv preprint arXiv:2305.08360*, abs/2305.08360,  
 704 2023a. URL <https://arxiv.org/abs/2305.08360>.
- 705  
 706 Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. Is your code generated by  
 707 chatgpt really correct? rigorous evaluation of large language models for code generation. *arXiv  
 708 preprint arXiv:2305.01210*, abs/2305.01210, 2023b. URL [https://arxiv.org/abs/2305.  
 709 01210](https://arxiv.org/abs/2305.01210).
- 710  
 711 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint  
 712 arXiv:1711.05101*, 2017. URL <https://arxiv.org/abs/1711.05101>.
- 713 Anton Lozhkov, Raymond Li, Loubna Ben Allal, Federico Cassano, Joel Lamy-Poirier, Nouamane  
 714 Tazi, Ao Tang, Dmytro Pykhtiar, Jiawei Liu, Yuxiang Wei, et al. Starcoder 2 and the stack v2: The  
 715 next generation. *arXiv preprint arXiv:2402.19173*, 2024. URL [https://arxiv.org/abs/  
 716 2402.19173](https://arxiv.org/abs/2402.19173).
- 717  
 718 Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xiubo Geng, Wenxiang Hu, Chongyang Tao, Jing  
 719 Ma, Qingwei Lin, and Daxin Jiang. WizardCoder: Empowering code large language models  
 720 with evol-instruct. *arXiv preprint arXiv:2306.08568*, 2023. URL [https://arxiv.org/abs/  
 721 2306.08568](https://arxiv.org/abs/2306.08568).
- 722  
 723 Mayank Mishra, Matt Stallone, Gaoyuan Zhang, Yikang Shen, Aditya Prasad, Adriana Meza So-  
 724 ria, Michele Merler, Parameswaran Selvam, Saptha Surendran, Shivdeep Singh, et al. Granite  
 725 code models: A family of open foundation models for code intelligence. *arXiv preprint  
 726 arXiv:2405.04324*, 2024.
- 727  
 728 Niklas Muennighoff, Qian Liu, Armel Zebaze, Qinkai Zheng, Binyuan Hui, Terry Yue Zhuo, Swayam  
 729 Singh, Xiangru Tang, Leandro von Werra, and Shayne Longpre. OctoPack: Instruction tuning  
 730 code large language models. *arXiv preprint arXiv:2308.07124*, abs/2308.07124, 2023. URL  
<https://arxiv.org/abs/2308.07124>.
- 731  
 732 Fionn Murtagh and Pedro Contreras. Algorithms for hierarchical clustering: an overview. *Wiley  
 733 Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 2(1):86–97, 2012. doi:  
[10.1002/WIDM.53](https://doi.org/10.1002/WIDM.53). URL <https://doi.org/10.1002/widm.53>.
- 734  
 735 Erik Nijkamp, Bo Pang, Hiroaki Hayashi, Lifu Tu, Huan Wang, Yingbo Zhou, Silvio Savarese, and  
 736 Caiming Xiong. CodeGen: An open large language model for code with multi-turn program  
 737 synthesis. In *The Eleventh International Conference on Learning Representations, ICLR 2023,  
 738 Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. URL [https://openreview.net/  
 739 forum?id=iaYcJKpY2B\\_](https://openreview.net/forum?id=iaYcJKpY2B_).
- 740  
 741 OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. URL <https://arxiv.org/abs/2303.08774>.
- 742  
 743 Gabriel Orlanski, Kefan Xiao, Xavier Garcia, Jeffrey Hui, Joshua Howland, Jonathan Malmaud, Jacob  
 744 Austin, Rishabh Singh, and Michele Catasta. Measuring the impact of programming language  
 745 distribution. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan  
 746 Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine  
 747 Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 26619–26645. PMLR, 23–  
 748 29 Jul 2023. URL <https://proceedings.mlr.press/v202/orlanski23a.html>.
- 749  
 750 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong  
 751 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser  
 752 Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F Christiano, Jan  
 753 Leike, and Ryan Lowe. Training language models to follow instructions with human feedback.  
 754 In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in  
 755 Neural Information Processing Systems*, volume 35, pp. 27730–27744. Curran Associates, Inc.,  
 2022. URL [https://proceedings.neurips.cc/paper\\_files/paper/2022/  
 hash/b1efde53be364a73914f58805a001731-Abstract-Conference.html](https://proceedings.neurips.cc/paper_files/paper/2022/hash/b1efde53be364a73914f58805a001731-Abstract-Conference.html).

- 756 Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic  
 757 evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association*  
 758 *for Computational Linguistics, July 6-12, 2002, Philadelphia, PA, USA*, pp. 311–318. ACL, 2002.  
 759 doi: 10.3115/1073083.1073135. URL <https://aclanthology.org/P02-1040/>.
- 760 Qiwei Peng, Yekun Chai, and Xuhong Li. Humaneval-xl: A multilingual code generation benchmark  
 761 for cross-lingual natural language generalization. *arXiv preprint arXiv:2402.16694*, 2024. URL  
 762 <https://arxiv.org/abs/2402.16694>.
- 763 Julian Aron Prenner and Romain Robbes. Runbugrun – an executable dataset for automated program  
 764 repair. *arXiv preprint arXiv:2304.01102*, 2023. URL <https://arxiv.org/abs/2304.01102>.
- 765 Laria Reynolds and Kyle McDonell. Prompt programming for large language models: Beyond the  
 766 few-shot paradigm. In *CHI '21: CHI Conference on Human Factors in Computing Systems, Virtual*  
 767 *Event / Yokohama Japan, May 8-13, 2021, Extended Abstracts*, pp. 314:1–314:7. ACM, 2021. doi:  
 768 10.1145/3411763.3451760. URL <https://doi.org/10.1145/3411763.3451760>.
- 769 Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi  
 770 Adi, Jingyu Liu, Tal Remez, Jérémie Rapin, et al. Code llama: Open foundation models for code.  
 771 *arXiv preprint arXiv:2308.12950*, 2023. URL <https://arxiv.org/abs/2308.12950>.
- 772 Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman  
 773 Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. Bloom: A 176b-  
 774 parameter open-access multilingual language model. *arXiv preprint arXiv:2211.05100*, 2022. URL  
 775 <https://arxiv.org/abs/2211.05100>.
- 776 Runchu Tian, Yining Ye, Yujia Qin, Xin Cong, Yankai Lin, Zhiyuan Liu, and Maosong Sun.  
 777 Debugbench: Evaluating debugging capability of large language models. *arXiv preprint*  
 778 *arXiv:2401.04621*, 2024. URL <https://arxiv.org/abs/2401.04621>.
- 779 Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of Machine*  
 780 *Learning Research*, 9(11), 2008. URL <https://www.jmlr.org/papers/volume9/vandermaaten08a/vandermaaten08a.pdf>.
- 781 Yue Wang, Weishi Wang, Shafiq Joty, and Steven CH Hoi. Codet5: Identifier-aware unified pre-trained  
 782 encoder-decoder models for code understanding and generation. *arXiv preprint arXiv:2109.00859*,  
 783 2021. URL <https://arxiv.org/abs/2109.00859>.
- 784 Zhiruo Wang, Shuyan Zhou, Daniel Fried, and Graham Neubig. Execution-based evaluation for  
 785 open-domain code generation. In *Findings of the Association for Computational Linguistics: EMNLP* 2023,  
 786 pp. 1271–1290, 2023. URL <https://arxiv.org/abs/2212.10481>.
- 787 Yuxiang Wei, Zhe Wang, Jiawei Liu, Yifeng Ding, and Lingming Zhang. Magicoder: Source code is  
 788 all you need. *arXiv preprint arXiv:2312.02120*, abs/2312.02120, 2023. doi: 10.48550/ARXIV.  
 789 2312.02120. URL <https://arxiv.org/abs/2312.02120>.
- 790 Rui Xie, Zhengran Zeng, Zhuohao Yu, Chang Gao, Shikun Zhang, and Wei Ye. Codeshell technical  
 791 report. *arXiv preprint arXiv:2403.15747*, 2024. URL <https://arxiv.org/abs/2403.15747>.
- 792 Weixiang Yan, Yuchen Tian, Yunzhe Li, Qian Chen, and Wen Wang. CodeTransOcean: A  
 793 comprehensive multilingual benchmark for code translation. In Houda Bouamor, Juan Pino,  
 794 and Kalika Bali (eds.), *Findings of the Association for Computational Linguistics: EMNLP* 2023,  
 795 pp. 5067–5089, Singapore, December 2023. Association for Computational Linguistics.  
 796 doi: 10.18653/v1/2023.findings-emnlp.337. URL <https://aclanthology.org/2023.findings-emnlp.337>.
- 797 Pengcheng Yin, Wen-Ding Li, Kefan Xiao, Abhishek Rao, Yeming Wen, Kensen Shi, Joshua  
 798 Howland, Paige Bailey, Michele Catasta, Henryk Michalewski, Oleksandr Polozov, and Charles  
 799 Sutton. Natural language to code generation in interactive data science notebooks. In Anna  
 800 Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual Meeting*  
 801 *of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 126–173, Toronto,

- 810 Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.9.  
 811 URL <https://aclanthology.org/2023.acl-long.9>.
- 812
- 813 Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng  
 814 Zhu, Jianqun Chen, Jing Chang, et al. Yi: Open foundation models by 01.ai. *arXiv preprint*  
 815 *arXiv:2403.04652*, 2024. URL <https://arxiv.org/abs/2403.04652>.
- 816 Hao Yu, Bo Shen, Dezhong Ran, Jiaxin Zhang, Qi Zhang, Yuchi Ma, Guangtai Liang, Ying Li, Qianxi-  
 817 ang Wang, and Tao Xie. Codereval: A benchmark of pragmatic code generation with generative  
 818 pre-trained models. In *Proceedings of the 46th IEEE/ACM International Conference on Soft-  
 819 ware Engineering*, pp. 1–12, 2024. URL <https://dl.acm.org/doi/abs/10.1145/3597503.3623316>.
- 820
- 821 Daoguang Zan, Ailun Yu, Bo Shen, Jiaxin Zhang, Taihong Chen, Bing Geng, Bei Chen, Jichuan  
 822 Ji, Yafen Yao, Yongji Wang, and Qianxiang Wang. Can programming languages boost each  
 823 other via instruction tuning? *arXiv preprint arXiv:2308.16824*, abs/2308.16824, 2023. URL  
 824 <https://arxiv.org/abs/2308.16824>.
- 825
- 826 Quanjun Zhang, Tongke Zhang, Juan Zhai, Chunrong Fang, Bowen Yu, Weisong Sun, and Zhenyu  
 827 Chen. A critical review of large language model on software engineering: An example from  
 828 chatgpt and automated program repair. *arXiv preprint arXiv:2310.08879*, 2023. URL <https://arxiv.org/abs/2310.08879>.
- 829
- 830 Qinkai Zheng, Xiao Xia, Xu Zou, Yuxiao Dong, Shan Wang, Yufei Xue, Zihan Wang, Lei Shen,  
 831 Andi Wang, Yang Li, Teng Su, Zhilin Yang, and Jie Tang. Codegeex: A pre-trained model for  
 832 code generation with multilingual evaluations on humaneval-x. *arXiv preprint arXiv:2303.17568*,  
 833 abs/2303.17568, 2023a. doi: 10.48550/ARXIV.2303.17568. URL <https://doi.org/10.48550/arXiv.2303.17568>.
- 834
- 835 Qinkai Zheng, Xiao Xia, Xu Zou, Yuxiao Dong, Shan Wang, Yufei Xue, Zihan Wang, Lei Shen,  
 836 Andi Wang, Yang Li, et al. Codegeex: A pre-trained model for code generation with multilingual  
 837 evaluations on humaneval-x. *arXiv preprint arXiv:2303.17568*, 2023b. URL <https://arxiv.org/abs/2303.17568>.
- 838
- 839 Tianyu Zheng, Ge Zhang, Tianhao Shen, Xueling Liu, Bill Yuchen Lin, Jie Fu, Wenhui Chen, and  
 840 Xiang Yue. Opencodeinterpreter: Integrating code generation with execution and refinement. *arXiv*  
 841 *preprint arXiv:2402.14658*, 2024. URL <https://arxiv.org/abs/2402.14658>.
- 842
- 843 Terry Yue Zhuo, Minh Chien Vu, Jenny Chim, Han Hu, Wenhao Yu, Ratnadira Widayasari, Imam  
 844 Nur Bani Yusuf, Haolan Zhan, Junda He, Indraneil Paul, et al. Bigcodebench: Benchmarking code  
 845 generation with diverse function calls and complex instructions. *arXiv preprint arXiv:2406.15877*,  
 846 2024.
- 847
- 848
- 849
- 850
- 851
- 852
- 853
- 854
- 855
- 856
- 857
- 858
- 859
- 860
- 861
- 862
- 863

864 A APPENDIX  
865866 A.1 LIMITATIONS  
867

868 **Cover more languages and tasks** Our work focuses on evaluating models on multilingual  
869 programming tasks, currently supporting assessments in 40 programming languages. Many models  
870 claim to support over 80 programming languages, so our work can continue to expand the range of  
871 programming languages, the number of test cases, and types of tasks to provide a more comprehensive  
872 evaluation of the models.

873 **Research on cross-lingual transfer** We do not delve deeply into the cross-lingual transfer capabilities  
874 of models, leaving room for exploring a wider variety of models and models of different sizes. In  
875 the future, we will explore the programming language transferability scaling law of the code LLMs  
876 with different sizes.

878 A.2 DATA ANNOTATION  
879880 A.2.1 HUMAN ANNOTATION  
881

882 To create the massively multilingual code evaluation benchmark MCEVAL, the annotation of multi-  
883 lingual code samples is conducted utilizing a comprehensive and systematic human annotation  
884 procedure, underpinned by rigorously defined guidelines to ensure accuracy and consistency.

885 We recruited annotators with backgrounds in computer software development from universities.  
886 During recruitment, we evaluated their programming and system operation skills to ensure they could  
887 handle tasks like question writing, code static analysis, and unit testing. Each annotator was assigned  
888 tasks based on their proficiency in specific programming languages.

889 Initially, 10 software developers in computer science are recruited as multilingual programming  
890 annotators with proven proficiency in the respective programming languages.

892 Following a detailed training session on the annotation protocol, annotators are tasked with creating  
893 problem definitions and the corresponding solution.

894 The guidelines for our annotation training session primarily cover the following aspects:  
895

- 896 • **Standardized Format:** We provide an annotation example for 40 programming languages.  
897 Annotators are required to adhere to this standardized format when annotating data.
- 898 • **Accessibility:** The reference data for our annotations is sourced from materials available  
899 under permissive licenses, allowing unrestricted use and redistribution for research purposes.
- 900 • **Difficulty Level:** We provide annotators with detailed guidelines on the difficulty clas-  
901 sification for each language. Annotators must strictly follow these guidelines to label  
902 problems according to their respective difficulty levels based on algorithmic complexity and  
903 functionality.
- 904 • **Self-Contained:** Annotators are required to thoroughly review their annotated problems  
905 to ensure that the problem descriptions include all necessary information for solving them  
906 without ambiguity. The provided example inputs and outputs must be correct, the reference  
907 answers must execute correctly, and the test cases written should comprehensively evaluate  
908 the accuracy of the functions.

910 A.2.2 REFERENCE DATA  
911

912 We use reference data (The annotators only draw the inspiration from the reference website, and  
913 create the question manually.) to draft questions and solutions with the help of GPT-4-Turbo. The  
914 draft questions initially may contain numerous issues, such as self-contained errors, inconsistent  
915 difficulty levels, overly simple unit tests, and incorrect unit tests. To address these issues, annotators  
916 revised and tested the draft questions to ensure that all questions and code were accurate and that the  
917 unit tests would pass. These reference data came from the following websites: For the algorithmic  
types of questions, we refer to the following websites:

- 918       • <https://www.dotcpp.com/>  
 919       • <https://www.luogu.com.cn>  
 920       • <https://www.codecademy.com/>  
 921       • <https://www.codewars.com/>

923       For markup language and design-type questions, we refer to the following websites:  
 924

- 925       • <https://www.runoob.com/>  
 926       • <https://www.w3schools.com/>

### 928       A.2.3 QUALITY CONTROL

930       We adopt a dual-pass system to ensure the quality of our benchmark MCEVAL. First, one annotator  
 931       labels the code snippets and their corresponding unit tests. Then, another annotator independently  
 932       reviews these annotations, verifying their correctness, code accuracy, and unit test integrity. Following  
 933       this, three senior annotators evaluate the overall unit test pass rate of the annotated dataset. If the  
 934       pass rate exceeds 90%, the senior annotators proceed to perform additional reviews and corrections.  
 935       Otherwise, the data is returned to the annotators for refinement. Finally, the canonical solution  
 936       (labeled by the annotator) passes all corresponding test cases to ensure the correctness of each created  
 937       problem (100% pass rate). This rigorous process ensures the creation of a high-quality, multilingual  
 938       programming benchmark that supports in-depth analysis and understanding of code examples across  
 939       diverse programming languages.

### 940       A.2.4 ANNOTATION COSTS

942       We paid all the annotators the equivalent of \$6 per question and provided them with a comfortable  
 943       working environment, free meals, and souvenirs. We also provided the computer equipment and  
 944       GPT-4 interface required for labeling. We labeled about 2,000 questions in total and employed them  
 945       to check the quality of the questions/answers, and the total cost was about \$12,000 in US dollars. The  
 946       annotators checked the derived tasks, including multilingual code explanation and code completion.

## 947       A.3 DETAILED ABOUT MCEVAL-INSTRUCT

949       Here we describe in detail the specific methods and processes of code sampling and quality control in  
 950       the process of constructing MCEVAL-INSTRUCT.

### 952       Code Sampling

- 954       • (1) We crawled a large amount of code data from GitHub and preprocessed the data according  
 955       to the StarCoder processing flow. Then, we can get a high-quality code dataset.  
 956       • (2) For each programming language, we randomly sampled the data to build an original  
 957       dataset containing 40 languages, where the instances of each language are sampled to the  
 958       same number of samples.  
 959       • (3) In addition, in the refined stage, we further used GPT-4 to refine the code snippets to  
 960       the clearer and more standardized code (the code with docstring, code comments, and clear  
 961       variables) and verified the refined code snippet through static code analysis and unit testing  
 962       methods to ensure the correctness and quality of the code.

963       **Quality Control** To check the correctness of the generated code from the GPT-4, we use static  
 964       analysis and unit testing to ensure the accuracy of generated data as much as possible:

- 966       • Code Correctness: We use static code analysis tools and Abstract Syntax Tree (AST) parsing  
 967       to ensure the correctness of the code.  
 968       • Automated Testing: We leverage an LLM to automatically generate unit tests, filtering out  
 969       code that fails to pass these tests.  
 970       • Harmful Content Filtering: We filter out code snippets containing harmful information using  
 971       keyword-based filters.

972 A.4 EXAMPLES IN MCEVAL  
973974 Figure 10 display three examples of multilingual generation.  
975976 In Figure 11, we show three examples of multilingual explanation.  
977978 In Figure 12, we display three examples of multilingual explanation. The three examples from left  
979 to right correspond to the span completion task, the single-line completion task, and the multi-line  
980 completion task.  
981

**I Instruction: Language Lisp**

```
Write a Common Lisp function `defun create-largest-number (numbers)` to solve the following problem:  
This function takes a list of positive integers and rearranges them to form the largest possible number.  
It does this by concatenating the integers in an order such that their digit-wise combination yields the maximum value.  
Example:  
>>> create-largest-number '(56 9 45)  
>>> create-largest-number '(5 50 56)  
>>> create-largest-number '(565589)
```

**Reference Solution**

```
(defun create-largest-number (numbers)
  (let ((ctr-numbers (mapcar #'number-to-string numbers)))
    (setf str-numbers (sort str-numbers #'larger-when-concatenated))
    (reduce #'lambda (acc x) (concatenate 'string acc x) str-numbers :initial-value ""))
  (defun number-to-string (number)
    (write-to-string number)
    (writer larger-when-concatenated a b)
    (string (concatenate 'string a b) (concatenate 'string b a))))
```

**Test Cases**

```
(defun test-create-largest-number ()
  (assert (string= (create-largest-number '(56 9 45)) "9565589")
          (assert (string= (create-largest-number '(5 50 56)) "5056565589")
          (assert (string= (create-largest-number '(3 34 382 50)) "50343823"))
          (assert (string= (create-largest-number '(10 2 23)) "2310"))
          (assert (string= (create-largest-number '(4 42 40 480)) "442440480")))))
```

**I Instruction: Language AWK**

Using the awk command in Linux, complete the following task: In data/awk/contribution.txt, where the data format is as follows:

```
# "Zhang Dandan 41117397 :250:100:175
# Zhang Liayou 30032015 :250:100:201
# Zhang Peng 80042301 :250:100:150
# Wu Weisai 70271111 :250:80:75
# Liu Bingbing 41117483 :250:100:175
# Wang Xiaoxi 3515064655 :250:95:135
# Zi Gege 1986787358 :250:168:200
# Lai Youliu 918391615 :175:75:300
# Lao Nama 918391655 :250:100:175'
```

**Reference Solution**

```
awk -F "[ :]" '$3~/^(41)/{print $1,$2,$3}' data/awk/contribution.txt
```

**Test Cases**

```
awk_command = < Code Being Evaluated >
ref_command = " awk -F '[ :]' '$3~/^(41)/{print $1,$2,$3}' data/awk/contribution.txt "
generate_result = subprocess.check_output(
    ref_command, shell=True, text=True)
ref_result = subprocess.check_output(
    awk_command, shell=True, text=True)
assert generate_result == reference_result
```

**I Instruction: Language R**

Write a R function `longest\_increasing\_subsequence <- function(sequence)` to solve the following problem: This function longest\_increasing\_subsequence takes a vector sequence representing a sequence of integers and returns the length of the longest increasing subsequence within it. An increasing subsequence is defined as a set of indices in the sequence that are in increasing order and taken from the original sequence without changing their order. The length of the subsequence is the number of elements it contains.

**Reference Solution**

```
longest_increasing_subsequence <- function(sequence) {
  n <- length(sequence)
  lis <- numeric(n)
  lis[1] <- 1
  for (i in 2:n) {
    max_val <- 0
    for (j in 1:(i-1)) {
      if (sequence[i] > sequence[j] & lis[j] > max_val)
        max_val <- lis[j]
    }
    lis[i] <- max_val + 1
  }
  max(lis)
}
```

**Test Cases**

```
main <- function() {
  stopifnot(longest_increasing_subsequence(c(1, 7, 3, 5, 9, 4, 8)) == 4)
  stopifnot(longest_increasing_subsequence(c(10, 22, 9, 33, 21, 50, 41, 60)) == 5)
}
main()
```

998 Figure 10: Examples of multilingual generation. The data mainly consists of an instruction part  
999 (including function name, function description, and function call cases), a reference solution, and  
1000 a test cases part. **Left.** Shows an example of the Lisp language. **Middle.** Shows a file processing  
1001 programming task in AWK language. During the evaluation, the corresponding file processing result  
1002 by the generated code will be compared with the reference answer. **Right.** Shows an example of the  
1003 R language.

**I Instruction: Language Kotlin**

```
fun findPrimePairs(maxNumber: Int): List<Pair<Int, Int>> {
    fun isPrime(num: Int): Boolean {
        if (num <= 1) return false
        for (i in 2 until num) {
            if (num % i == 0) return false
        }
        return true
    }

    val pairs = mutableListOf<Pair<Int, Int>>()
    for (i in 2..maxNumber - 2) {
        if (isPrime(i) && isPrime(i + 2)) {
            pairs.add(Pair(i, i + 2))
        }
    }
    return pairs
}
```

Provide a concise natural language description (docstring) of the Kotlin code in English using at most 500 characters.

**Reference Explanation**

Finds all prime pairs where each prime is less than or equal to a given number and the pair differs by 2. A prime pair is defined as two prime numbers where the difference between them is exactly 2.

Example:  
>>> findPrimePairs(10)  
[(3, 5), (5, 7)]  
>>> findPrimePairs(100)  
[(3, 5), (5, 7), (11, 13), (17, 19), (29, 31), (41, 43), (59, 61), (71, 73)]

**I Instruction: Language Lua**

```
function addDigits(num)
    while num >= 10 do
        local sum = 0
        while num > 0 do
            sum = sum + (num % 10)
            num = math.floor(num / 10)
        end
        num = sum
    end
    return num
end
```

Provide a concise natural language description (docstring) of the Lua code in English using at most 500 characters.

**Reference Explanation**

Given a non-negative integer num, repeatedly add all its digits until the result has only one digit.

For example:  
>>> addDigits(38)  
2  
Because 3 + 8 = 11, and 1 + 1 = 2. Since 2 has only one digit, 2 is the result.

**I Instruction: Language HTML**

```
<table>
  <tr>
    <th><strong>col_1</strong></th>
    <th><strong>col_2</strong></th>
  </tr>
  <tr>
    <td>row_1; col_1</td>
    <td>row_1; col_2</td>
  </tr>
  <tr>
    <td>row_2; col_1</td>
    <td>row_2; col_2</td>
  </tr>
</table>
```

Provide a concise natural language description (docstring) of the HTML code in English using at most 500 characters.

**Reference Explanation**

a table with two rows and two columns, with column names col\_1 and col\_2. The table content includes its corresponding row and column numbers, such as row\_1; col\_1, and bold the header text

1022 Figure 11: Examples of multilingual explanation. The data mainly consists of an instruction part  
1023 (including a complete function), a reference Explanation. **Left.** Shows an example of the Kotlin  
1024 language. **Middle.** Shows an example of the Lua language. **Right.** Shows an example of the HTML  
1025 language.

**I Instruction: Language C++**

```
Below is a explanation of CPP code and incomplete code implementation.

* DocString:
You are given three integers a, b, c, where two of them are equal, and the third is different from the other two. Your task is to find the value that occurs exactly once.

Examples:
extraNumber(0, 0, 1) returns 1
extraNumber(4, 3, 4) returns 3

* Incomplete Code:
int extraNumber(int a, int b, int c){
    [MASK]a == c) {
        return b;
    }
    else {
        return a;
    }
}

Please fill the [MASK] (multiple lines of code may be masked out) and write the complete function.
```

**Reference Solution**

```
int extraNumber(int a, int b, int c){
    if (a == b) {
        return c;
    }
    else if (a == c) {
        return b;
    }
    else {
        return a;
    }
}
```

**Test Cases**

```
#include<assert.h>
#include<vector>
using namespace std;

int main() {
    assert(extraNumber(2, 7, 2) == 7);
    assert(extraNumber(3, 2, 2) == 3);
    assert(extraNumber(5, 5, 1) == 1);
    assert(extraNumber(500000000, 3, 500000000) == 3);
    assert(extraNumber(500000000, 500000000, 3) == 3);
    return 0;
}
```

**I Instruction: Language Rust**

```
Below is a explanation of Rust code and incomplete code implementation.

* DocString:
You are given a list of deposit and withdrawal operations on a bank account that starts with zero balance. Your task is to detect if at any point the balance of account falls below zero, and at that point function should return True. Otherwise it should return False.

* Incomplete Code:
fn below_zero(operations:Vec<i32>) -> bool{
    let mut balance:i32 = 0;
    for op in operations {
        [MASK]
        if balance < 0 {
            return true;
        }
    }
    return false;
}

Please fill the [MASK] (multiple lines of code may be masked out) and write the complete function.
```

**Reference Solution**

```
fn below_zero(operations:Vec<i32>) -> bool{
    let mut balance:i32 = 0;
    for op in operations {
        balance = balance + op;
        if balance < 0 {
            return true;
        }
    }
    return false;
}
```

**Test Cases**

```
#![test]
mod tests {
    use super::*;

    #[test]
    fn test_below_zero() {
        assert_eq!(below_zero(vec![]), false);
        assert_eq!(below_zero(vec![1, 2, -3, 1, 2, -3]), false);
        assert_eq!(below_zero(vec![1, 2, -4, 5, 6]), true);
        assert_eq!(below_zero(vec![1, -1, 2, -2, 5, -5, 4, -4]), false);
        assert_eq!(below_zero(vec![1, -2, 2, -2, 5, -5, 4, -4]), true);
        assert_eq!(below_zero(vec![1, -2, 2, -2, 5, -5, 4, -4]), true);
    }
}
```

**I Instruction: Language Shell**

```
Below is a explanation of Shell code and incomplete code implementation.

* Docstring:
This function calculates the maximum number of pens that can be bought with a given amount of money. The price of one pen is 1 Yuan and 9 Jiao (1.9 Yuan). The function takes two integers, a and b, as input where 'a' represents the Yuan and 'b' represents the Jiao part of the total money available. It returns the maximum number of pens that can be purchased. For example, if a=5 and b=0, the function will return 2, as the total money is 5 Yuan, and two pens cost 3.8 Yuan.
```

**Reference Solution**

```
calculate_max_pens() {
    [MASK]
    [MASK]
    local total_jiao=$((yuan * 10 + jiao))
    [MASK]
    local max_pens=$((total_jiao / price_per_pen))
    echo "$max_pens"
}
```

**Test Cases**

```
test_calculate_max_pens() {
    local result
    result=$(calculate_max_pens 5 5)
    [[ "$result" -eq 2 ]] || { echo "Test 1 failed: Expected 2, got $result"; exit 1; }
    result=$(calculate_max_pens 28)
    [[ "$result" -eq 14 ]] || { echo "Test 2 failed: Expected 14, got $result"; exit 1; }
    result=$(calculate_max_pens 38)
    [[ "$result" -eq 20 ]] || { echo "Test 3 failed: Expected 20, got $result"; exit 1; }
    result=$(calculate_max_pens 11 0)
    [[ "$result" -eq 5 ]] || { echo "Test 4 failed: Expected 5, got $result"; exit 1; }
}
```

Figure 12: Examples of multilingual completion. The data mainly consists of an instruction part (including a incomplete function ), a reference complete code solution and test cases. **Left**. Shows an span completion example of the C++ language. **Middle**. Shows an single line completion example of the Rust language. **Right**. Shows an multiple line completion example of the Shell language.

**I Instruction Language Json**

```
Create a JSON represents a collection of websites, each defined with two key details: their name and URL. The sites array holds two objects. The first object represents the website "Google" with the URL "www.google.com". The second object describes the website "Weibo" with its URL "www.weibo.com". This structure efficiently organizes these popular websites by their most essential identifiers.
```

**Reference Solution**

```
{
    "sites": [
        {
            "name": "Google",
            "url": "www.google.com"
        },
        {
            "name": "Weibo",
            "url": "www.weibo.com"
        }
    ]
}
```

**Test**

```
{
    "sites": [
        {
            "url": "www.google.com",
            "name": "Google"
        },
        {
            "url": "www.weibo.com",
            "name": "Weibo"
        }
    ]
}
```

All subcomponents match exactly, test passes !

Figure 13: Examples of Markup language (Json) generation task evaluation.

## A.5 EVALUATION

For programming languages other than markup languages, we use an execution-based correctness metric by running the code with the provided test cases. For markup languages, we use the Exact

Table 6: Runtime environments for different programming languages.

| Language     | Runtime Environments   |
|--------------|--|
| AWK          | GNU bash, version 4.4.20(1)-release (x86_64-pc-linux-gnu)                          |
| C            | gcc (Ubuntu 7.5.0-3ubuntu118.04) 7.5.0   |
| C#           | dotnet 8.0.100   |
| CPP          | g++ (Ubuntu 7.5.0-3ubuntu118.04) 7.5.0   |
| CoffeeScript | CoffeeScript version 1.12.7  |
| Common Lisp  | SBCL 1.4.5.debian  |
| Dart         | Dart SDK version: 3.3.1 (stable)   |
| Elixir       | elixir 1.3.3   |
| Emacs Lisp   | GNU Emacs 25.2.2   |
| Erlang       | Erlang/OTP 20 [erts-9.2]   |
| F#           | dotnet 8.0.100   |
| Fortran      | GNU Fortran (Ubuntu 7.5.0-3ubuntu118.04) 7.5.0                                     |
| Go           | go version go1.18.4 linux/amd64  |
| Groovy       | Groovy Version: 4.0.16 JVM: 17.0.9 Vendor: Oracle Corporation OS: Linux            |
| HTML         | -  |
| Haskell      | The Glorious Glasgow Haskell Compilation System, version 9.4.7                     |
| Json         | -  |
| Java         | javac 11.0.19  |
| JavaScript   | Node.js v16.14.0   |
| Julia        | julia v1.9.4   |
| Kotlin       | kotlinc-jvm 1.9.21 (JRE 17.0.9+11-LTS-201)   |
| Lua          | Lua 5.4.6 Copyright (C) 1994-2023 Lua.org, PUC-Rio                                 |
| Markdown     | -  |
| PHP          | PHP 7.2.24-0ubuntu0.18.04.17 (cli) (built: Feb 23 2023 13:29:25) ( NTS )           |
| Pascal       | Free Pascal Compiler version 3.2.2 [2021/05/16] for x86_64                         |
| Perl         | perl 5, version 26, subversion 1 (v5.26.1) built for x86_64-linux-gnu-thread-multi |
| PowerShell   | PowerShell 7.4.0   |
| Python       | Python 3.8.12  |
| R            | R version 3.4.4  |
| Racket       | Racket v6.11   |
| Ruby         | ruby 2.5.1p57 (2018-03-29 revision 63029) [x86_64-linux-gnu]                       |
| Rust         | rustc 1.74.0 (79e9716c9 2023-11-13)  |
| Scala        | Scala code runner version 3.3.1 – Copyright 2002-2023, LAMP/EPFL                   |
| Scheme       | Racket v6.11   |
| Shell        | GNU bash, version 4.4.20(1)-release (x86_64-pc-linux-gnu)                          |
| Swift        | Swift version 5.9.2 (swift-5.9.2-RELEASE)  |
| Tcl          | tclsh 8.6.11   |
| TypeScript   | tsc Version 5.3.3  |
| VimScript    | VIM - Vi IMproved 9.0 (2022 Jun 28, compiled Dec 20 2023 18:57:50)                 |
| Visual Basic | dotnet 8.0.100   |

Match metric for evaluation. Taking Json as an example, we parse all subcomponents in Json. If the model result is exactly the same as the subcomponent of the reference solution, the model generation result is considered correct. An example of Markup language (Json) is shown in Figure 13.

We adopt the greedy Pass@1 (%) metric (Kulal et al., 2019; Chen et al., 2021) for our evaluations. For closed-source models, we generate answers through the official API service. For open-source models, we prioritize using vLLM (Kwon et al., 2023) for faster inference if the model is supported by vLLM. Otherwise, we perform inference with the Distributed Data Parallel (DDP) module from PyTorch. For the code generation and code completion tasks, we extract the functional part of the code from the model outputs and combine it with corresponding test cases to form compilable and executable code. For the code explanation task, we adopt a two-pass generation approach (Code-to-Natural-Language and Natural-Language-to-Code). The extraction and execution process for this task is consistent with the previous two tasks. We conduct all evaluations in a Docker environment. Detailed information on the code compilation and execution environment are displayed in Table 6. We have uploaded the Docker image to docker hub to facilitate the reproduction of results and the evaluation of new models.

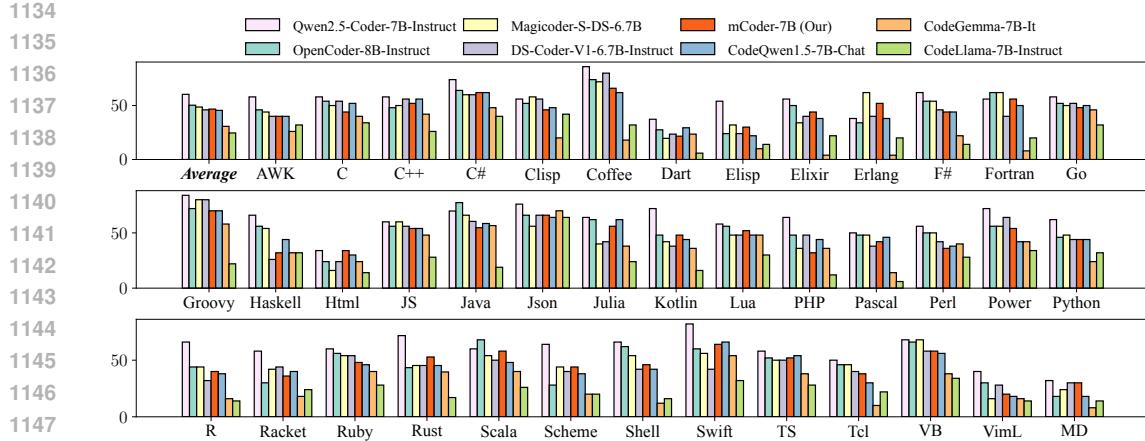


Figure 14: Pass@1 (%) scores of different code LLMs (<10B) for multilingual code generation tasks on MCEVAL. “AVG” represents the average scores of all code languages.

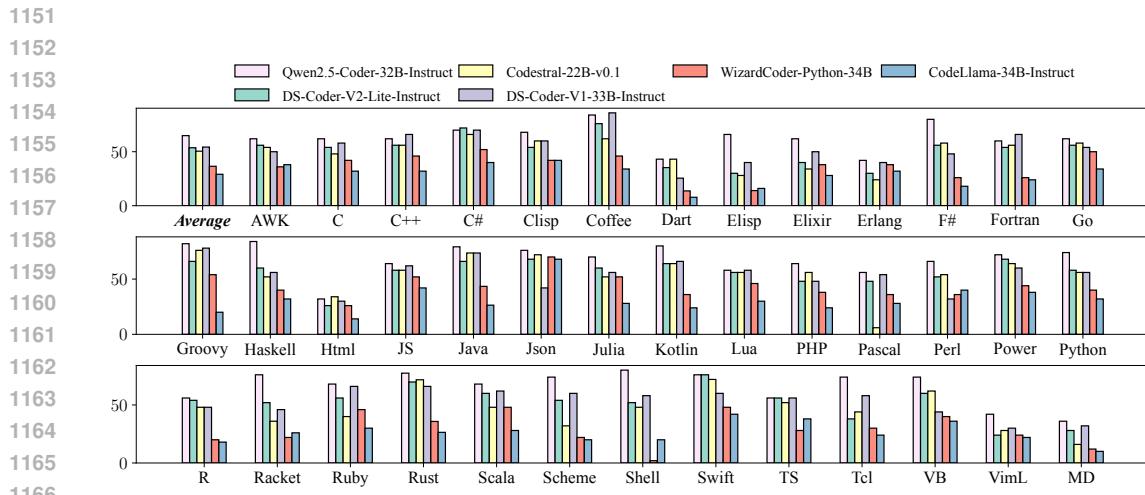


Figure 15: Pass@1 (%) scores of different code LLMs (10B to 40B) for multilingual code generation tasks on MCEVAL. “AVG” represents the average scores of all code languages.

## A.6 OPTIMIZATION DETAILS

All MCODER models are fine-tuned using 8 NVIDIA A800-80GB GPUs. The models are trained for 2 epochs with a cosine scheduler, starting at a learning rate of 2e-5 and incorporating a 3% warmup phase. Training a model takes about 5 hours. We used AdamW (Loshchilov & Hutter, 2017) as the optimizer and a batch size of 512 with a sequence truncation length of 4096. We use PyTorch’s Fully Sharded Data Parallel (FSDP) to perform distributed training of the model, and use gradient checkpointing technology and gradient accumulation to save memory and achieve training with a larger batch size.

## A.7 EXTRA RESULTS

## A.8 PROGRAMMING CLASSIFICATION

As shown in Table 7 and Table 8, we comprehensively display the code generation performance of the models we tested across various programming paradigms and application scenarios.

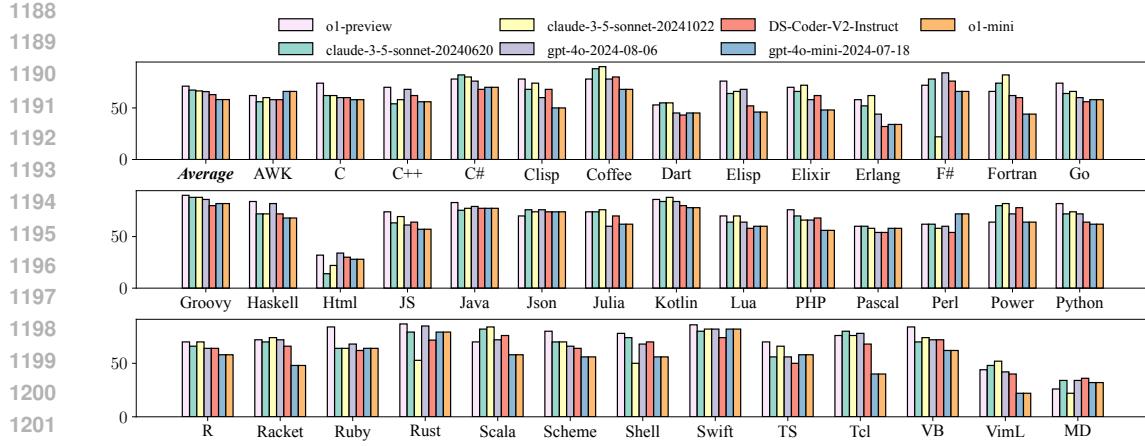


Figure 16: Pass@1 (%) scores of different code LLMs (Closed Source & 200B+) for multilingual code generation tasks on MCEVAL. “AVG” represents the average scores of all code languages.

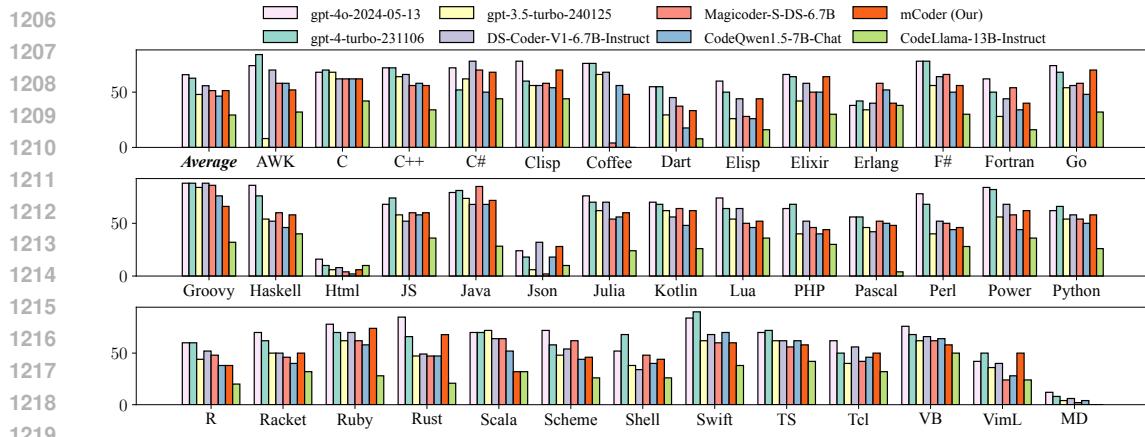


Figure 17: Pass@1 (%) scores of different code LLMs for multilingual code explain tasks on MCEVAL. “AVG” represents the average scores of all code languages.

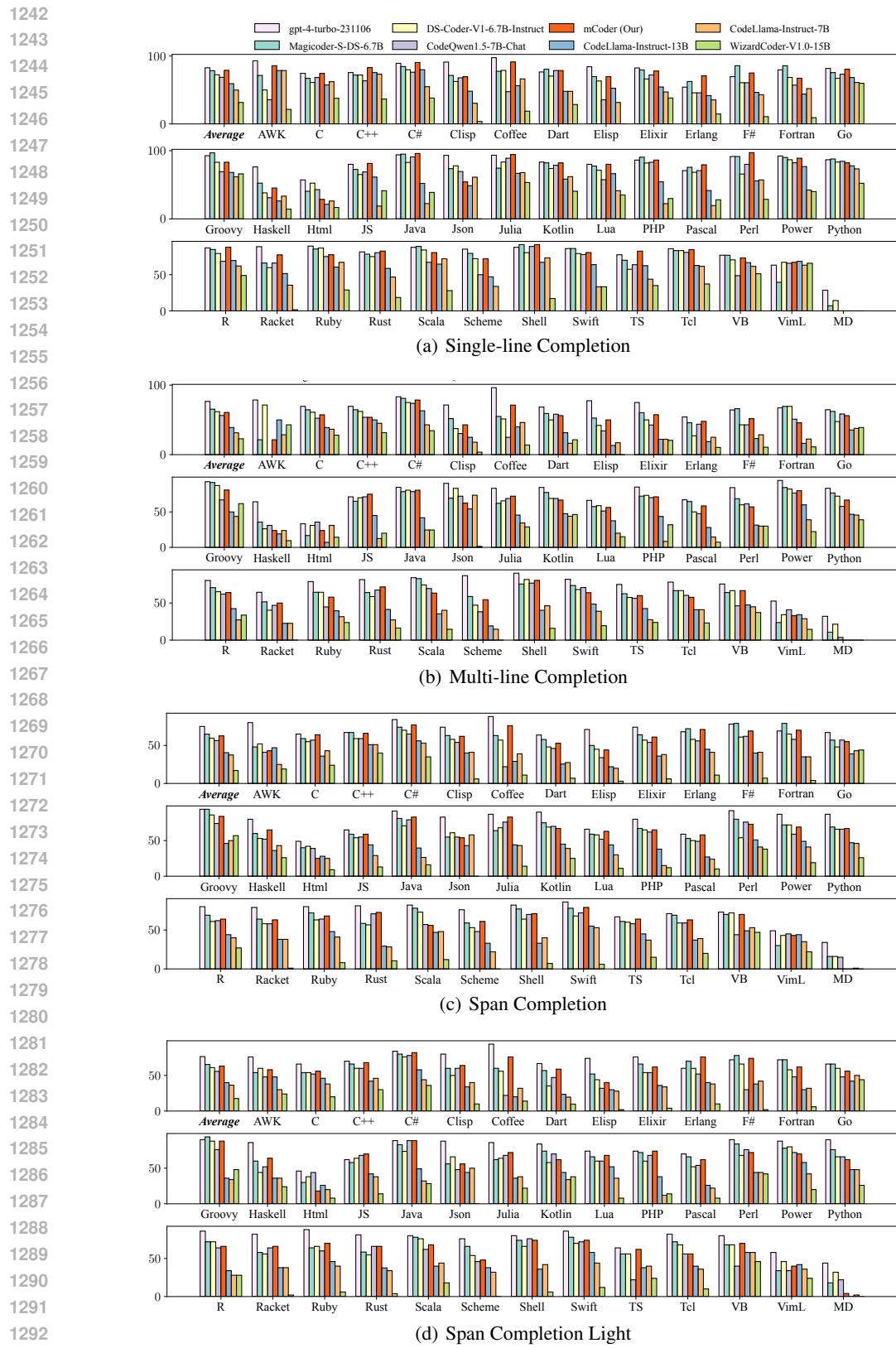
### A.9 MCODER RESULT

In Table 9, we show some extra MCODER Pass@1 (%) results on multilingual code generation tasks. We evaluate the base models CodeQwen-1.5 and DeepsSeek-Coder-1.5 respectively. In addition to CodeQwen-1.5, we also selected DeepSeek-Coder-1.5-base as the base model for fine-tuning.

### A.10 PARALLEL QUESTIONS ACROSS LANGUAGES & PROGRAMMING GRAMMAR

Due to the large number of languages, it is difficult to ensure parallel problem annotation. For most language annotations, we follow the characteristics of the language and perform independent annotations. For example, structured languages such as Markdown and HTML need independent annotations. For some similar languages, such as Typescript and Javascript, we use parallel annotation on some data.

As shown in Figure 19, we analyzed the programming languages in the MCEVAL from the representation perspective. We used CodeBERT (Feng et al., 2020) to extract code representations from code snippets in MCEVAL. These representations were visualized using t-SNE (Van der Maaten & Hinton, 2008) and hierarchical clustering (Murtagh & Contreras, 2012) methods. The figure clearly shows that languages with similar syntax have closely related representations. For example, other functional programming languages similar to Common Lisp, as well as C, C++, Java, and scripting languages, exhibit high grammar similarity.



1294  
1295  
Figure 18: Pass@1 (%) scores of different models for multilingual code completion tasks on MCEVAL.  
“Avg” represents the average scores of all code languages.

1296 Table 7: Pass@1(%) results of code generation performance of across various programming paradigms  
1297

| Method                        | Procedural  | Object Oriented | Multiple Paradigms | Functional  | Markup Language |
|-------------------------------|-------------|-----------------|--------------------|-------------|-----------------|
| GPT-4o (240517)               | <b>58.0</b> | <b>79.8</b>     | <b>65.9</b>        | <b>67.0</b> | 46.0            |
| GPT-4 Turbo (231106)          | 56.7        | 78.7            | 65.2               | 59.3        | <b>46.7</b>     |
| GPT-3.5-Turbo (240125)        | 38.7        | 66.8            | 57.6               | 44.3        | 39.3            |
| Codegemma-7b-it               | 19.3        | 46.6            | 34.0               | 16.3        | 34.0            |
| CodeLlama-13b-Instruct        | 21.3        | 32.0            | 27.0               | 32.3        | 28.0            |
| CodeLlama-34b-Instruct        | 27.3        | 33.6            | 28.0               | 30.0        | 30.7            |
| CodeLlama-7b                  | 20.3        | 28.1            | 23.4               | 26.7        | 30.7            |
| CodeQwen-1.5-7b-Chat          | 41.3        | 57.3            | 46.3               | 41.0        | 37.3            |
| Codeshell-7b-chat             | 16.0        | 24.1            | 25.7               | 14.0        | 34.7            |
| Codestral-22B-v0.1            | 40.0        | 67.6            | 54.1               | 39.7        | 40.7            |
| DeepSeekCoder-33b-instruct    | 52.7        | 62.8            | 56.3               | 52.0        | 34.7            |
| DeepSeekCoder-1.5-7b-instruct | 39.0        | 51.8            | 48.8               | 41.0        | 40.0            |
| Magicoder-S-DS-6.7B           | 45.7        | 58.5            | 49.4               | 49.0        | 32.0            |
| Llama-3-8B-Instruct           | 27.3        | 44.7            | 38.0               | 32.0        | 33.3            |
| Nxcode-CQ-7B-orpo             | 40.7        | 54.9            | 45.5               | 41.3        | 36.7            |
| OCTOCODER                     | 20.7        | 28.9            | 21.9               | 25.0        | 25.3            |
| OpenCodeInterpreter-DS-6.7B   | 40.7        | 57.7            | 46.4               | 42.0        | 42.0            |
| Phi-3-medium-4k-instruct      | 32.3        | 43.1            | 36.6               | 26.7        | 35.3            |
| Qwen1.5-72B-Chat              | 38.3        | 37.2            | 36.2               | 29.3        | 39.3            |
| WizardCoder-15B-V1.0          | 19.0        | 31.6            | 34.2               | 24.0        | 6.7             |
| WizardCoder-Python-34B        | 27.7        | 43.9            | 38.2               | 33.7        | 36.0            |
| mCODER                        | 41.3        | 57.3            | 47.4               | 42.3        | 43.3            |

1316 Table 8: Pass@1(%) results of code generation performance of across various application scenarios  
1317

| Method                        | Mobile      | Cross       | Desktop     | Frontend    | Backend     | Scientific  | General     | Content     | Education   | Scripts     | Editor      |
|-------------------------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| GPT-4o (230517)               | <b>84.0</b> | <b>68.3</b> | <b>75.0</b> | <b>66.7</b> | 64.6        | <b>71.6</b> | <b>57.6</b> | 46.0        | <b>72.7</b> | <b>65.7</b> | <b>52.0</b> |
| GPT-4 Turbo (231106)          | 81.0        | 64.4        | 74.0        | 64.0        | <b>66.6</b> | 66.8        | <b>57.6</b> | <b>46.7</b> | 60.7        | <b>65.7</b> | 50.0        |
| GPT-3.5 (240125)              | 60.0        | 56.7        | 71.0        | 63.3        | 57.5        | 55.6        | 50.4        | 39.3        | 45.3        | 50.0        | 25.0        |
| Codegemma-7b-it               | 45.0        | 40.4        | 43.0        | 34.7        | 37.7        | 21.6        | 24.8        | 34.0        | 22.0        | 29.7        | 13.0        |
| code-Llama-13b                | 30.0        | 15.4        | 39.0        | 34.7        | 28.0        | 23.2        | 34.8        | 28.0        | 24.0        | 27.7        | 13.0        |
| CodeLlama-34b-Instruct        | 33.0        | 17.3        | 38.0        | 38.0        | 27.2        | 24.0        | 32.8        | 30.7        | 26.7        | 31.7        | 19.0        |
| Code-Llama-7b-Instruct        | 24.0        | 12.5        | 37.0        | 29.3        | 22.7        | 20.8        | 29.2        | 30.7        | 19.3        | 27.0        | 14.0        |
| CodeQwen-1.5-7b               | 55.0        | 44.2        | 59.0        | 56.7        | 48.7        | 47.6        | 46.8        | 37.3        | 42.7        | 40.0        | 20.0        |
| Codeshell-7b-chat             | 23.0        | 14.4        | 26.0        | 40.7        | 26.1        | 17.2        | 21.2        | 34.7        | 13.3        | 22.7        | 8.0         |
| Codestral-22B-v0.1            | 68.0        | 58.7        | 64.0        | 57.3        | 55.0        | 54.0        | 44.8        | 40.7        | 30.0        | 53.3        | 28.0        |
| DeepSeekCoder-33b-instruct    | 63.0        | 50.0        | 57.0        | 68.0        | 60.6        | 54.8        | 54.0        | 34.7        | 56.7        | 52.7        | 35.0        |
| DeepSeekCoder-1.5-7b-instruct | 40.0        | 42.3        | 59.0        | 62.0        | 52.7        | 40.8        | 50.0        | 40.0        | 34.7        | 46.0        | 22.0        |
| Magicoder-S-DS-6.7B           | 49.0        | 43.3        | 64.0        | 60.7        | 50.4        | 49.6        | 52.4        | 32.0        | 48.7        | 49.7        | 24.0        |
| Llama-3-8B-Instruct           | 41.0        | 30.8        | 48.0        | 50.7        | 40.5        | 30.0        | 37.2        | 33.3        | 34.0        | 33.0        | 15.0        |
| Nxcode-CQ-7B-orpo             | 54.0        | 40.4        | 55.0        | 53.3        | 48.4        | 48.0        | 46.8        | 36.7        | 42.7        | 39.7        | 20.0        |
| OCTOCODER                     | 22.0        | 20.2        | 33.0        | 28.7        | 21.8        | 16.4        | 27.2        | 25.3        | 16.0        | 29.0        | 14.0        |
| OpenCodeInterpreter-DS-6.7B   | 47.0        | 42.3        | 64.0        | 58.0        | 45.9        | 47.6        | 46.4        | 42.0        | 43.3        | 44.0        | 24.0        |
| Phi-3-medium-4k-instruct      | 40.0        | 26.9        | 48.0        | 48.7        | 30.3        | 39.2        | 31.6        | 35.3        | 33.3        | 39.0        | 13.0        |
| Qwen1.5-72B-Chat              | 30.0        | 29.8        | 43.0        | 44.7        | 36.3        | 30.4        | 38.0        | 39.3        | 32.7        | 40.0        | 21.0        |
| WizardCoder-15B-V1.0          | 28.0        | 24.0        | 36.0        | 48.0        | 37.1        | 29.2        | 27.2        | 6.7         | 20.7        | 26.3        | 9.0         |
| WizardCoder-Python-34B        | 42.0        | 28.8        | 46.0        | 42.0        | 44.2        | 32.8        | 38.0        | 36.0        | 32.7        | 32.3        | 19.0        |
| mCODER                        | 56.0        | 38.5        | 60.0        | 57.3        | 50.4        | 48.0        | 46.0        | 43.3        | 39.3        | 44.3        | 25.0        |

## 1334 A.11 DETAILED RELATED WORK

1335 **Code Large Language Model.** In recent years, numerous large language models (LLMs) have  
1336 been developed specifically for code-related tasks. For the field of soft engineering, code LLMs (Feng  
1337 et al., 2020; Chen et al., 2021; Scao et al., 2022; Li et al., 2022; Allal et al., 2023; Fried et al.,  
1338 2022; Wang et al., 2021; Zheng et al., 2024; Guo et al., 2024) pre-trained on billions of code  
1339 snippets, such as StarCoder (Li et al., 2023; Lozhkov et al., 2024), CodeLlama (Rozière et al.,  
1340 2023), DeepSeekCoder (Guo et al., 2024), and Code-Qwen (Bai et al., 2023). The development and  
1341 refinement of code LLMs have been pivotal in automating software development tasks, providing  
1342 code suggestions, and supporting code generation/translation.

1343 To improve the performance of code generation, researchers used optimized prompts (Liu et al., 2023a;  
1344 Reynolds & McDonell, 2021; Zan et al., 2023; Beurer-Kellner et al., 2023), bring test cases (Chen  
1345 et al., 2023) and collaborative roles (Dong et al., 2023). There are also some related studies on using  
1346 large language models for other code tasks, such as dynamic programming (Dagan et al., 2023),  
1347 compiler optimization (Cummins et al., 2023), multi-lingual prompts (Di et al., 2023), and Program  
1348 of Thoughts (Chen et al., 2022).

Table 9: Additional MCODER Pass@1 (%) results on multilingual code generation tasks. “Avg<sub>all</sub>” represents the average Pass@1 scores across all programming languages in the MCEVAL. Here, MCODER-DS indicates that the fine-tuned base model is DeepSeekCoder-1.5-7b-base.

| Method                    | Size   | AWK  | C    | C++    | Cf   | Clisp | Coffee | Dart   | Erlang | Elixir | Erlang | Fortran | F#    | Go     | Groovy | Haskell | HTML | Java | JS   | Json | Julia              |
|---------------------------|--------|------|------|--------|------|-------|--------|--------|--------|--------|--------|---------|-------|--------|--------|---------|------|------|------|------|--------------------|
| DeepSeekCoder-1.5-base    | 7B     | 30.0 | 36.0 | 38.0   | 40.0 | 40.0  | 58.0   | 0.0    | 18.0   | 2.0    | 14.0   | 50.0    | 44.0  | 4.0    | 26.0   | 2.0     | 4.0  | 49.1 | 32.0 | 16.0 | 34.0               |
| CodeQwen-1.5              | 7B     | 38.0 | 40.0 | 46.0   | 42.0 | 28.0  | 56.0   | 2.0    | 14.0   | 14.0   | 0.0    | 40.0    | 46.0  | 44.0   | 32.0   | 2.0     | 0.0  | 47.2 | 52.0 | 30.0 | 60.0               |
| CodeQwen-1.5-Python       | 7B     | 42.0 | 48.0 | 48.0   | 12.0 | 52.0  | 68.0   | 23.5   | 22.0   | 40.0   | 62.0   | 50.0    | 42.0  | 48.0   | 68.0   | 56.0    | 32.0 | 67.9 | 52.0 | 52.0 | 54.0               |
| MCODER-DS                 | 7B     | 34.0 | 46.0 | 50.0   | 26.0 | 30.0  | 72.0   | 19.6   | 6.0    | 26.0   | 24.0   | 58.0    | 30.0  | 48.0   | 12.0   | 26.0    | 28.0 | 67.9 | 48.0 | 62.0 | 48.0               |
| MCODER                    | 7B     | 40.0 | 44.0 | 52.0   | 62.0 | 46.0  | 66.0   | 21.6   | 30.0   | 44.0   | 52.0   | 56.0    | 44.0  | 48.0   | 70.0   | 32.0    | 34.0 | 54.7 | 54.0 | 66.0 | 56.0               |
| Method                    | Kotlin | Lua  | MD   | Pascal | Perl | PHP   | Power  | Python | R      | Racket | Ruby   | Rust    | Scala | Scheme | Shell  | Swift   | Tcl  | TS   | VB   | VimL | Avg <sub>all</sub> |
| DeepSeekCoder-1.5-7B-base | 42.0   | 20.0 | 0.0  | 24.0   | 24.0 | 36.0  | 42.0   | 54.0   | 24.0   | 20.0   | 38.0   | 39.6    | 44.0  | 20.0   | 18.0   | 10.0    | 44.0 | 22.0 | 22.0 | 28.9 |                    |
| CodeQwen-1.5              | 58.0   | 50.0 | 0.0  | 14.0   | 20.0 | 10.0  | 48.0   | 38.0   | 30.0   | 24.0   | 36.0   | 52.8    | 32.0  | 34.0   | 42.0   | 46.0    | 30.0 | 52.0 | 54.0 | 22.0 | 33.2               |
| CodeQwen-1.5-Python       | 46.0   | 46.0 | 24.0 | 42.0   | 36.0 | 36.0  | 54.0   | 44.0   | 36.0   | 40.0   | 46.0   | 52.8    | 58.0  | 40.0   | 42.0   | 62.0    | 48.0 | 52.0 | 58.0 | 18.0 | 45.5               |
| MCODER-DS                 | 36.0   | 42.0 | 22.0 | 34.0   | 8.0  | 34.0  | 46.0   | 42.0   | 22.0   | 40.0   | 56.0   | 45.3    | 48.0  | 30.0   | 38.0   | 48.0    | 34.0 | 46.0 | 50.0 | 28.0 | 37.8               |
| MCODER                    | 48.0   | 52.0 | 30.0 | 42.0   | 36.0 | 32.0  | 54.0   | 44.0   | 40.0   | 36.0   | 48.0   | 52.8    | 58.0  | 44.0   | 46.0   | 64.0    | 38.0 | 52.0 | 58.0 | 20.0 | 46.7               |

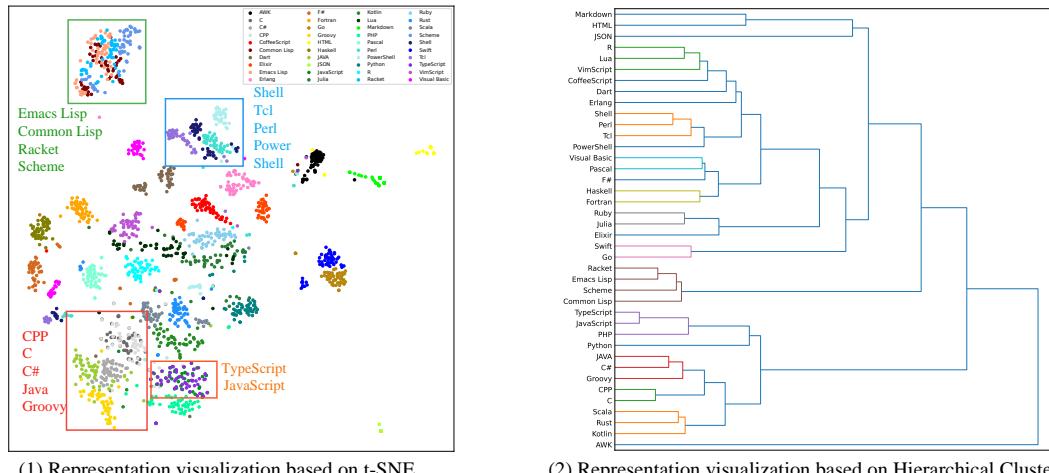


Figure 19: Analysis from the representation perspective on MCEVAL. Languages with similar syntax have closely related representations

**Code Evaluation.** In the domain of code evaluation, a rich tapestry of benchmarks (Zheng et al., 2023b; Yu et al., 2024; Yin et al., 2023; Peng et al., 2024; Khan et al., 2023; Orlanski et al., 2023) has been woven to address the challenges of accurately assessing code quality, functionality, and efficiency, such as HumanEval (Chen et al., 2021), MBPP (Austin et al., 2021), their upgraded version EvalPlus (Liu et al., 2023b). 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 understanding, function calling (Zhuo et al., 2024), code repair (Lin et al., 2017; Tian et al., 2024; Jimenez et al., 2023; Zhang et al., 2023; Prenner & Robbes, 2023; He et al., 2022), code translation (Yan et al., 2023). Recently, many works (Wei et al., 2023; Zhuo et al., 2024) have leveraged LLMs to construct large-scale evaluation datasets and instruction-tuning corpora, further enhancing the evaluation and performance of code models. In our work, we used a similar approach to construct an instruction dataset and proposed the Cross-lingual Code Transfer method to expand the number of languages to 40. Some recent works pay attention to the multilingual scenarios (Cassano et al., 2023; Wang et al., 2023; Athiwaratkun et al., 2023; Zheng et al., 2023a; Peng et al., 2024; Zheng et al., 2023b) by extending the existing python-only HumanEval or MBPP benchmark, such as MultiPL-E (Cassano et al., 2023) and MBXP (Athiwaratkun et al., 2023), which is challenged by the number of the covering languages and data leaking problem (Li et al., 2023; Jain et al., 2024).