



# McEval: Massively Multilingual Code Evaluation

Linzheng Chai<sup>1\*</sup>, Shukai Liu<sup>1\*</sup>, Jian Yang<sup>1\*,†</sup>, Yuwei Yin<sup>2</sup>, Ke Jin<sup>1</sup>, Jiaheng Liu<sup>1</sup>, Tao Sun<sup>1</sup>, Ge Zhang<sup>3</sup>, Changyu Ren<sup>1</sup>, Hongcheng Guo<sup>1</sup>, Zekun Wang<sup>1</sup>, Boyang Wang<sup>1</sup>, Xianjie Wu<sup>1</sup>, Bing Wang<sup>1</sup>, Tongliang Li<sup>4</sup>, Liqun Yang<sup>1</sup>, Sufeng Duan<sup>5</sup>, Zhoujun Li<sup>1</sup>

<sup>1</sup>CCSE, Beihang University, <sup>2</sup>University of British Columbia, <sup>3</sup>University of Waterloo

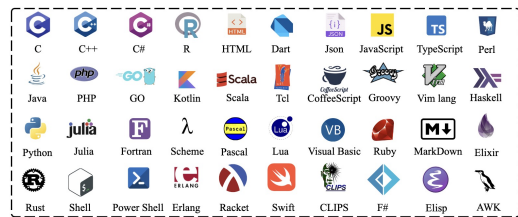
<sup>4</sup>Beijing Information Science and Technology University, <sup>5</sup>Shanghai Jiao Tong University



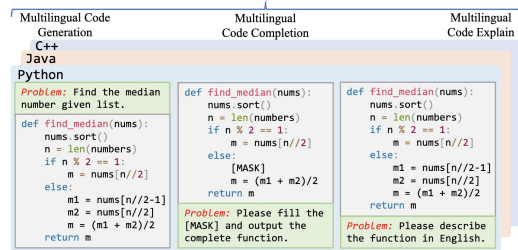
# ICLR

## Introduction

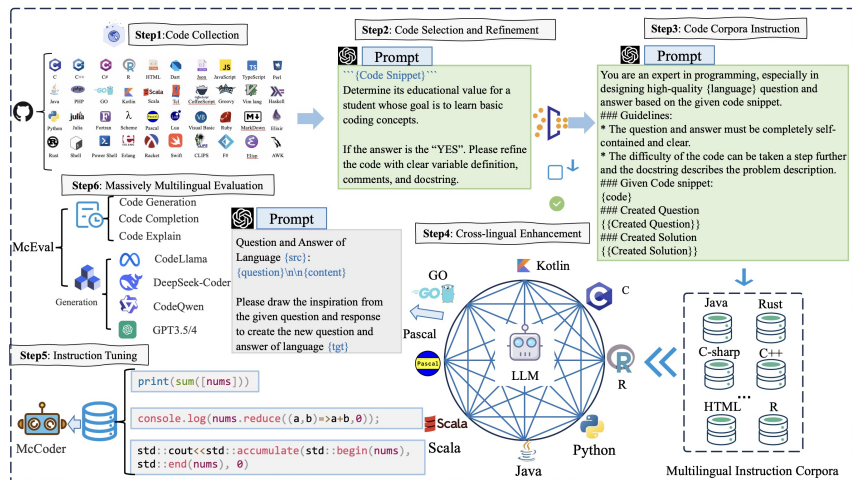
- **McEval**: A massively **multilingual code benchmark** covering **40** programming languages.



Statistics	Value
Questions	
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
Difficulty Level	
- Easy	1,221
- Medium	401
- Hard	385
Length	
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

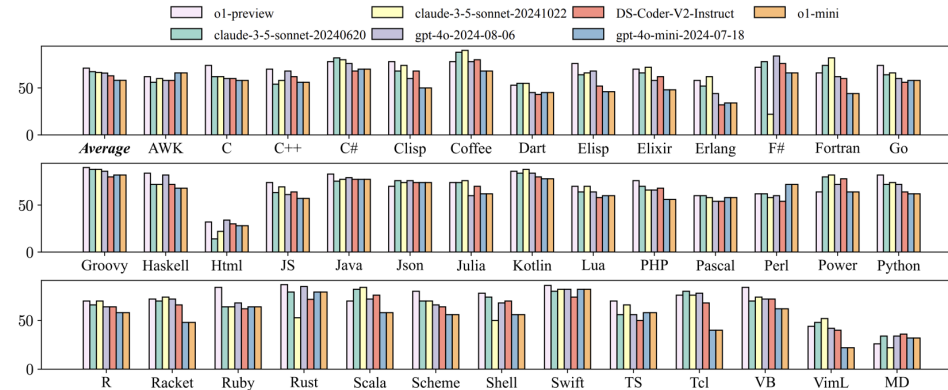


- **McEval-Instruct**: covering 40 languages from code snippets to fine-tune mCoder.



Home Page: <https://mceval.github.io/>

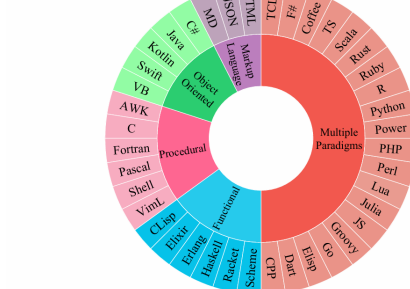
## Evaluation



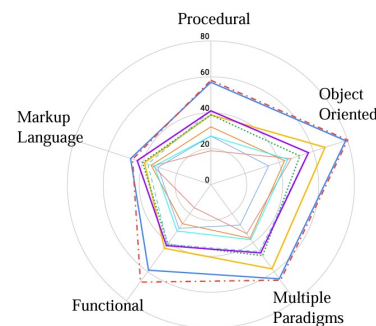
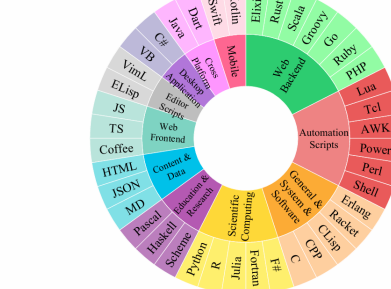
## Model Performance in Code Completion Tasks

We categorize the McEval into **5 paradigms** and **11 application scenarios** and summarize the performance of code LLMs on the code generation task.

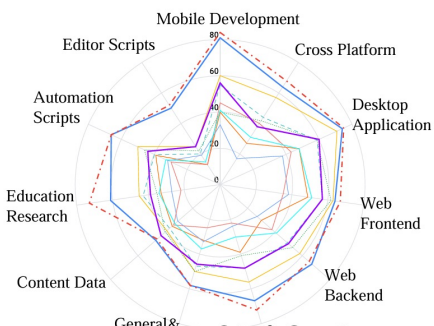
(1) Programming Paradigm



(2) Application Scenario



(1) Performance on Programming Paradigm Categories

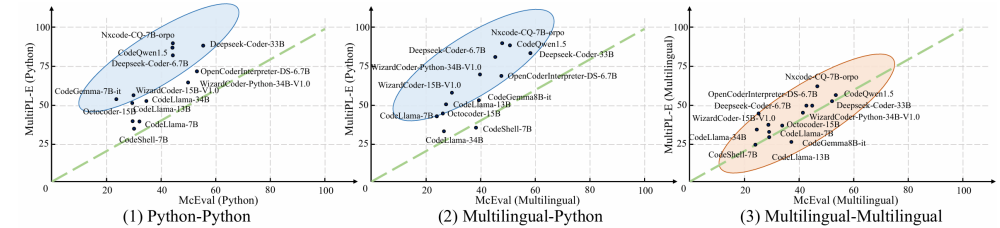


(2) Performance on Application Scenario Categories

## Further Analysis

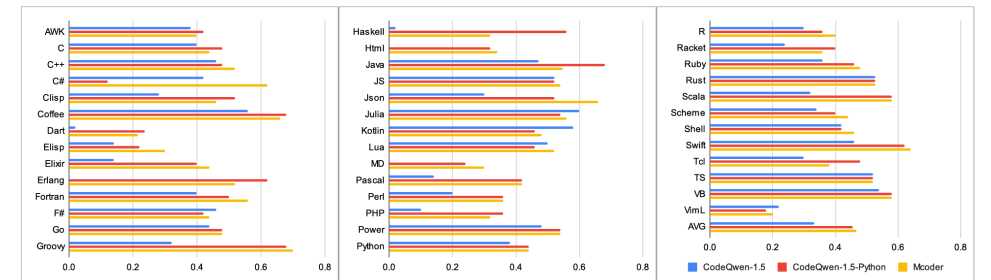
### Unbalance in Different Languages

We compare the results of several open-source models on the MultiPL-E multilingual benchmark with corresponding languages on McEval.

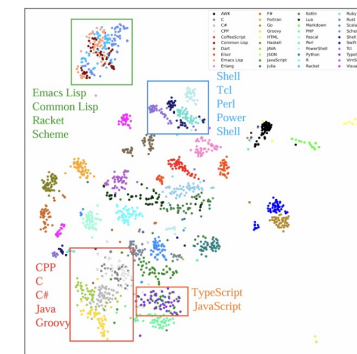


### Cross-lingual Transfer

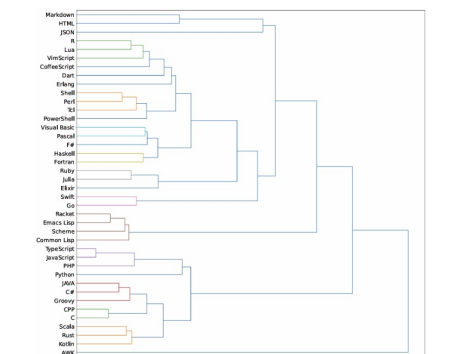
**Fine-tuning with Python-only data** can still effectively transfer instruction-following abilities to other languages, resulting in superior multilingual performance.



### Analysis of Language Representations



(1) Representation visualization based on t-SNE



(2) Representation visualization based on Hierarchical Cluster