# OOP: Object-Oriented Programming Evaluation Benchmark for Large Language Models

**Shuai Wang**[1]  **Liang Ding**[2]  **Li Shen**[3]  **Yong Luo**[1]  **Bo Du**[1]  **Dacheng Tao**[2]

[1]Wuhan University  [2]The University of Sydney  [3]JD Explore Academy

wangshuai123@whu.edu.cn, liangding.liam@gmail.com

## Abstract

Advancing automated programming necessitates robust and comprehensive code generation benchmarks, yet current evaluation frameworks largely neglect **o**bject-**o**riented **p**rogramming (**OOP**) in favour of functional programming (FP), e.g., HumanEval and MBPP. To address this, ❶ our study introduces **a pioneering OOP-focused benchmark**, featuring 431 Python programs that encompass essential OOP concepts and features like classes and encapsulation methods. ❷ We propose **a novel evaluation metric**, *pass@o*, tailored for OOP, enhancing traditional *pass@k* metric. ❸ Our evaluation of 23 leading large language models (LLMs), including both general and code-specialized models, **reveals three key insights**: 1) *pass@o* offers a more relevant and comprehensive assessment for OOP code generation; 2) Despite excelling in FP, code-specialized LLMs like WizardCoder lag in OOP compared to models like ChatGPT; 3) The poor performance of all advanced LLMs on our OOP benchmark highlights a critical need for improvements in this field. Our benchmark and scripts are publicly released at: https://github.com/alphadl/OOP-eval.

## 1 Introduction

Large language models (LLMs, Ouyang et al., 2022a; Touvron et al., 2023), consisting of billions or even trillions of parameters' Transformer blocks (Vaswani et al., 2017), have emerged like mushrooms after the rain, especially since the emergence of ChatGPT[1]. In comparison to small models, LLMs exhibit stronger generalization and reasoning capabilities (Wei et al., 2022). Currently, LLMs are playing a crucial role in various tasks, e.g., code generation (Chen et al., 2021; Li et al., 2022; Roziere et al., 2023), language understanding (Zhong et al., 2023), human-computer interaction (Tolomei et al., 2023; Moslem et al., 2023),
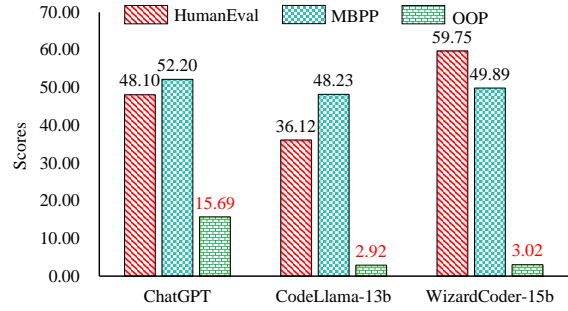
---

[1]https://chat.openai.com



Figure 1: **The performance comparison of widely-used code language models on functional programming (FP) and object-oriented programming (OOP) code generation benchmarks**, in terms of *pass@1* scores. We see that all models perform relatively well on FP benchmarks, i.e., Humaneval (Chen et al., 2021) and MBPP (Austin et al., 2021), while exhibiting poor performance on our OOP benchmark.

and translation (Peng et al., 2023; Lu et al., 2023).

The process of code generation entails crafting code in a suitable programming language from natural language descriptions of problems or requirements, aiming to effectively solve the problems or fulfill the requirements. Given that hiring professional programmers to write code consumes a significant amount of human and material resources, the importance of automated programming becomes particularly evident. Currently, the question of how to use the rising LLMs to generate more accurate automated programming codes based on problems or requirements stated by actual natural language has become an important research topic (Liu et al., 2023; Zhong and Wang, 2023). In the research process, code generation evaluation is crucial. Code generation evaluation not only needs to objectively and impartially reflect the current performance of LLMs in programming but also should disclose the shortcomings in LLM programming to further enhance its potential.

**[Importance of OOP]** According to the Novem-

| Benchmark | Number | NL | PL | Task Type |
|---|---|---|---|---|
| HumanEval (Chen et al., 2021) | 164 | en | Python | Function Programming |
| MBPP (Austin et al., 2021) | 974 | en | Python | Function Programming |
| APPS (Hendrycks et al., 2021) | 5000 | en | Python | Function Programming |
| CodeContests (Li et al., 2022) | 165 | en | Multi | Function Programming |
| MultiPL-MBPP (Cassano et al., 2023) | 974* | Multi | Multi | Function Programming |
| HumanEval-X (Zheng et al., 2023) | 164* | en | Multi | Function Programming |
| MultiPL-HumanEval (Cassano et al., 2023) | 164* | en | Multi | Function Programming |
| MTPB (Nijkamp et al., 2022) | 115 | en | Python | Function Programming |
| ODEX (Wang et al., 2022) | 945 | Multi | Python | Function Programming |
| PandasEval (Zan et al., 2022) | 101 | en | Python | Function Programming |
| BIG-Bench (Srivastava et al., 2022) | 32 | en | Python | Function Programming |
| CodeApex* (Fu et al., 2023) | 476* | zh&en | C++ | Function Programming |
| OOP (Our) | 431 | en | Python | Object-Oriented Programming |

Table 1: **Overview of existing code evaluation benchmarks**. ("NL" denotes natural language describing the problem or requirements; "PL" represents the generated programming language; "en" and "zh" denote English and Chinese, respectively, and "Multi" means containing multiple NLs or PLs; "*" indicates the number of samples for each language; "★" means that in CodeApex, we only considered code generation tasks.)

ber programming language rankings by TIOBE [2], four out of the top five programming languages are OOP, which reflects the importance of OOP languages. OOP is centred on designing code around data or objects rather than organizing it based on functionality and logic (Stroustrup, 1988; Stefik and Bobrow, 1985). OOP focuses more on the programming paradigm of class and object (Wegner, 1990). Functions are commonly referred to as methods in OOP.

**[Motivation]** However, existing code generation evaluation benchmarks primarily focus on the evaluation of FP, and lack the evaluation of relevant concepts and features of OOP, e.g., class, inheritance, encapsulation methods, etc. If using existing benchmarks in Table 1 for evaluation can only show the performance of LLMs in FP, it fails to reflect their potential in OOP, as illustrated in Figure 1.

**[OOP benchmark and metric]** Considering the limitations of current code generation evaluation FP benchmarks and the widespread use of the Python programming language, we propose the first OOP evaluation benchmark based on Python. OOP benchmark consists of 431 Python programs, covering key concepts and features of OOP, including class, inheritance, encapsulation methods, etc. Furthermore, to prevent the issue where LLMs may not generate concepts and features of OOP, we have optimized the $pass@k$ (Kulal et al., 2019; Chen et al., 2021) metric by matching key points in natu-

ral language with key points in the programming language, i.e., the class names and private function names, etc, for natural language requirements are matched with the class names and private function names, etc., in the programming language. Our main **contributions** are summarized as follows:

1. We construct and release the first OOP evaluation benchmark, which encompasses concepts and features of OOP, e.g., class, polymorphism, encapsulation methods, etc.

2. We devise a new metric $pass@o$ based on conventional $pass@k$, tailored for the OOP code generation task, by matching key points in natural language and programming language.

3. We extensively evaluated our OOP with 23 advanced LLMs, demonstrating that *i)* there is still significant room for improving the OOP tasks, *ii)* our benchmark could serve as a robust and fair indicator that helps the community quantify LLMs' OOP performance.

## 2 Related work

**Code Evaluation Benchmark** In the early days of LLMs, researchers from Google and OpenAI launched artificial handwritten code evaluation benchmarks, namely MBPP (Austin et al., 2021) and HumanEval (Chen et al., 2021), respectively. MBPP and HumanEval are currently the mainstream code generation evaluation benchmarks, but both of them are based on the Python programming

language. Subsequently, MultiPL-MBPP (Cassano et al., 2023) and MultiPL-HumanEval (Cassano et al., 2023) expanded upon these two benchmarks by translating the Python programming language into eighteen other programming languages, e.g., Java, C++, PHP, etc, to evaluate the performance of LLMs across others programming languages. Additionally, HumanEval-X (Zheng et al., 2023) incorporated multiple test cases into the HumanEval benchmark. Apart from the extensions made to these two benchmarks, other benchmarks like CodeApex (Fu et al., 2023) and ODEX (Wang et al., 2022) exhibit distinctive features across different natural languages and task types. Unlike existing code evaluation benchmarks, our proposed OOP benchmark primarily focuses on the concepts and features of OOP, e.g., class, inheritance, etc. These works are summarized in Table 1.

**Code Evaluation Metrics** Existing evaluation metrics can be broadly categorized into two types: dynamic evaluation metrics and static evaluation metrics. Dynamic evaluation metrics evaluate the executability of generated codes by using test cases, with $pass@k$ (Kulal et al., 2019; Chen et al., 2021) serving as the primary representative. The calculation process for $pass@k$ is shown in Appendix A. Additionally, this category of metrics includes $n@k$ (Li et al., 2022). Static evaluation metrics calculate BLUE (Papineni et al., 2002), ROUGE (Lin, 2004), Codescore (Dong et al., 2023) and Code-BLEU (Ren et al., 2020) among manually written examples and generated programs. However, these code evaluation metrics do not specifically focus on evaluating the concepts and features of OOP. Therefore, we further optimized the $pass@k$ metric based on the evaluation benchmark for OOP.

## 3 Evaluation Framework

### 3.1 Overview

Existing code generation benchmarks in Table 1 for are confined to FP and do not involve essential concepts and features of OOP. We take the frequently used benchmarks, HumanEval (Chen et al., 2021) and MBPP (Austin et al., 2021) in Table 1, as examples. They primarily evaluate the capabilities of LLMs in FP. The detailed descriptions of HumanEval and MBPP are provided in Appendix B. If we use existing benchmarks in Table 1 for evaluation, it does not show the capability of LLMs in OOP, as illustrated in Figure 1, that is, the seem-



Figure 2: **The generation of private functions cannot be evaluated using** $pass@k$. (We instructed ChatGPT (Ouyang et al., 2022b; OpenAI, 2023) model to generate the class *class SS*, public function *public_Shortest_subarray*, and private function *def __private_Shortest_subarray* based on a given prompt and implement the corresponding requirements within the functions. However, ChatGPT does not generate the private functions named *private_Shortest_subarray* outlined in the <mark>red box</mark>.)

ingly decent LLMs (on FP tasks) perform relatively worse on OOP tasks. In addition, existing code generation evaluation metrics primarily use $pass@k$ to evaluate the executability of the generated code. However, using the $pass@k$ metric can not reflect whether LLMs generate concepts and features related to OOP, as illustrated in Figure 2. Therefore, $pass@k$ can not objectively and fairly reflect the OOP capabilities of LLMs.

As a result, we established an OOP benchmark and proposed the evaluation metric $pass@k$ for OOP. The process for constructing the OOP benchmark is illustrated in Figure 3.

### 3.2 Building OOP Benchmarks

**Data Filtering.** The training data for current LLMs mostly comes from the internet. If we directly evaluate LLMs using existing OOP data from the web, it would not reflect the OOP capabilities of LLMs. Therefore, we first **rigorously selected 500 natural language description-based problems or requirements based on Python** from plat-
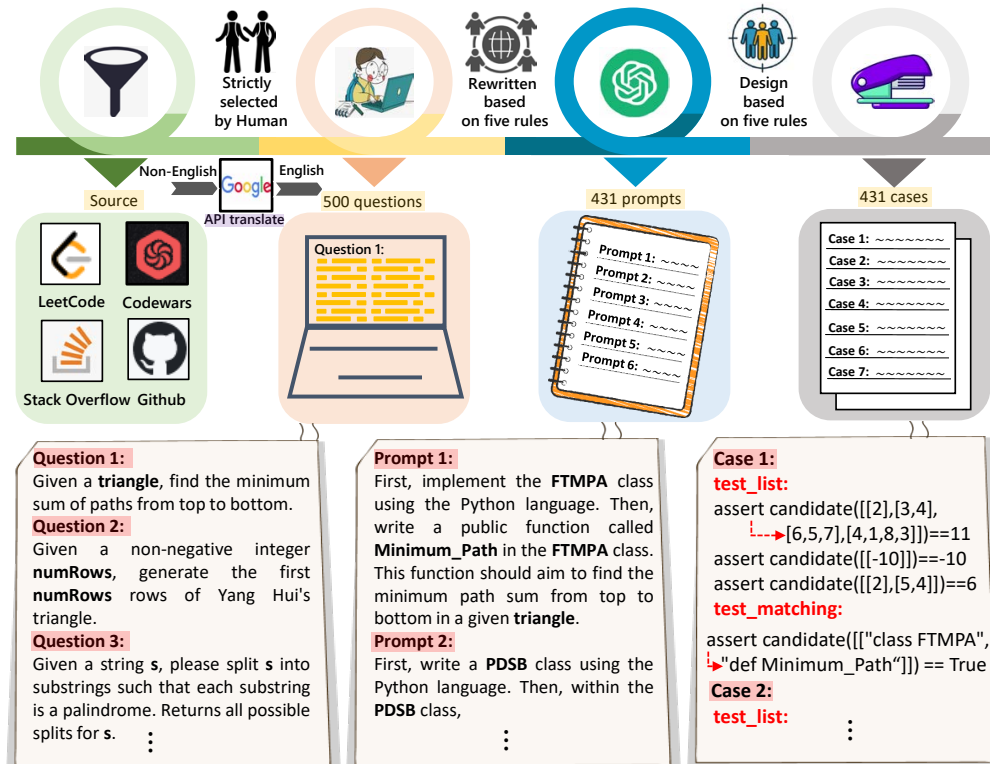
Figure 3: **The construction process of our object-oriented programming (OOP) benchmark**.

forms like LeetCode [3], open-source repositories on GitHub [4], Stack Overflow [5], and Codewars [6]. These 500 questions or requirements only are limited to FP and do not involve concepts and features related to OOP.

**Human Rewritten.** Subsequently, we manually rewrite the collected 500 questions or requirements by adhering to the following rules:

1. Designing, based on the problems or requirements, with relevant OOP concepts and features, e.g., class names, inheritance name (i.e., parent class name), encapsulation methods name (i.e., public function name and private function names), etc.

2. Related problems or requirements are implemented within the public function and private function of the class while ensuring the encapsulation of that implementation.

3. Convert the variables associated with problems or requirements into class attribute vari-

ables, ensuring that these variables are accessible in both public and private functions.

4. If the implementation of problems or requirements is placed within the private function of the class, it is necessary to design a corresponding public function for access.

5. The rewritten OOP relevant problems or requirements can be successfully implemented and accessed through objects.

Following the five rules mentioned above, we conducted a standardized rewriting of the 500 Python-based problems or requirements.

**Case Design.** Finally, we designed corresponding test cases to evaluate OOP. Finally, we **obtained 431 samples of OOP**, as shown in Figure 3. The specific construction details of OOP are provided in Appendix C.

**Level Classification.** Given the difficulty nature of programming, we divided the designed OOP benchmark into three levels: **Simple-level OOP**, **Moderate-level OOP**, and **Difficult-level OOP**, as shown in Figure 7.

Simple-level OOP has 77 program samples, and includes only class, and public function. Moderate-level OOP builds upon simple-level OOP by adding

attribute variables and private functions, and has 179 program samples. Nevertheless, the difficult-level OOP is based on the Simple-level of OOP, and adds inheritance, polymorphism and other related concepts and features of OOP. There are a total of 175 program samples for difficult-level OOP. Although private functions are not involved in the difficulty level, the problems or requirements in difficult-level OOP are more complex and varied. Using such a level of classification, we can not only evaluate the performance of existing LLMs in OOP but also analyze the shortcomings of LLMs, which allows us to better unearth the potential of LLMs in OOP. Using this approach makes it more convenient for us to improve the OOP performance of LLMs.

### 3.3 Evaluation Metrics *Pass@o*

To evaluate whether LLMs generate concepts and features related to OOP, i.e., generated subclass name, parent class name, private function name and public function name, etc, in the programming language, we proposed a *pass@o* metric based on OOP. The *pass@o* metric adds keyword points matching between natural language with programming language based on the *pass@k*, i.e.,

$$\alpha = \sum_{i=1}^{n} f(X_i),$$

$$where f(X_i) = \qquad (1)$$

$$\begin{cases} 1, if\ utf\ (X_i)\ passed\ and\ \sum_{j}^{m} x_j \exists X_i, \\ 0, otherwise \end{cases}$$

$$pass@o := \mathop{\mathbb{E}}_{Problems} \left[ 1 - \frac{\binom{n-\alpha}{o}}{\binom{n}{o}} \right], \qquad (2)$$

In Eq. (1), $n$ represents the number of code generations for a given problem; $X_i$ represents the $i$-th generated program code; $\alpha$ represents the quantity of $n$ generated codes passing tests and matches; $ut(\cdot)$ denotes the unit test function; $m$ represents the number of keyword points in the current $prompt$; and $x_j$ represents the $j$-th keyword points in natural language. In Eq. (2), $o \leq n$.

The *pass@o* metric not only optimizes the limitations of *pass@k* evaluation but also objectively and fairly reflects the OOP performance of LLMs.

## 4 Experiments

### 4.1 Experimental Setup

**Evaluated LLMs** In the OOP task, we conduct experiments on 23 mainstream LLMs. These models include both general LLMs, i.g., Chat-GPT (Ouyang et al., 2022b; OpenAI, 2023), Llama2 (Touvron et al., 2023), InternLm (Team, 2023a), MPT (Team, 2023b), DeepSeek (Team, 2024), Falcon (Almazrouei et al., 2023), Qwen (Bai et al., 2023), Yi [7] and code-specialized LLMs, e.g., CodeLlama (Roziere et al., 2023), Wizard-Coder (Luo et al., 2023), StarCoder (Li et al., 2023), as shown in Table 6. The details description of 24 LLMs are shown in Appendix D.

**Parameter Settings.** In the experiment, we followed the settings on Llama2 (Touvron et al., 2023), configuring the temperature to 0.1 and 0.8 for code generation. The remaining parameters $(top-p = 0.95, n = 200, o \leq n)$, consistently remained unchanged. We evaluate the OOP benchmark on eight NVIDIA A100 GPUs using the vllm (Kwon et al., 2023) 0.2.1.post1 framework [8].

**Metrics.** In terms of evaluation metrics, we use for *pass@k* and the proposed *pass@o* metrics.

### 4.2 Overall Evaluation Result

The evaluation results of the LLMs with temperatures set to 0.1 and 0.8 are presented in Table 2, respectively. From the experimental results, We have obtained the following conclusions:

**The OOP capabilities of the existing LLMs fall far short of the ideal state** In Table 2, we can observe that LLMs with strong coding capabilities (e.g., WizardCoder-15b, CodeLlama-7b-Python, and CodeLlama-13b, achieved scores of 58.12, 40.48, and 35.07, respectively, in the HumanEval code leaderboard [9]), exhibit performance in OOP benchmarks that falls significantly short of the ideal state. WizardCoder-15b, CodeLlama-7b-Python, and CodeLlama-13b scored 3.02, 1.27, and 2.92, respectively, on the OOP benchmark at *pass@1*. Their scores on *pass@100* were also 10.01, 14.07, and 13.11, respectively. Even the current ChatGPT model with strong general capabilities scores 15.69

---

| | Model | **1** | | | **80** | | | **100** | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | pass@k | pass@o | Δ (↓) | pass@k | pass@o | Δ (↓) | pass@k | pass@o | Δ (↓) |
| General | Falcon-7b | 0.01 | 0.00 | **-0.01** | 0.37 | 0.19 | **-0.18** | 0.47 | 0.23 | **-0.24** |
| | Falcon-40b | 0.01 | 0.00 | **-0.01** | 2.90 | 1.11 | **-1.79** | 3.42 | 1.26 | **-2.16** |
| | Llama2-7b | 0.01 | 0.01 | **-0.00** | 4.02 | 1.72 | **-2.30** | 4.62 | 1.94 | **-2.68** |
| | InternLm-7b | 0.03 | 0.02 | **-0.01** | 1.04 | 0.52 | **-0.52** | 1.22 | 0.58 | **-0.64** |
| | Yi-6b | 0.07 | 0.01 | **-0.06** | 5.07 | 1.67 | **-3.40** | 6.00 | 1.98 | **-4.02** |
| | Llama2-13b | 0.09 | 0.06 | **-0.03** | 7.28 | 2.17 | **-5.11** | 8.24 | 2.41 | **-5.83** |
| | MPT-7b | 0.28 | 0.02 | **-0.26** | 4.77 | 1.27 | **-3.50** | 5.50 | 1.46 | **-4.04** |
| | Qwen-7b | 0.94 | 0.61 | **-0.33** | 15.02 | 5.68 | **-9.34** | 16.35 | 5.83 | **-10.52** |
| | Qwen-14b | 1.52 | 0.75 | **-0.77** | 26.28 | 10.58 | **-15.70** | 28.10 | 11.48 | **-16.62** |
| | DeepSeek-7b | 1.53 | 0.50 | **-1.03** | 16.83 | 7.72 | **-9.11** | 18.70 | 8.70 | **-10.00** |
| | Yi-34b | 2.20 | 1.09 | **-1.11** | 21.96 | 8.43 | **-13.53** | 23.68 | 9.22 | **-14.46** |
| | Llama2-70b | 3.55 | 1.25 | **-2.30** | 21.01 | 9.97 | **-11.04** | 23.14 | 11.16 | **-11.98** |
| | DeepSeek-67b | 8.02 | 3.71 | **-3.95** | 49.31 | 27.42 | **<u>-21.89</u>** | 51.60 | 29.47 | **<u>-22.13</u>** |
| | Qwen-72b | 11.20 | 4.62 | **-6.58** | 57.48 | 35.70 | **-21.78** | 59.52 | 37.83 | **-21.69** |
| | ChatGPT | **42.88** | **15.69** | **<u>-27.19</u>** | **75.71** | **58.28** | **-17.43** | **76.20** | **59.80** | **-16.40** |
| Specialized | GPT_BigCode | 0.10 | 0.06 | **-0.04** | 7.00 | 2.58 | **-4.42** | 8.01 | 2.92 | **-5.09** |
| | CodeLlama-7b | 2.67 | 1.20 | **-1.47** | 24.09 | 9.16 | **-14.93** | 25.72 | 9.92 | **-15.80** |
| | CodeLlama-13b-Python | 2.80 | 1.03 | **-1.77** | 36.34 | 17.22 | **-19.12** | 38.75 | 18.96 | **-19.79** |
| | StarCoder | 4.61 | 1.26 | **-3.35** | 28.67 | 10.05 | **-18.62** | 30.44 | 10.88 | **-19.56** |
| | CodeLlama-7b-Python | 4.68 | 1.27 | **-3.41** | 28.68 | 12.90 | **-15.78** | 30.33 | 14.07 | **-16.26** |
| | CodeLlama-34b | 6.24 | 1.58 | **-4.66** | **46.31** | **22.59** | **<u>-23.72</u>** | **49.01** | **24.68** | **<u>-24.33</u>** |
| | WizardCoder-15b | 6.83 | **3.02** | **-3.81** | 28.10 | 9.50 | **-18.60** | 29.41 | 10.01 | **-19.40** |
| | CodeLlama-13b | **6.87** | 2.92 | **-3.95** | 32.69 | 12.20 | **-20.49** | 34.53 | 13.11 | **-21.42** |

Table 2: **Performance of** 23 **large language models (LLMs) on object-oriented programming (OOP) tasks**. We also reported the differences in evaluation results between *pass@k* and *pass@o*. (All LLMs are evaluated in zero-shot fashion. For *pass@*100 and *pass@*80 scores, we use a temperature of 0.8 and top-*p*=0.95. For *pass@*1 scores, we use a temperature of 0.1 and top-*p*=0.95. The best results are highlighted in black bold; <mark>Red</mark> indicates the differences evaluated using the *pass@o* and *pass@k* metrics; <u>Underlined</u> indicates the maximum disparities evaluated between *pass@o* and *pass@k* metrics; <mark>Gray</mark> indicates models with a larger number of parameters.)

on *pass@*1 and 59.80 on *pass@*100. The results indicate that the untapped potential of existing LLMs in OOP has not been fully explored.

**Limitations of *pass@k* evaluated OOP** The scores from Table 2 indicate that using *pass@k* does not objectively reflect the OOP performance of LLMs, e.g., the WizardCoder-15b model achieves scores of 6.83, 28.10, and 29.41 using *pass@k*, while its scores drop to 3.02, 9.50, and 10.01 when using *pass@o*. The evaluation scores of other LLMs using *pass@o* in Table 2 showed a decline, once again proving the limitations of *pass@k* in evaluation OOP.

In addition, we also observed a significant phenomenon, e.g., when evaluated using *pass@k*, Qwen-14b (score 26.28) scored lower than WizardCoder-15b (score 28.10) on *pass@*80. However, when evaluated using *pass@o*, Qwen-14b (score 10.58) scored higher than WizardCoder-15b (score 9.50) on *pass@*80. Analyzing the experimental results of Qwen-14b and WizardCoder-15b, we observed that when evaluated using *pass@o*, Qwen-14b outperforms WizardCoder-15b in terms

of the ability to correctly generate OOP concepts and feature keywords, as illustrated in Figure 4. It also reiterates that *pass@k* cannot objectively and fairly reflect the evaluation results of OOP.

**A larger model scale does not necessarily perform better on *pass@*1** In Table 2, when evaluated using *pass@o*, CodeLlama-34b scores 1.58 on *pass@*1, whereas CodeLlama-13b scores 2.92 on *pass@*1. Additionally, CodeLlama-13b-Python scores 1.03 on *pass@*1, while the corresponding CodeLlama-7b-Python scores 1.27 on *pass@*1. However, CodeLlama-7b scores 1.20 on *pass@*1, which is lower than the score achieved by CodeLlama-13b. The scores of CodeLlama-7b, CodeLlama-13b, CodeLlama-34b, CodeLlama-7b-Python, and CodeLlama-13b-Python on *pass@*1 indicate that a larger model scale does not necessarily result in the highest scores on *pass@*1.

### 4.3 Different-level Evaluation Results

Following the classification of OOP benchmarks in Section 3.2, we conducted evaluations for three levels of OOP benchmarks, and the results are pre-
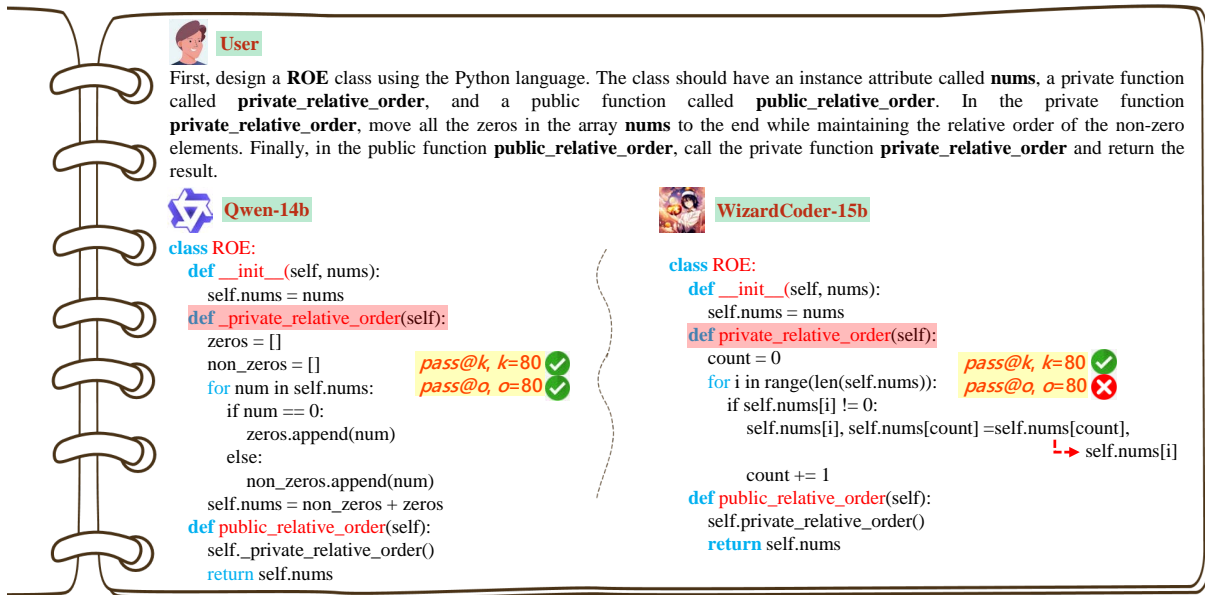
Figure 4: **The case comparison of generation results between Qwen-14b and WizardCoder-15b in the OOP benchmark**. We see: 1) Qwen-14b can accurately generate private functions, while WizardCoder-15b cannot accurately generate private functions; 2) The results generated by Qwen-14b and WizardCoder-15b can both pass the evaluation using $pass@k$; 3) The results generated by Qwen-14b can pass the evaluation using $pass@o$, but the results generated by WizardCoder-15b cannot pass the evaluation using $pass@o$.

sented in Tables 3, 4, and 5, respectively. we have drawn the following conclusions:

**LLMs perform better at the simple-level OOP compared to the moderate-level and difficult-level OOP**   From the simple-level OOP evaluation results in Table 3, we can see that the evaluation results using $pass@k$ and $pass@o$ are the same. It also indicates that LLMs can comprehend the fundamental concepts and features of OOP, e.g., class, and encapsulation methods (i.e., public function). However, in Tables 4 and 5, LLMs exhibit a weaker understanding of concepts and features related to OOP, e.g., encapsulation methods (i.e., private function), inheritance, and polymorphism, and are unable to generate corresponding code accurately. Detailed descriptions are in Appendix E.

**ChatGPT has large gap in moderate-level usage using $pass@k$ and $pass@o$**   From the results in Table 4, We observe that with $pass@k$ evaluation, ChatGPT scores are 51.71, 83.30, and 83.67, but with $pass@o$ evaluation, ChatGPT only achieves scores of 2.53, 51.54, and 54.78. We analyzed the moderate-level OOP results for ChatGPT, and found that its understanding of private functions is relatively poor. When evaluated using $pass@k$, a total of 5551 codes can pass the test cases correctly. However, when evaluated using $pass@o$, only 272

codes can successfully pass the test cases. Among them, 5279 codes fail to match the $pass@o$ criteria. Upon careful examination, we found that all these 5279 codes resulted from errors generated by private functions To further validate the authenticity of the experimental results, we randomly selected prompts corresponding to three error results. Subsequently, we input prompts of the erroneous results into the web version of ChatGPT for code generation, as illustrated in Figure 13, 14 and 15. We found that the code generated by online ChatGPT [10] is also private function error.

**ChatGPT outperforms moderate-level in difficult-level evaluation results**   According to the evaluation results from Tables 3 and 4, we observe that the performance of ChatGPT at the difficult level is stronger than at the moderate level. At the difficult-level OOP, ChatGPT scores are 19.70, 71.83, and 73.37, whereas at the moderate-level OOP, ChatGPT scores are only 2.53, 51.54, and 54.78.

## 5   Discussion

In this section, we will explore the reasons behind the generally lower scores of LLMs in OOP, as well as the applicability of the Chain-of-Thoughts

---
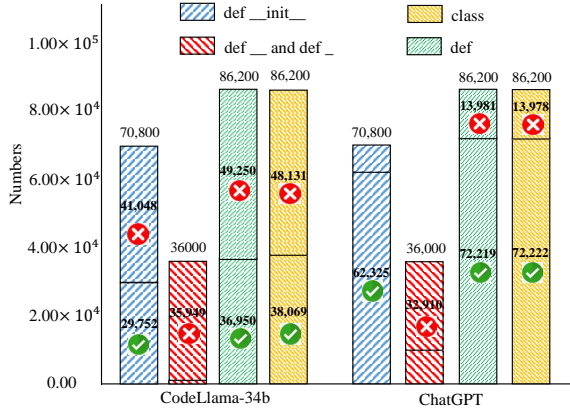[10]https://chat.openai.com/

Figure 5: **Distribution of search results for ChatGPT and CodeLlama-34b.** (In program, "**class**" serves as the indicator for program *class names*. If the program does not contain a "**class**", it signifies an error in the generation of *class names* by the LLM. Similarly, it can be deduced that "**def _**" and "**def __**" serve as indicators for *private function names*; "**def**" signifies a *public function name*; and "**def __init__**" represents the indicator for *attribute variables name*. Moreover, In our OOP benchmark, the LLM should ideally generate at least 86,200 "**class**", 36,000 "**def __**" or "**def _**", 86,200 "**def**", and 70,800 "**def __init__**".)

(CoT) method to OOP.

**Why LLMs score lower in OOP benchmarks?**
We use the experimental results of ChatGPT and CodeLlama-34b on *pass@1* as examples for analysis. As we instruct LLMs to generate relevant class names, private function names, public function names, etc., We conducted searches using simple keywords, e.g., "**class**", "**def _**", "**def __**", "**def __init__**", and "**def**" on both CodeLlama-34b and ChatGPT results. A detailed description of the retrieval process is provided in Appendix F. We compiled and analyzed the distribution of retrieval "**class**", "**def _**", "**def __**", "**def __init__**", and "**def**", as shown in Figure 5, concluding that: 1) Weak knowledge, e.g., class, encapsulation methods, etc, of OOP in LLMs; 2) LLMs particularly lack cognition of private functions; 3) There is a certain degree of gap between CodeLlama-34b and ChatGPT. Specific example is shown in Figure 11.

**The applicability of CoT in OOP.** Taking CodeLlama-13b, StarCoder, and WizardCoder-15b as examples, we respectively incorporate the few-shot, zero-shot CoT, and few-shot CoT methods to validate whether CoT approaches demonstrate applicability in OOP, as shown in Table 7. We observed a significant improvement in the scores

of LLMs in OOP when using the few-shot approach, e.g., CodeLlama-13b achieved scores of 14.50, 48.13, and 49.85 using the few-shot method, representing improvements of 396.58%, 294.51%, and 280.24%, respectively, compared to the zero-shot method. In Table 7, we also observe that CodeLlama-13b achieves scores of 1.33, 13.31, and 14.62 in zero-shot CoT, but its score at *pass@1* is lower at 2.92 compared to zero-shot. Additionally, StarCoder scores 0.25, 6.58, and 7.07 in zero-shot CoT, which are lower than StarCoder scores in zero-shot at 1.26, 10.05, and 10.88, respectively. The scores of the CodeLlama-13b, StarCoder, and WizardCoder-15b models on few-shot CoT are also lower than their scores on few-shot. We analyzed the experimental results of zero-shot and zero-shot CoT and found that using the CoT method introduces an illusion to the model, preventing it from directly generating the corresponding code, as illustrated in Figure 12. Therefore, it is necessary to integrate the concepts and features of OOP to design appropriate CoT strategies in order to enhance the effectiveness of generating OOP by LLMs. Appendix G provides detailed prompts for few-shot, zero-shot CoT, and few-shot CoT.

## 6 Conclusion

In this paper, we propose the first OOP evaluation benchmark based on Python, consisting of 431 Python programs, encompassing key concepts and features of OOP, e.g., class, encapsulation methods, etc. Simultaneously, we propose the evaluation metric *pass@o* for the OOP benchmark. *pass@o* improves upon the limitations of *pass@k* by matching keyword points between natural language with program language. We evaluate 23 mainstream LLMs using the proposed OOP benchmark and *pass@o* metric. Experimental results show that the current OOP of LLMs is far from ideal, which also reveals that LLMs have room for further improvement. Furthermore, Existing LLMs have a certain gap with ChatGPT in OOP. Moreover, we also investigate that applying some of the current improvement strategies directly to the OOP benchmark does not show significant improvement. In the future, we need to further strengthen the OOP knowledge of LLMs, especially regarding private functions. At the same time, we also hope that more researchers can contribute to the advancement of research in OOP.

## Limitations

Our OOP benchmark has several limitations: (1) Our proposed OOP benchmark is based on the Python programming language and does not cover other OOP languages. (2) Given the incorporation of crucial concepts like polymorphism and inheritance in the OOP benchmark, it does not specifically address challenges associated with more intricate scenarios, e.g., multiple inheritance and overloading. (3) While OOP languages hold a significant share, non-OOP languages, e.g., C and Go languages, also play irreplaceable roles. In future work, we plan to consider expanding the OOP benchmark to cover a broader spectrum. Additionally, we encourage researchers to explore the potential of LLMs through evaluations based on the OOP benchmark.

## Ethics Statement

We take ethical considerations very seriously. This paper focuses on establishing benchmarks for OOP to analyze the performance of existing LLMs. Our research reveals that existing LLMs fall far short of ideal performance in OOP. We conducted experiments on open and publicly available LLMs and accurately and objectively report the findings and conclusions of this paper. Therefore, we believe that this study does not raise ethical concerns.

## References

Loubna Ben Allal, Raymond Li, Denis Kocetkov, Chenghao Mou, Christopher Akiki, Carlos Munoz Ferrandis, Niklas Muennighoff, Mayank Mishra, Alex Gu, Manan Dey, et al. 2023. Santacoder: don't reach for the stars! *arXiv preprint*.

Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, Mérouane Debbah, Étienne Goffinet, Daniel Hesslow, Julien Launay, Quentin Malartic, Daniele Mazzotta, Badreddine Noune, Baptiste Pannier, and Guilherme Penedo. 2023. The falcon series of open language models. *arXiv preprint*.

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

Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. 2023. Qwen technical report. *arXiv preprint*.

Federico Cassano, John Gouwar, Daniel Nguyen, Sydney Nguyen, Luna Phipps-Costin, Donald Pinckney, Ming-Ho Yee, Yangtian Zi, Carolyn Jane Anderson, Molly Q Feldman, et al. 2023. Multipl-e: a scalable and polyglot approach to benchmarking neural code generation. *IEEE TSE*.

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. Evaluating large language models trained on code. *arXiv preprint*.

Yihong Dong, Jiazheng Ding, Xue Jiang, Zhuo Li, Ge Li, and Zhi Jin. 2023. Codescore: Evaluating code generation by learning code execution. *arXiv preprint*.

Lingyue Fu, Huacan Chai, Shuang Luo, Kounianhua Du, Weiming Zhang, Longteng Fan, Jiayi Lei, Renting Rui, Jianghao Lin, Yuchen Fang, et al. 2023. Codeapex: A bilingual programming evaluation benchmark for large language models. *arXiv preprint*.

Dan Hendrycks, Steven Basart, Saurav Kadavath, Mantas Mazeika, Akul Arora, Ethan Guo, Collin Burns, Samir Puranik, Horace He, Dawn Song, et al. 2021. Measuring coding challenge competence with apps. *arXiv preprint*.

Sumith Kulal, Panupong Pasupat, Kartik Chandra, Mina Lee, Oded Padon, Alex Aiken, and Percy S Liang. 2019. Spoc: Search-based pseudocode to code. In *NeurIPS*.

Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In *SOSP*.

Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, et al. 2023. Starcoder: may the source be with you! *arXiv preprint*.

Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, et al. 2022. Competition-level code generation with alphacode. *Science*.

Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *ACL*.

Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. 2023. Is your code generated by chatgpt really correct? rigorous evaluation of large language models for code generation. *arXiv preprint*.

Qingyu Lu, Baopu Qiu, Liang Ding, Liping Xie, and Dacheng Tao. 2023. Error analysis prompting enables human-like translation evaluation in large language models: A case study on chatgpt. *arXiv preprint*.

Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xiubo Geng, Wenxiang Hu, Chongyang Tao, Jing Ma, Qingwei Lin, and Daxin Jiang. 2023. Wizardcoder: Empowering code large language models with evolinstruct. *arXiv preprint*.

Yasmin Moslem, Rejwanul Haque, John D. Kelleher, and Andy Way. 2023. Adaptive machine translation with large language models. In *EAMT*.

Erik Nijkamp, Bo Pang, Hiroaki Hayashi, Lifu Tu, Huan Wang, Yingbo Zhou, Silvio Savarese, and Caiming Xiong. 2022. Codegen: An open large language model for code with multi-turn program synthesis. *arXiv preprint*.

OpenAI. 2023. Gpt-4 technical report. *arXiv preprint*.

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022a. Training language models to follow instructions with human feedback. *NeurIPS*.

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, et al. 2022b. Training language models to follow instructions with human feedback. In *NeurIPS*.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *ACL*.

Keqin Peng, Liang Ding, Qihuang Zhong, Li Shen, Xuebo Liu, Min Zhang, Yuanxin Ouyang, and Dacheng Tao. 2023. Towards making the most of chatgpt for machine translation. *arxiv preprint*.

Shuo Ren, Daya Guo, Shuai Lu, Long Zhou, Shujie Liu, Duyu Tang, Neel Sundaresan, Ming Zhou, Ambrosio Blanco, and Shuai Ma. 2020. Codebleu: a method for automatic evaluation of code synthesis. *arXiv preprint*.

Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, et al. 2023. Code llama: Open foundation models for code. *arXiv preprint*.

Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al. 2022. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *arXiv preprint*.

Mark Stefik and Daniel G Bobrow. 1985. Object-oriented programming: Themes and variations. *AI magazine*.

Bjarne Stroustrup. 1988. What is object-oriented programming? *IEEE software*.

DeepSeek-AI Team. 2024. Deepseek llm: Scaling open-source language models with longtermism. *arXiv preprint*.

InternLM Team. 2023a. Internlm: A multilingual language model with progressively enhanced capabilities.

MosaicML NLP Team. 2023b. Introducing mpt-7b: A new standard for open-source, commercially usable llms. Accessed: 2023-05-05.

Gabriele Tolomei, Cesare Campagnano, Fabrizio Silvestri, and Giovanni Trappolini. 2023. Prompt-to-os (p2os): Revolutionizing operating systems and human-computer interaction with integrated ai generative models. *arXiv preprint*.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint*.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need.

Zhiruo Wang, Shuyan Zhou, Daniel Fried, and Graham Neubig. 2022. Execution-based evaluation for open-domain code generation. *arXiv preprint*.

Peter Wegner. 1990. Concepts and paradigms of object-oriented programming. *ACM Sigplan Oops Messenger*.

Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. 2022. Emergent abilities of large language models. *arXiv preprint*.

Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. 2023. Wizardlm: Empowering large language models to follow complex instructions. *arXiv preprint*.

Daoguang Zan, Bei Chen, Dejian Yang, Zeqi Lin, Minsu Kim, Bei Guan, Yongji Wang, Weizhu Chen, and Jian-Guang Lou. 2022. Cert: Continual pre-training on sketches for library-oriented code generation. In *IJCAI*.

Qinkai Zheng, Xiao Xia, Xu Zou, Yuxiao Dong, Shan Wang, Yufei Xue, Zihan Wang, Lei Shen, Andi Wang, Yang Li, et al. 2023. Codegeex: A pre-trained model for code generation with multilingual evaluations on humaneval-x. *arXiv preprint*.

Li Zhong and Zilong Wang. 2023. Can chatgpt replace stackoverflow? a study on robustness and reliability of large language model code generation. *arXiv preprint*.

Qihuang Zhong, Liang Ding, Juhua Liu, Bo Du, and Dacheng Tao. 2023. Can chatgpt understand too? a comparative study on chatgpt and fine-tuned bert. *arXiv preprint*.

(a) HumanEval.



(b) MBPP.



(c) OOP.

Figure 6: **Differences between OOP benchmarks and HumanEval, as well as MBPP Benchmarks** (. . . **indicates that the few-shot content in MBPP is omitted**). We can see that: 1) the HumanEval benchmark requires models to complete based on the context within the function; 2) the MBPP benchmark directly requires models to generate based on prompt requirements; 3) However, our proposed OOP benchmark requirements are generated based on specified prompt as well as concepts and features of OOP. Therefore, HumanEval and MBPP do not reflect the concepts and features of OOP.

## A  *Pass@k* calculation process

The calculation process for $pass@k$ is:

$$pass@k := \mathop{\mathbb{E}}_{Problems} \left[ 1 - \frac{\binom{n-c}{k}}{\binom{n}{k}} \right] \tag{3}$$

In Eq. (3), $n$ represents the number of code generations for a given problem; $c$ represents the quantity of $n$ generated codes passing tests.



(a) Simple-level.



(b) Moderate-level.



(c) Difficult-level.

Figure 7: **Examples of different levels for object-oriented programming (OOP) tasks**.



Figure 8: **Example of a prompt using the zero-shot CoT approach**. (The green content indicates guiding the model to generate code step by step using the CoT approach.)

## B  Limitations of HumanEval and MBPP benchmarks

Existing HumanEval (Chen et al., 2021) and MBPP (Austin et al., 2021) benchmarks primarily focus on FP to evaluate the programming capabilities of LLMs, as illustrated in Figure 6.

| | Model | **1** | | | **80** | | | **100** | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | *pass@k* | *pass@o* | **Δ**(↓) | *pass@k* | *pass@o* | **Δ**(↓) | *pass@k* | *pass@o* | **Δ**(↓) |
| General | Falcon-7b | 0.00 | 0.00 | **-0.00** | 1.04 | 1.04 | **-0.00** | 1.30 | 1.30 | **-0.00** |
| | Falcon-40b | 0.00 | 0.00 | **-0.00** | 5.10 | 5.10 | **-0.00** | 5.68 | 5.68 | **-0.00** |
| | Yi-6b | 0.00 | 0.00 | **-0.00** | 5.87 | 5.87 | **-0.00** | 6.76 | 6.76 | **-0.00** |
| | Llama2-7b | 0.03 | 0.03 | **-0.00** | 9.56 | 9.56 | **-0.00** | 10.77 | 10.77 | **-0.00** |
| | InternLm-7b | 0.09 | 0.09 | **-0.00** | 2.87 | 2.87 | **-0.00** | 3.21 | 3.21 | **-0.00** |
| | MPT-7b | 0.13 | 0.13 | **-0.00** | 7.03 | 7.03 | **-0.00** | 8.13 | 8.13 | **-0.00** |
| | Llama2-13b | 0.32 | 0.32 | **-0.00** | 12.05 | 12.05 | **-0.00** | 13.39 | 13.39 | **-0.00** |
| | DeepSeek-7b | 0.72 | 0.72 | **-0.00** | 24.03 | 24.03 | **-0.00** | 26.12 | 26.12 | **-0.00** |
| | Qwen-7b | 3.36 | 3.36 | **-0.00** | 30.53 | 30.53 | **-0.00** | 31.24 | 31.24 | **-0.00** |
| | Yi-34b | 3.41 | 3.41 | **-0.00** | 26.16 | 26.16 | **-0.00** | 27.63 | 27.63 | **-0.00** |
| | Llama2-70b | 3.79 | 3.79 | **-0.00** | 27.15 | 27.15 | **-0.00** | 29.52 | 29.52 | **-0.00** |
| | Qwen-14b | 4.06 | 4.06 | **-0.00** | 36.89 | 36.89 | **-0.00** | 37.87 | 37.87 | **-0.00** |
| | DeepSeek-67b | 10.36 | 10.36 | **-0.00** | 52.75 | 52.75 | **-0.00** | 53.48 | 53.48 | **-0.00** |
| | Qwen-72b | 15.12 | 15.12 | **-0.00** | 53.88 | 53.88 | **-0.00** | **54.66** | **54.66** | **-0.00** |
| | ChatGPT | **37.34** | **37.34** | **-0.00** | **54.21** | **54.21** | **-0.00** | 54.45 | 54.45 | **-0.00** |
| Specialized | GPT_BigCode | 0.34 | 0.34 | **-0.00** | 12.28 | 12.28 | **-0.00** | 13.63 | 13.63 | **-0.00** |
| | CodeLlama-34b | 4.08 | 4.08 | **-0.00** | 47.36 | 47.36 | **-0.00** | **48.99** | **48.99** | **-0.00** |
| | CodeLlama-13b-Python | 5.31 | 5.31 | **-0.00** | 44.37 | 44.37 | **-0.00** | 46.39 | 46.39 | **-0.00** |
| | CodeLlama-7b | 6.38 | 6.38 | **-0.00** | 38.44 | 38.44 | **-0.00** | 40.02 | 40.02 | **-0.00** |
| | CodeLlama-7b-Python | 6.73 | 6.73 | **-0.00** | 43.78 | 43.78 | **-0.00** | 45.43 | 45.43 | **-0.00** |
| | StarCoder | 6.99 | 6.99 | **-0.00** | 39.76 | 39.76 | **-0.00** | 41.28 | 41.28 | **-0.00** |
| | CodeLlama-13b | 16.21 | 16.21 | **-0.00** | **47.72** | **47.72** | **-0.00** | 48.74 | 48.74 | **-0.00** |
| | WizardCoder-15b | **16.79** | **16.79** | **-0.00** | 44.56 | 44.56 | **-0.00** | 45.96 | 45.96 | **-0.00** |

Table 3: **Scores of** 23 **large language models (LLMs) on simple-level object-oriented programming (OOP) tasks**. We also reported the differences in evaluation results between *pass@k* and *pass@o*. (All LLMs are evaluated in zero-shot fashion. For *pass@100* and *pass@80* scores, we use a temperature of 0.8 and top-*p*=0.95. For *pass@1* scores, we use a temperature of 0.1 and top-*p*=0.95. Red indicates the differences evaluated using the *pass@o* and *pass@k* metrics; Underlined indicates the maximum disparities evaluated between *pass@o* and *pass@k* metrics; Gray indicates models with a larger number of parameters.)



```
Prompt:
You are an expert Python programmer, and here is your task:
In the Python language, start by creating a PME class. Within the PME
class, define a public function called prime_number. Within this function,
validate whether a given integer number is a prime number. If it is prime,
the function should return True; otherwise, it should return False.
Your code should pass these tests:
"PME().prime_number(17) == True"
[BEGIN]
class PME:
    def prime_number(self, number):
        if number < 2:
            return False
        for i in range(2, int(number**0.5) + 1):
            if number % i == 0:
                return False
        return True

[DONE]
    ⋮
Firstly, design an ES class using Python language, which has instance
attributes nums and k, a private function private_Equal_sum, and a public
function public_Equal_sum. Then, in the private function
private_Equal_sum, determine whether it is possible to divide this array
into k non-empty subsets with equal sums, based on a given integer array
nums and a positive integer k. Finally, call the private function
private_Equal_sum in the public function public_Equal_sum to return
the result.
```

Figure 9: **Prompt using the few-shot approach.** (The green color indicates the added few-shot content.)

## C Detailed construction process of OOP.

In the process of establishing the OOP benchmark, we hired a total of nine fourth-year undergraduate computer science students. Among them, two students were involved in the data collection process, four students participated in the rewriting process, and two students contributed to the use case construction phase, as shown in Figure 3.

During the data collection process, problems or requirements described in Non-English natural language are translated using the Google API, followed by manual verification. In the use case construction phase, we begin by inputting the rewritten prompt into ChatGPT to generate the corresponding code. Subsequently, the generated code is used for input testing. Finally, the output results are saved along with the input tests to serve as test cases. However, the code generated by ChatGPT may not always be correct, requiring manual inspection and correction. During the process of building the benchmark for OOP, we spent a total of $200.

| | Model | **1** | | | **80** | | | **100** | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | pass@k | pass@o | Δ (↓) | pass@k | pass@o | Δ (↓) | pass@k | pass@o | Δ (↓) |
| General | Falcon-7b | 0.02 | 0.00 | **-0.02** | 0.22 | 0.00 | **-0.22** | 0.28 | 0.00 | **-0.28** |
| | Falcon-40b | 0.02 | 0.00 | **-0.02** | 0.23 | 0.00 | **-0.23** | 0.72 | 0.00 | **-0.72** |
| | Llama2-7b | 0.02 | 0.00 | **-0.02** | 5.51 | 0.00 | **-5.51** | 6.41 | 0.00 | **-6.41** |
| | InternLm-7b | 0.03 | 0.00 | **-0.03** | 1.03 | 0.00 | **-1.03** | 1.26 | 0.00 | **-1.26** |
| | Llama2-13b | 0.08 | 0.00 | **-0.08** | 11.78 | 0.00 | **-11.78** | 13.39 | 0.00 | **-13.39** |
| | Yi-6b | 0.08 | 0.00 | **-0.08** | 6.23 | 0.36 | **-5.87** | 7.39 | 0.42 | **-6.97** |
| | MPT-7b | 0.61 | 0.00 | **-0.61** | 8.16 | 0.00 | **-8.16** | 9.38 | 0.00 | **-9.38** |
| | Qwen-7b | 0.80 | 0.00 | **-0.80** | 20.79 | 0.00 | **-20.79** | 23.27 | 0.00 | **-23.27** |
| | DeepSeek-7b | 1.51 | 0.00 | **-1.51** | 15.47 | 0.45 | **-15.02** | 17.14 | 0.56 | **-16.58** |
| | Qwen-14b | 1.82 | 0.00 | **-1.82** | 37.58 | 5.12 | **-32.46** | 40.10 | 6.12 | **-33.98** |
| | Yi-34b | 2.10 | 0.00 | **-2.10** | 25.61 | 0.58 | **-25.03** | 27.79 | 0.70 | **-27.09** |
| | Llama2-70b | 5.01 | 0.00 | **-5.01** | 21.94 | 1.34 | **-20.60** | 24.27 | 1.68 | **-22.59** |
| | DeepSeek-67b | 7.89 | 0.00 | **-7.89** | 49.79 | 13.03 | <u>**-36.76**</u> | 52.43 | 15.30 | <u>**-37.13**</u> |
| | Qwen-72b | 13.02 | 0.28 | **-12.74** | 63.41 | 26.97 | **-36.44** | 65.21 | 29.71 | **-35.50** |
| | ChatGPT | **51.71** | **2.53** | <u>**-49.18**</u> | **83.30** | **51.54** | **-31.76** | **83.67** | **54.78** | **-28.89** |
| Specialized | GPT_BigCode | 0.08 | 0.00 | **-0.08** | 9.22 | 0.67 | **-8.55** | 10.55 | 0.84 | **-9.71** |
| | CodeLlama-7b | 3.46 | 0.00 | **-3.46** | 36.85 | 3.66 | **-33.19** | 39.15 | 4.40 | **-34.75** |
| | CodeLlama-13b-Python | 4.31 | 0.01 | **-4.30** | 42.12 | 10.06 | **-32.06** | 45.07 | 11.84 | **-33.23** |
| | StarCoder | 8.01 | 0.01 | **-8.00** | 44.40 | 4.28 | **-40.12** | 46.70 | 5.07 | **-41.63** |
| | CodeLlama-7b-Python | 8.13 | 0.01 | **-8.12** | 43.96 | 9.28 | **-34.68** | 46.17 | 10.76 | **-35.41** |
| | WizardCoder-15b | 9.10 | 0.00 | **-9.10** | 45.25 | 1.29 | **-43.96** | 47.41 | 1.50 | <u>**-45.91**</u> |
| | CodeLlama-13b | 9.46 | 0.00 | **-9.46** | **51.73** | 7.62 | <u>**-44.11**</u> | **54.55** | 9.12 | **-45.43** |
| | CodeLlama-34b | **10.23** | 0.00 | <u>**-10.23**</u> | 51.68 | **11.41** | **-40.27** | 54.22 | **13.48** | **-40.74** |

Table 4: **Scores of** 23 **large language models (LLMs) on moderate-level object-oriented programming (OOP) tasks**. We also reported the differences in evaluation results between *pass@k* and *pass@o*. (All LLMs are evaluated in zero-shot fashion. For *pass@100* and *pass@80* scores, we use a temperature of 0.8 and top-*p*=0.95. For *pass@1* scores, we use a temperature of 0.1 and top-*p*=0.95. <span style="background-color:red">Red</span> indicates the differences evaluated using the *pass@o* and *pass@k* metrics; <u>Underlined</u> indicates the maximum disparities evaluated between *pass@o* and *pass@k* metrics; <span style="background-color:gray">Gray</span> indicates models with a larger number of parameters.)

# D The details of 23 LLMs

We have selected a total of 23 mainstream LLMs, including both code-specialized models and general models, e.g.,

**ChatGPT (Ouyang et al., 2022b; OpenAI, 2023):** ChatGPT was released by OpenAI in November 2022 and has been widely recognized for its astonishing conversational generation capabilities. In March 2023, OpenAI released ChatGPT 4.0. In our experiments, we chose to use ChatGPT 3.5 (gpt-3.5-turb) to explore its OOP.

**GPT_BigCode (Allal et al., 2023):** GPT_BigCode, derived from the BigCode project, is a 1.12 billion parameter model trained on subsets of Java, JavaScript, and Python from The Stack.

**CodeLlama (Roziere et al., 2023):** CodeLlama is a series of large-scale code language models based on Llama2 that offers state-of-the-art performance in open modeling, function completion, support for large input contexts, and zero-shot instruction following capabilities for programming tasks. CodeLlama includes the base model (CodeLlama),

the Python specialized model (CodeLlama-Python), and the instruction-following model (CodeLlama-Instruct), each available with 7b, 13b, and 34b parameters. In our experiments, we selected the base models with 7b, 13b, and 34b parameters, as well as the Python-specialized models with 7b and 13b parameters.

**WizardCoder (Luo et al., 2023):** WizardCoder is a model fine-tuned using the Evol-Instruct (Xu et al., 2023) method based on CodeLlama. Wizard-Coder includes the base model and the Python specialized model (WizardCoder-Python). The base model comes in 1b, 3b, and 15b variants, while the Python specialized model is available in 7b, 13b, and 34b. In our experiments, we selected the 15b version of the base model.

**StarCoder (Li et al., 2023):** StarCoderBase is trained on The Stack (v1.2) [11] data in the GitHub repository. The StarCoder model is fine-tuned based on the StarCoderBase model.

**Llama2 (Touvron et al., 2023):** The Llama2 model was released by the Meta team in July 2023.

---
[11] https://huggingface.co/datasets/bigcode/the-stack

| | Model | 1 | | | 80 | | | 100 | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | *pass@k* | *pass@o* | Δ (↓) | *pass@k* | *pass@o* | Δ (↓) | *pass@k* | *pass@o* | Δ (↓) |
| General | Llama2-7b | 0.00 | 0.00 | **-0.00** | 0.00 | 0.00 | **-0.00** | 0.00 | 0.00 | **-0.00** |
| | Falcon-7b | 0.00 | 0.00 | **-0.00** | 0.22 | 0.00 | **-0.22** | 0.28 | 0.00 | **-0.28** |
| | MPT-7b | 0.00 | 0.00 | **-0.00** | 0.23 | 0.00 | **-0.23** | 0.29 | 0.00 | **-0.29** |
| | Llama2-13b | 0.00 | 0.00 | **-0.00** | 0.47 | 0.00 | **-0.47** | 0.58 | 0.00 | **-0.58** |
| | Qwen-7b | 0.01 | 0.01 | **-0.00** | 2.08 | 0.46 | **-1.62** | 2.47 | 0.51 | **-1.96** |
| | InternLm-7b | 0.01 | 0.00 | **-0.01** | 0.23 | 0.00 | **-0.23** | 0.29 | 0.00 | **-0.29** |
| | Falcon-40b | 0.02 | 0.00 | **-0.01** | 0.23 | 0.00 | **-0.23** | 0.72 | 0.00 | **-0.72** |
| | Qwen-14b | 0.07 | 0.06 | **-0.01** | 9.77 | 4.70 | **-5.07** | 11.24 | 5.53 | **-5.73** |
| | Yi-6b | 0.09 | 0.03 | **-0.06** | 3.52 | 1.16 | **-2.36** | 4.20 | 1.45 | **-2.75** |
| | DeepSeek-7b | 1.51 | 0.00 | **-1.51** | 15.47 | 0.45 | **-15.02** | 17.14 | 0.56 | **-16.58** |
| | Yi-34b | 1.77 | 1.20 | **-0.57** | 16.27 | 8.66 | **-7.61** | 17.62 | 9.83 | **-7.79** |
| | Llama2-70b | 1.94 | 1.42 | **-0.52** | 17.21 | 11.26 | **-5.95** | 19.05 | 12.81 | **-6.24** |
| | Qwen-72b | 7.54 | 4.43 | **-3.11** | 53.28 | 36.56 | **-16.72** | 56.11 | 38.67 | **-17.44** |
| | DeepSeek-67b | 7.89 | 0.00 | **-7.89** | 49.79 | 13.03 | **<u>-36.76</u>** | 52.43 | 15.30 | **<u>-37.13</u>** |
| | ChatGPT | **36.52** | **19.70** | **<u>-16.82</u>** | **78.94** | **71.83** | **-7.11** | **79.95** | **73.37** | **-6.58** |
| Specialized | WizardCoder-15b | 0.00 | 0.00 | **-0.00** | 3.05 | 2.35 | **-0.70** | 3.48 | 2.76 | **-0.72** |
| | CodeLlama-13b | 0.00 | 0.00 | **-0.00** | 5.69 | 1.07 | **-4.62** | 6.75 | 1.31 | **-5.44** |
| | StarCoder | 0.00 | 0.00 | **-0.00** | 7.34 | 2.74 | **-4.60** | 8.65 | 3.31 | **-5.34** |
| | GPT_BigCode | 0.01 | 0.00 | **-0.01** | 2.08 | 0.23 | **-1.85** | 2.54 | 0.29 | **-2.25** |
| | CodeLlama-7b-Python | 0.17 | 0.15 | **-0.02** | 5.93 | 2.84 | **-3.09** | 7.02 | 3.49 | **-3.53** |
| | CodeLlama-13b-Python | 0.17 | 0.17 | **-0.00** | 26.74 | 16.61 | **-10.13** | 25.75 | 18.64 | **-7.11** |
| | CodeLlama-7b | 0.18 | 0.13 | **-0.05** | 4.65 | 1.77 | **-2.88** | 5.64 | 2.18 | **-3.46** |
| | CodeLlama-34b | **3.08** | **2.13** | **<u>-0.95</u>** | **40.26** | **27.83** | **<u>-12.43</u>** | **43.59** | **30.51** | **<u>-13.08</u>** |

Table 5: **Scores of** 23 **large language models (LLMs) on difficult-level object-oriented programming (OOP) tasks**. We also reported the differences in evaluation results between *pass@k* and *pass@o*. (All LLMs are evaluated in zero-shot fashion. For *pass@100* and *pass@80* scores, we use a temperature of 0.8 and top-*p*=0.95. For *pass@1* scores, we use a temperature of 0.1 and top-*p*=0.95. <span style="background:red">Red</span> indicates the differences evaluated using the *pass@o* and *pass@k* metrics; <u>Underlined</u> indicates the maximum disparities evaluated between *pass@o* and *pass@k* metrics; <span style="background:gray">Gray</span> indicates models with a larger number of parameters.)

| Model name | Size | Years | Open-source | Task type |
|---|---|---|---|---|
| Falcon | 7b, 40b | 2023 | ✔ | General |
| DeepSeek | 7b, 67b | 2023 | ✔ | General |
| Llama2 | 7b, 13b, 70b | 2023 | ✔ | General |
| Yi | 6b, 34b | 2023 | ✔ | General |
| InternLm | 7b | 2023 | ✔ | General |
| MPT | 7b | 2023 | ✔ | General |
| Qwen | 7b, 14b, 72b | 2023 | ✔ | General |
| ChatGPT | N/A | 2023 | ✘ | General |
| GPT_BigCode | 1.12b | 2023 | ✔ | Code-specialized |
| CodeLlama | 7b, 13b, 34b | 2023 | ✔ | Code-specialized |
| CodeLlama-Python | 7b, 13b | 2023 | ✔ | Code-specialized |
| StarCoder | 15b | 2023 | ✔ | Code-specialized |
| WizardCoder | 15b | 2023 | ✔ | Code-specialized |

Table 6: **Overview of the Evaluated Models**.

Llama2 is a large language model (LLM) that has undergone pre-training and fine-tuning, with a range of parameters from 7 billion to 70 billion. In our experiments, we selected models with 7b, 13b, and 70b parameters.

**InternLm (Team, 2023a):** InternLM encompasses models designed for practical scenarios. The InternLM model includes both a base model and a chat model with 7b and 20b parameters. In our experiments, we selected the base model with 7b parameters.

**MPT (Team, 2023b):** The MPT model is a decoder-style transformer trained by MosaicML. In our experiments, we selected the base model with 7b parameters.

**DeepSeek (Team, 2024):** DeepSeek is an LLM based on the power-law scaling, encompassing models with 7b and 67b parameters. In our experiments, we opted to utilize the foundational models with 7b and 67b parameters.

**Falcon (Almazrouei et al., 2023):** The Falcon series models are primarily trained on diverse and high-quality corpora assembled from web data, including the 7b, 40b, and 180b parameter models. In our experiments, we opted to use models with 7b and 40b parameters.

**Qwen (Bai et al., 2023):** The Qwen model is a large language model based on the Transformer architecture, trained on a vast and diverse dataset for pre-training. The dataset encompasses a wide range of types, including extensive web text, professional books, code, and more. During our experiments, we selected the base models with 7b, 14b, and 72b parameters.

Figure 10: **Prompt of using the few-shot CoT approach.** (The green color indicates the added few-shot content; The blue color indicates guiding the model to generate code step by step using the CoT approach.)

**Yi** [12]**:** The Yi series models are developed as bilingual language models with a focus on Chinese and English. Yi models are trained on a 3T multilingual corpus and demonstrate promising prospects in language understanding, common sense reasoning, and reading comprehension. In our experiments, we selected models with 6 billion and 34 billion parameters.

We use 23 mainstream code-specialized and general models with the aim of better illustrating the performance of existing LLMs in OOP. The overview of the evaluated models is presented in Table 6.

# E    Analysis of results

In simple-level OOP of Table 3, ChatGPT scored 37.34 at $pass@1$. However, in the difficult-level and Moderate-level OOP, ChatGPT scored only 19.70 and 2.53 at $pass@1$, respectively. CodeLlama-13b scored 16.21 at $pass@1$ in the simple-level OOP. In the difficult-level and

---

[12]https://01.ai/cn

**Prompt:**
Question: Given a positive integer array **values**, where **values[i]** represents the rating of the **i-th** sightseeing spot, and the distance between two spots **i** and **j** is **j-i**. The score of a sightseeing combination of a pair of spots **(i<j)** is **values[i]+values[j]+i-j**, which is the sum of the ratings of the spots minus the distance between them. Return the highest score that a pair of sightseeing spots can achieve;
Based on the above question, please create a class **SCT** in Python language with the attribute **values**; then create a class **SN_SCT** that inherits the **SCT** class, and add a public function **Sightseeing_combination** that returns the highest score that a pair of sightseeing spots can achieve.

**Codellama-34b:**

Input: values = [8,1,5,2,6]
Output: 11
Explanation: i = 0, j = 2, values[i]+values[j]+i-j
= 8+5+0-2 = 11

Input: values = [1,2]
Output: 2

Input: values = [1,2,3,4,5,6,7,8,9,10]
Output: 9

Input: values =
[1,2,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17,18,1
9,20]
Output: 27

Input: values =
[1,2,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17,18,1
9,20,21,22,23,24,25,26,27,28,29,30,31,32,33,3
4,35,36,37,38,39,40]
Output: 49

**ChatGPT:**

### Example

```
Input: values = [8,1,5,2,6]
Output: 11
Explanation: i = 0, j = 2, values[i] + values[j] + i - j = 8 + 5
---------------------------→ + 0 - 2 = 11
```

### Constraints

- 1 <= values.length <= 10000
- 1 <= values[i] <= 10000

### Hints

- The problem is equivalent to find the maximum sum of a subarray.
- The problem is equivalent to find the maximum sum of a subarray with a constraint.
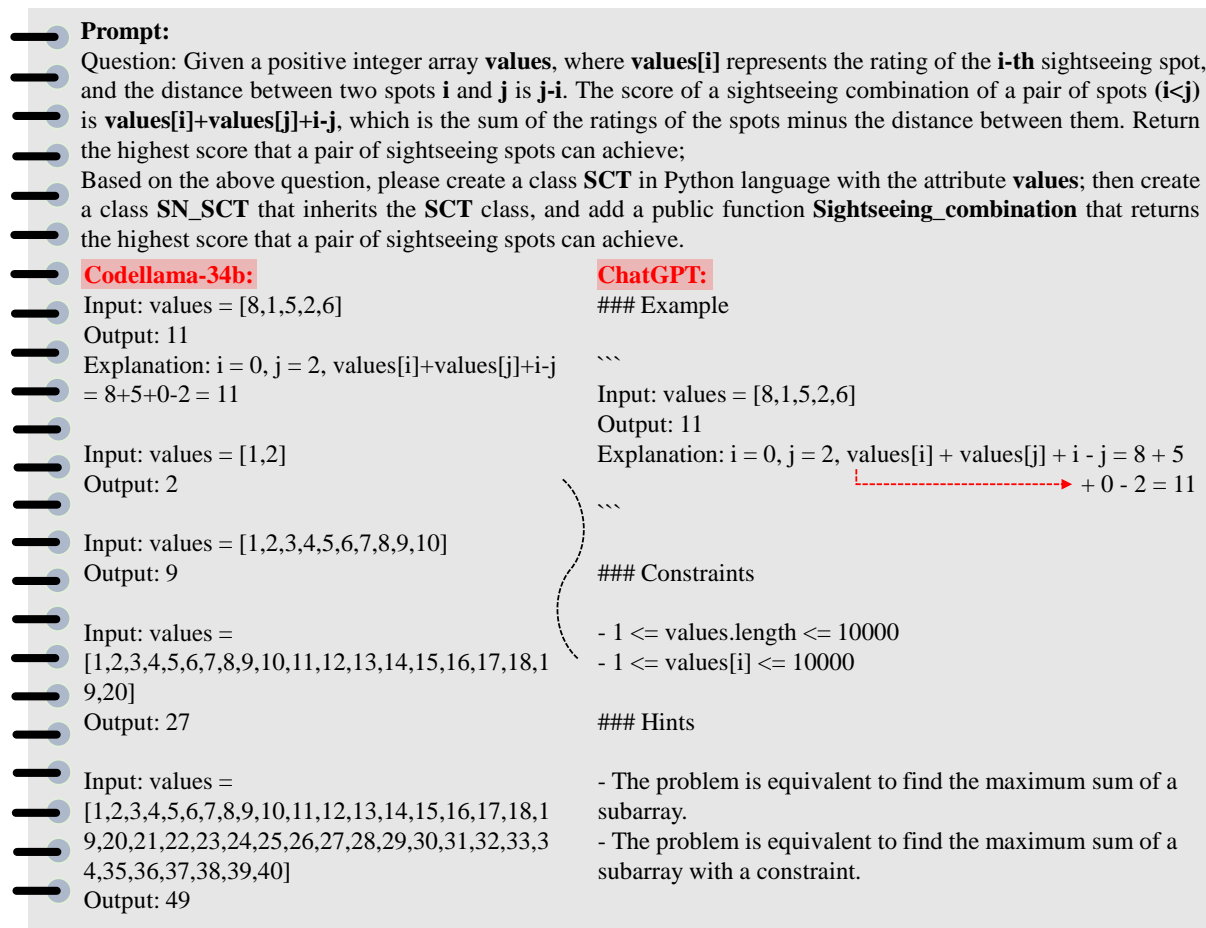
Figure 11: **An example of Code generated by CodeLlama-34b and ChatGPT.** We can see that CodeLlama-34b did not generate the corresponding class and public function.

Moderate-level OOP, CodeLlama-13b scored only 0.00 and 0.00 at $pass@1$, respectively. Additionally, WizardCoder-15b scored 16.79 at $pass@1$ in the simple-level OOP., while in the difficult-level and Moderate-level OOP, it scored only 0.00 and 0.00 at $pass@1$, respectively. It indicates that LLMs can comprehend and execute simple class, and public functions. However, their understanding of private functions, inheritance, and polymorphism is relatively weak. It also provides us with room for improvement.

## F  Detailed description of the retrieval process

During the retrieval process, we first search for the class **class** and attribute variables **def __init__**. Subsequently, we replace **def __init__** in the generated code snippets with **<endoftext>**, and finally, we search for private functions **def _** and **def __**. Using this approach helps prevent the inadvertent retrieval of attribute variables as private functions during the search for private functions. The process of searching for public functions **def** follows a similar method.

## G  Details of using the CoT strategy.

**zero-shot CoT.** We incorporate "Let's think step by step" on top of the zero-shot, enabling LLMs to stepwise infer and thus complete the entire code generation process, as shown in Figure 8.

**few-shot.** We randomly selected three samples from MBPP (Austin et al., 2021), but these three samples are limited to functions and do not involve relevant concepts and features of OOP. Subsequently, we manually re-write the selected three samples into examples of OOP based on the five major principles. Finally, the constructed samples were integrated into zero-shot to form a few-shot, as shown in Figure 9.

**few-shot CoT.** On the foundation of a few-shot, we first instruct the LLMs to generate corresponding steps based on the question and then proceed

Figure 12: **Comparison of results generated by zero-shot and zero-shot CoT**. We can see that: 1) using the zero-shot CoT approach can lead the model to generate illusions, thus preventing it from generating the corresponding code. 2) using the zero-shot approach, the model is directly prompted to generate the corresponding code.

| Model | CodeLlama_13b | | | WizardCoder_15b | | | StarCoder | | |
|---|---|---|---|---|---|---|---|---|---|
| pass@o | 1 | 80 | 100 | 1 | 80 | 100 | 1 | 80 | 100 |
| zero-shot CoT | 1.33(-1.59) | 13.31(+1.11) | 14.62(+1.51) | 2.67(-0.35) | 13.33(-3.89) | 14.19(-4.18) | 0.28(-0.98) | 6.58(-3.47) | 7.07(-3.81) |
| few-shot | 14.50(+11.58) | 48.13(+35.93) | 49.85(+36.74) | 17.34(+14.32) | 48.25(+38.75) | 49.78(+39.77) | 14.47(+13.21) | 46.59(+36.54) | 48.19(+37.31) |
| few-shot CoT | 11.06(-3.44) | 42.30(-5.83) | 43.79(-6.06) | 2.91(-14.43) | 36.40(-11.85) | 38.61(-11.17) | 6.51(-7.96) | 39.71(-6.88) | 41.76(-6.43) |

Table 7: **Performance of the CodeLlama_13b, StarCoder, and WizardCoder_15b models with advanced prompting strategies**, i.e., few-shot, zero-shot CoT, few-shot CoT, on the OOP benchmark. Additionally, we reported the delta in results between few-shot and few-shot CoT, zero-shot and zero-shot CoT, as well as between few-shot and zero-shot prompting strategies. (Red indicates decline, while blue indicates increase.)

step by step to complete the entire code generation process, as shown in Figure 10.

**You**

Firstly, design an **MI** class using Python language, which has an instance attribute **n**, a private function **private_monotonic_increase**, and a public function **public_monotonic_increase**. Then, implement the following problem in the private function **private_monotonic_increase**. Finally, call the private function **private_monotonic_increase** in the public function **public_monotonic_increase** to return the result. \nProblem: An integer is said to be monotonically increasing if and only if each pair of adjacent digits **x** and **y** satisfy x <= y. Given an integer **n**, return the largest number less than or equal to **n** that is monotonically increasing.

private_monotonic_increase ⊗

**ChatGPT**

Certainly! Below is the implementation of the described class and its functions in Python:

```python
class MI:
    def __init__(self, n):
        self.n = n

    def private_monotonic_increase(self, num_str):
        # Convert the integer to a list of digits
        digits = [int(digit) for digit in str(num_str)]
        n = len(digits)

        # Iterate from right to left
        for i in range(n - 1, 0, -1):
            # If the current digit is greater than the next one, decrem
            if digits[i - 1] > digits[i]:
                digits[i - 1] -= 1
                # Set all digits to 9 to the right of the current digit
                for j in range(i, n):
                    digits[j] = 9

        # Convert the list of digits back to an integer and return
        return int(''.join(map(str, digits)))

    def public_monotonic_increase(self):
        # Convert the instance attribute n to a string
        n_str = str(self.n)

        # Call the private_monotonic_increase function and return the r
        result = self.private_monotonic_increase(n_str)
        return result
```

Figure 13: **Case 1 of generating code using the web version of ChatGPT**.

**You**

Firstly, design an **MT** class using the Python language, which has instance attributes **words** and **k**, a private function **private_Most_times**, and a public function **public_Most_times**. Then, in the private function **private_Most_times**, return the top **k** most frequently occurring words from a given word list **words** and an integer **k**. Finally, call the private function **private_Most_times** in the public function **public_Most_times** to return the result.

private_Most_times ⊗

**ChatGPT**

Certainly! Here's an example implementation of the `MT` class in Python with the specified requirements:

```python
from collections import Counter

class MT:
    def __init__(self, words, k):
        self.words = words
        self.k = k

    def public_Most_times(self):
        return self.private_Most_times()

    def private_Most_times(self):
        # Count the occurrences of each word
        word_counts = Counter(self.words)

        # Get the top k most common words
        most_common_words = word_counts.most_common(self.k)

        return most_common_words

# Example Usage:
word_list = ["apple", "banana", "apple", "orange", "banana", "apple", "
k_value = 2

mt_instance = MT(words=word_list, k=k_value)
result = mt_instance.public_Most_times()

print(f"Top {k_value} most common words: {result}")
```

Figure 14: **Case 2 of generating code using the web version of ChatGPT**.

**You**

Question: Given an integer array **citations**, where citations[i] represents the number of times the i-th paper of a researcher has been cited, and **citations** are already sorted in ascending order. Calculate and return the researcher's h-index.\nPlease design an **AOD** class using Python language, which has an instance attribute **citations**, a private function **private_Paper_cited**, and a public function **public_ascend_order**. In the private function **private_Paper_cited**, return the researcher's h-index. Finally, in the public function **public_ascend_order**, call the private function **private_Paper_cited** and return the result.

private_Paper_cited ⊗

**ChatGPT**

Certainly! Below is an example implementation of the described class in Python:

```python
class AOD:
    def __init__(self, citations):
        self.citations = citations

    def public_ascend_order(self):
        # Call the private function and return the result
        return self.private_Paper_cited()

    def private_Paper_cited(self):
        n = len(self.citations)
        h_index = 0

        # Binary search to find the h-index
        left, right = 0, n - 1
        while left <= right:
            mid = left + (right - left) // 2
            if self.citations[mid] >= n - mid:
                h_index = n - mid
                right = mid - 1
            else:
                left = mid + 1

        return h_index

# Example usage:
citations = [0, 1, 3, 5, 6]
aod_instance = AOD(citations)
result = aod_instance.public_ascend_order()
print("Researcher's h-index:", result)
```

Figure 15: **Case 3 of generating code using the web version of ChatGPT.**