

Concerned with Data Contamination? Assessing Countermeasures in Code Language Model

Jialun Cao
The Hong Kong University of
Science and Technology
Hong Kong, China
jcaoap@cse.ust.hk

Wuqi Zhang
The Hong Kong University of
Science and Technology
Hong Kong, China
wzhangcb@cse.ust.hk

Shing-Chi Cheung
The Hong Kong University of
Science and Technology
Hong Kong, China
scc@cse.ust.hk

Abstract—Various techniques have been proposed to leverage the capabilities of code language models (CLMs) for software engineering tasks. While these techniques typically evaluate their effectiveness using publicly available datasets, the evaluation can be subject to data contamination threats where the evaluation datasets have already been used to train the concerned CLMs. This can significantly affect the reliability of the evaluation. Different countermeasures have been suggested to mitigate the data contamination threat. Countermeasures include using more recent data, curating new data, and refactoring existing data are introduced, yet it is unclear whether these countermeasures could really mitigate data contamination threats to model evaluation. To fill the gap, we systematically study to quantify the impacts of these countermeasures on CLMs’ performance.

To facilitate the study, we collected 2,493,174 (over 2 million) Python functions with timestamps ranging from January 1st, 2018, to December 31st, 2023. The data created before the models’ cut-off date are considered “contaminated data”, while the data where the countermeasures are taken are regarded as “cleansed data”. We study the impact of these countermeasures by investigating the difference in CLMs’ performance on contaminated and cleansed data derived from different countermeasures. Our experiments yield several interesting observations. For instance, CLMs do not necessarily perform worse on data after the models’ cut-off date; on the contrary, they sometimes perform better. In addition, refactoring did not always result in decreased performance; it could lead to improvements instead. Furthermore, existing metrics such as perplexity are incapable of distinguishing contaminated/cleansed data. We hope that the results and observations could help deepen the understanding of CLMs’ capabilities and inform the community about data contamination when validating the effectiveness of methods.

Index Terms—Code Language Model, Empirical Study, Code Clone, Data Contamination

I. INTRODUCTION

Large language models (LLMs) are increasingly utilized by state-of-the-art techniques to solve challenging software engineering problems [1], [2], [3], [4]. Since such models were trained on a large amount of data from various code corpus, it is hard to tell whether the evaluation results of these techniques overfit the training data. In other words, the data used for evaluation may have already been inadvertently seen in the training data, resulting in exaggerated performance. The

situation is identified as *data contamination* [5], [6], [7], [8]¹.

Challenge – Identifying the existence of data contamination is difficult. In traditional machine learning scenario where the dataset scale is not that large, the contamination could be avoid by splitting training/validation/testing sets and avoid leaking target variables during training. However, in the era of LLMs, where models are trained on vast corpora of data spanning a diverse range of topics and sources, *avoiding the data contamination becomes expensive and sometimes unrealistic*. First, *Opacity of training source*. Model developers are *unwilling to disclose* the training data and implementation details [14], [15], [16] due to various reasons (*i.e.*, intellectual property, data privacy, license, and commercial competition), leaving no clues to the training data. Second, *Indirect contamination*. The data pipelines for LLMs are much more complex. The preprocessing steps, such as tokenization and data augmentation, can introduce contamination indirectly. Third, *Scale of training data*. Even if the training source is disclosed, the sheer volume of data makes it impractical to traverse the entire training set. Moreover, simply scanning the training data may not be reliable due to the differences in data format (*e.g.*, XML, CSV) and their preprocessing (*e.g.*, tokenization, and normalization).

Research gap – Given the uncontrollability of training process of LLMs, researchers shifted to proposing **countermeasures** to ensure the reliability of evaluation. Various countermeasures are introduced. For example, ❶ Using *recent data* (*i.e.*, the data released after the deployed model’s training cut-off date) will be less subject to contamination. These data are less likely to be used for model training. ❷ The community keeps releasing new calibrated-crafted code benchmarks [17], [18], [19], [20], [21], [22] for fairer and more comprehensive evaluation of model capabilities. ❸ Various *code refactoring operations* [9], [23] (*e.g.*, renaming the identifiers or adjusting code structures) are adopted to mutate the code before the cut-off date to avoid LLMs’ memorization [9]. While these initiatives are promising, there is *no systematic, quantitative study to assess the impact of these countermeasures*.

¹Other terms like “data leakage” [9], [10], [11], [12], “memorization” [13], “task contamination” [7] are used in related works. This paper uses “data contamination” to represent uniformly.

To bridge the gap, we conducted a systematic study assessing the efficacy of three countermeasures on code language models (CLMs). We focus on CLMs rather than general large language models due to their *relevancy* to software engineering tasks. To determine whether data is contaminated or not, we use *time* as a subjective criterion. Such temporal separation is common in time-sensitive domains like financial forecasting [24], [25]. The data created before the models’ cut-off date (*i.e.*, the date when the data collection for training concluded) are considered as “*contaminated data*”, where the cut-off dates vary from models. On the contrary, the data where the *countermeasures* are adopted are regarded as “*cleansed data*”. The efficacy of countermeasures then is assessed by investigating the difference in CLMs’ performance on contaminated/cleansed data, in line with related works [5], [15].

To facilitate the study, we collected 2 million Python codes within six years (January 1st, 2018, to December 31st, 2023). The *contaminated data* is a form of the codes whose timestamps precede models’ cut-off dates. The *cleansed data* are derived from implementing countermeasures ❶ ~ ❸. Furthermore, we use three code similarities (*e.g.*, Jaccard distance, Levenshtein distance as used in related work [26]) to quantify the overlap degree between contaminated/cleansed data, showing that the *overlap is at a relatively minimal level* (Section IV-A). In addition, we compute the code similarities between contaminated and cleansed data to make sure that the overlap between them is under a relatively small level. We study three research questions (RQs):

- **RQ1. How does CLMs’ performance differ on contaminated / *recent* data** (Countermeasure ❶)? We investigate the impact of recent data (created after models’ cut-off date) as a countermeasure. We compare the CLMs’ performance on recent data compared with that on contaminated data, seeing how the performance changes over time.
- **RQ2. How does CLMs’ performance differ on contaminated / *curated* data** (Countermeasure ❷)? We also study whether curate data could serve as an effective countermeasure. In particular, we consider coding benchmarks HumanEval and CoderEval because of their popularity. Then, we analyzed the CLMs’ performance differences in contaminated and curated data.
- **RQ3. How does CLMs’ performance differ on contaminated / *refactored* data** (Countermeasure ❸)? There are several commonly adopted code refactoring operators whose impacts on CLMs’ performance are worth exploring. We then apply five operators to contaminated data and see how CLMs reacts to the refactored data.

Finally, related tasks such as membership inference attacks (MIA) [14], [27] and privacy extraction [28] introduce various metrics such as perplexity [29]. While these metrics have been proven effective in related tasks, it is still unknown whether they are effective in coding tasks and CLMs. So, we further designed RQ4 to investigate MIA-related metrics for distinguishing contaminated/cleansed data.

- **RQ4. Could existing metrics distinguish contaminated/-**

cleansed data? We investigate six related metrics on contaminated and three kinds of cleansed data, checking whether they are useful in distinguishing contaminated/cleansed data.

Findings – Our study yields several interesting findings.

- CLMs, in general, *perform better on countermeasure-applied data* (*e.g.*, recent data, curated data, and syntactically refactored data) compared with the performance on contaminated data, indicating that the current *countermeasures may not be effective* in alleviating data contamination.
- Existing MIA-related metrics, such as Perplexity, Zlib Compression Entropy, and MIN-K% PROB can *hardly distinguish* the contaminated/cleansed data.
- The popularity of AI programming assistants such as Copilot may further exacerbate data contamination in the evaluation.
- *Changing code structures* may not be useful to alleviate data contamination. On the contrary, refactoring the code structure could even upgrade the models’ performance.
- *Semantic refactoring operators* such as identifier renaming and appending special parameters have a greater impact on data and may be more useful to alleviate data contamination during evaluation.

Contributions – Our contribution is summarized as follows.

- **Novelty** We present the first study of the mitigation effect of various countermeasures for CLM data contamination. In particular, we group Python codes according to how they are collected/curated (*e.g.*, crawled from online resources, manually crafted, or refactored from other codes) and compare how CLMs perform in these code groups.
- **Significance** We decompose the study on the data contamination problem into four-fold, allowing for a more holistic analysis. Existing work either directly builds new datasets or applies automatic or manual code refactoring without systematically and quantitatively studying whether doing so can mitigate the threat of data contamination.
- **Impact.** We collected 2 million Python codes within six years (January 1st, 2018, to December 31st, 2023), then sampled 384 Python functions each year, with a sampling confidence level of 95%, culminating in 2304 Python functions. It serves not only as the foundation for our rigorous evaluation but also as a dataset for future research endeavors.

II. EXPERIMENT DESIGN

A. Experiment Design

To facilitate the study, we construct several groups of code (**code groups**) into contaminated and cleansed data. Comparing the CLM’s performance on the *contaminated* and *cleansed* data can infer the extent to which the CLMs rely on memorization versus their genuine ability and thus act as *a hint for the degree of data contamination* during the evaluation. Fig. 1 visualizes an overview design for RQ1-RQ3. The *cleansed data* (highlighted in yellow) could be collected after model release (RQ1), curated manually (RQ2), or mutated from contaminated data (RQ3). These three RQs study how CLMs’ performance changes over contaminated/cleansed data. Finally, RQ4 explores whether there are metrics that could distinguish them.

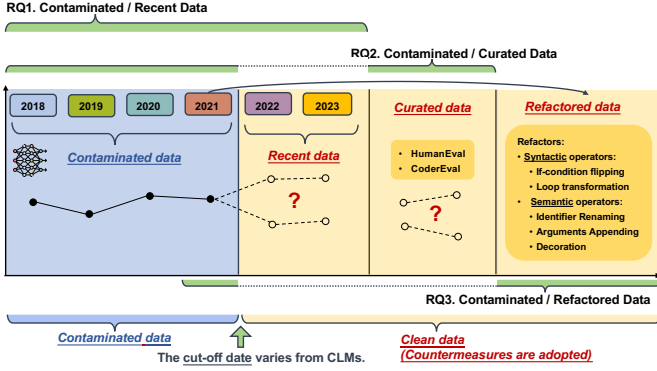


Fig. 1: **Experiment Design for RQ1-3.** We assess CLMs’ performance on contaminated and cleaned data. Three countermeasures (*i.e.*, using recent data, using curated data or refactored contaminated data) are considered.

1) **Experiment Design for RQ1:** In Fig. 1 RQ1, codes could be partitioned into distinct code groups based on their creation *times* (*e.g.*, commit time in Github) because time serves as an objective indicator [24], [25], [30]. This **chronological grouping** aims to separate the code that a CLM could have encountered during its training (*contaminated data*) from the code generated after its training was completed (*recent data*).

For clarity, we designate code within a specific year with corresponding labels, such as Code-2018 for snippets that fall in 2018 and similarly for other annual collections. Note that the pivot time of contaminated/recent data varies from CLMs (See Section III-A).

After the temporal splitting, we apply CLMs to these code groups, allowing us to systematically investigate how CLMs’ performance changes chronologically, implying the degree of data contamination within them.

2) **Experiment Design for RQ2:** Researchers and developers introduced new datasets [17], [20], [19], [21], [31], [22] to fairly and comprehensively evaluate the effectiveness of CLMs. We investigate the performance of CLMs on them.

Among the newly proposed datasets, we select two representative and popular-used function-level Python coding datasets, *i.e.*, HumanEval [17] and CoderEval [20], which we refer to as Code-HumE and Code-CodE, respectively. We focus on function-level programs to eliminate the influence of extra-long contexts [19]. In addition, since Python is the most well-trained programming language, if there are any data contamination issues, they should be more evident in Python. Also, most datasets and methods based on CLMs are designed specifically for Python. Thus, investigating data contamination threats in Python may hold broader significance than studying in other programming languages. As shown in Fig. 1 RQ2, we demonstrate the CLMs’ performance in these two datasets and compare them with that in contaminated data.

The release time of HumanEval and CoderEval are Jul 8, 2021 [32] and Jan 28, 2023 [33], respectively, according to the commit date shown in Github. We then compare these commit dates with models’ cut-off dates to determine whether these datasets contain contaminated data regarding a CLM.

3) **Experiment Design for RQ3:** As proposed by existing works [9], a more applicable way to alleviate the threat of data contamination in evaluating CLM-based approaches is to refactor the contaminated code. If CLMs maintain performance on refactored code, this would imply that they truly grasp the coding tasks rather than relying on rote memorization. In contrast, if there is a substantial fluctuation in CLMs’ performance between contaminated and refactored data, it indicates that CLMs are likely threatened by data contamination via memorizing specific code patterns or formats in the training data. Conversely, if a refactoring operator can induce significant variability (primarily decreases) in CLMs’ performance, it may mitigate the CLM’s exposure to data contamination.

As shown in Fig. 1 RQ3, we adopted two *syntactic* and three *semantic code refactoring operations* (See Section II-E) to contaminated code. In principle, we choose the operators that could be *automated* to avoid subjectivity brought by manual operation, reduce manual effort, and increase the *practicality and reusability* of these refactor operators.

On the one hand, *syntactic operators* check whether CLMs rely on rote memorization of pattern recognition. The refactoring preserves the semantic purpose but alters its syntactic structure [34], [35]. If CLMs truly understand the semantics of the code, their performance should not significantly decline when faced with semantically equivalent but syntactically varied code snippets. On the other hand, the *semantic operators* ensure CLMs understand code semantics rather than keyword/string matching. Replacing with synonyms or adding more semantics should not significantly affect model performance.

The experiment design of RQ3 aims to identify refactoring operators that are more likely to *disrupt the memorized patterns*. The performance of CLMs on each refactoring operator is compared to assess which operators cause the greatest decline in performance, hence indicating a higher effectiveness in alleviating data contamination threats.

4) **Experiment Design on the Assumption of RQ1-RQ3:** Recap that using CLMs’ performance as a hint for data contamination has an assumption, *i.e.*, **the contaminated/cleaned data are on similar difficulty level**. If some code groups have a significantly higher or lower complexity than others, it is hard to explain CLMs’ differential performance on code groups, obscuring true data contamination effects. Therefore, we quantify the code complexity of different code groups used in RQ1-RQ3 using two code complexity metrics (see Section II-C).

Besides, we also compared the similarity between code groups horizontally using various code similarities covering syntactic and semantic information (see Section II-D), to check if any code group is significantly different from other code groups. As illustrated in Fig. 5, the average code similarity between two code groups is calculated pairwise.

5) **Experiment Design for RQ4:** Various techniques and algorithms have been proposed for membership inference attack (MIA) [14], [36] or indicate how likely the model is to predict the given sequence [29], [15], [37]. Though these metrics are not designed for measuring data contamination, they could still serve as indicators. For example, perplexity indicates how

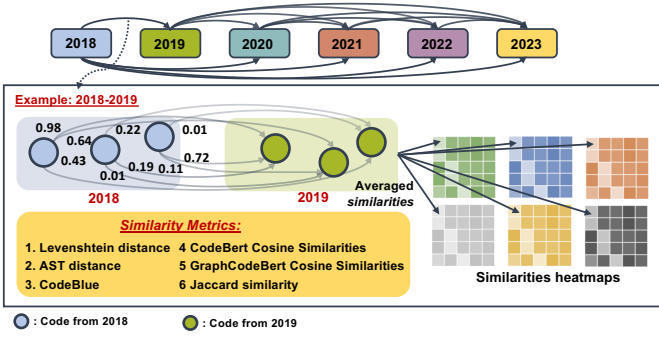


Fig. 2: Pairwise Code Similarity Comparison

well the given data fits the models' training data. In the study, we adopt three types of MIA-related metrics (see Section II-F) with different parameters, resulting in six MIA-related metrics. Under each metric, we compare the scores of different code groups (e.g., contaminated/recent, contaminated/refactored), investigating whether these MIA-related metrics could distinguish the code groups.

In addition, though these metrics could not directly serve as data contamination measurements, they could still indicate *how close the data distribution is* between the code group and the model's training data. In particular, if a code group achieves an average lower score than other groups, it means that the distribution of the code group is closer to that of training data.

B. Coding Task

We chose *code completion* because it aligns well with the characteristics of CLMs (i.e., generative models that receive prefixes and complete suffixes) [38], [39]. Also, this task simulates a common scenario in software development where a programmer starts writing a piece of code and relies on tooling to suggest how to continue. In addition, all CLMs support the code completion task. In comparison, tasks such as code infill (i.e., infilling the missing code inside the code snippets) and *code generation* (i.e., given textual descriptions, generating code) are not supported by all CLMs. Studying code completion tasks may achieve broader generalization.

Another reason for choosing code completion is that it could be conducted on the original collected code *without human intervention*. To minimize the impact of human factors on this study, it is prudent to preserve the authenticity of the collected code. Introducing artificial elements (e.g., natural language descriptions, code comments, and docstrings that are necessary for code generation) may *inadvertently influence the model's performance* by providing additional context or clues that are not inherently part of the code's logic or structure. Therefore, we consistently automate the collection and preprocessing for all code groups, ensuring fair data processing and curation.

C. Code Complexity Metrics

To check the validity of the assumption of RQ1-RQ3, we consider two metrics to evaluate code complexity.

1) **Cyclomatic complexity**: Cyclomatic complexity [40], [41] is a metric used to measure the complexity of a piece of software code. It quantifies the number of decision points

or branches in a program, indicating the potential number of unique paths through the code. It is calculated based on the control flow graph of the code, which represents the flow of control between different statements, loops, conditionals, and function calls. It is determined by *counting the number of decision points in the graph*, including conditional statements (if, switch), loops (for, while), and logical operators. We leverage Radon [42] to calculate it on Python functions.

2) **Cognitive complexity**: Unlike cyclomatic complexity, which counts the number of decision points and branches in code, cognitive complexity [43] considers factors such as nesting levels, logical operators, and boolean expressions that can make code harder to understand. It focuses on the structural complexity of the code and how it impacts human comprehension. In the implementation, we count the summation of the number of nesting levels, branches, and exceptions leveraging an off-the-shelf package [44].

D. Code Similarity Metrics

Apart from code complexity metrics, we also consider six code similarity metrics to calculate the distance between code groups, further validating the assumption of RQ1-RQ3.

1) **Levenshtein Similarity**: Levenshtein distance measures the minimum number of operations required to transform one string into another. To project this distance into a similarity between 0 and 1, we normalize the distance by dividing the maximum string lengths of two strings.

2) **CodeBleu Similarity**: Code Bilingual Evaluation Understudy [45] (CodeBLEU) measures the quality of machine-generated code against human reference code. Using n-grams of different lengths, CodeBLEU can assess the lexical structural and syntactic similarity between generated and reference codes.

3) **Jaccard Similarity**: Jaccard Similarity is a measure used to determine the similarity between two sets. It calculates the ratio of the size of the intersection of the sets to the size of their union. In the implementation, we follow the practice of existing works [46], which use MinHash [47] to approximate the Jaccard Similarity between code snippets to boost efficiency.

E. Code Refactoring Operators

For RQ3, we consider five refactoring operators, covering *syntactic* and *semantic* refactoring.

1) **Syntactic Refactoring Operators**: alters the code syntactic structure while the identifiers' names are untouched.

- **If-condition flipping (IFF)**. We flip the condition of `if` statements and their true and false branches to syntactically restructure the program in our dataset while the original program logic and semantics are preserved.
- **Loop transformation (Loop)**. Similarly, we replace `while` loops with equivalent `for` loop structures, and vice versa, to syntactically refactor the code.

2) **Semantic Refactoring Operators**: changes identifier names or adds additional context while the code's functionality remains unchanged.

- **Identifier Renaming (Renm)**. To prevent CLMs from simply searching variable names, we leverage wordhoard [48] to

fetch the synonyms of variable names and rename the variables in the Python code. Note that we rename identifiers to synonyms instead of arbitrary names in that the variable names often indicate the code’s functionality to be completed and serve as a necessary hint for CLMs to interpret the coding task and generate effective completions.

- **Special Parameter Appending (Param).** We refactor functions by appending unnamed positional parameters (`*args`) and keyword parameters (`**kwargs`) in the parameter list of function declarations if they do not exist. The appended parameters are unused in the function. Such refactoring does not change the behaviors of the function logic or any of its call sites while altering the function signature into a different form from what was seen in the CLMs’ training set.
- **Performance Measurement Decoration (Deco).** Since measuring the performance of a function, *e.g.*, memory usage and execution time do not change the function’s behavior, we also add performance measuring decorators to functions: `@timing`, which measures the total execution time of the decorated function, and `@measure_memory_usage`, which measures the memory usage of the function.

F. MIA (Membership Inference Attack)-related Metrics

For RQ4, we consider three existing metrics that may be used as indicators for data contamination, exploring whether they can distinguish different code groups. Note that these metrics are **model-dependent**. Different models assign probabilities to the same input data, resulting in different scores.

1) **Perplexity (*ppl*):** Perplexity [29] is a measure used to evaluate the predictability of a sequence of words. Usually, the perplexity measures how “surprised” the model is to see a given value [36]. It is defined as the exponentiated average negative log-likelihood of a sequence. Given a tokenized sequence X with N tokens $= (x_0, x_1, \dots, x_N)$, the perplexity of X is:

$$ppl(x) = \exp\left\{-\frac{1}{N} \sum_{i=1}^N \log p_{\theta}(x_i | x_0 \dots x_{i-1})\right\} \quad (1)$$

where $\log p_{\theta}$ is the log-likelihood of x_i conditioned on the prefix $x_0 \dots x_{i-1}$ according to the model. The lower the perplexity, the more natural the X to the model.

2) **Perplexity_{lower} (*ppl_{lower}*):** A variation [14] of the perplexity measure can be implemented by converting the input data to lowercase before calculating perplexity by removing the variability introduced by capitalization.

3) **Zlib Compression Entropy:** The Zlib of a given input x denotes the number of bits after compressing x with zlib. Zlib entropy [49] of the text (*i.e.*, the number of bits after compression with zlib) could be used as an indicator [15], [50]

$$zlib_{en}(x) = ppl(x) / \ln(zlib(x)) \quad (2)$$

4) **MIN-K% PROB:** MIN-K% PROB [14] was proposed to for membership inference attacks. It computes the score

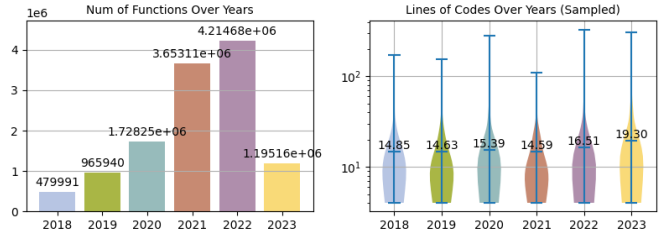


Fig. 3: Statistics of Code Over Years

using the $k\%$ of tokens with the lowest likelihoods.

$$MIN_{K\%} Prob(x) = \frac{1}{E} \sum_{x_i \in Min_{K\%}(x)} \log p(x_i | x_1, \dots, x_{i-1}) \quad (3)$$

where E is the size of the $Min_{K\%}(x)$ set. In this study, we set K as 5.0, 10.0, and 20.0 as set in previous work [14].

III. EXPERIMENT PREPARATION

A. Model and Training Data

We select the state-of-the-art CLMs widely studied in recent code generation work. Table I shows the **model information** (including model name, base model, model size, and the first release time) and **training data information** (including training source, the time span of the training data, and the languages the models support). Note that since the majority of models do not specify the exact time of their training data, we then elaborate on how we determine the model’s first release time, the data source, and the time span.

On the one hand, we determined a model’s **first release time** according to the commit time when the model was uploaded. Note that since the first commit usually only updates the `README.md` without uploading the model, also, other files (*e.g.*, licenses) could be uploaded before the model upload, so we carefully traced through the commit history in chronological order from past to present and located the first commit when the model checkpoints were uploaded. On the other hand, the **training source** and time span of the training data are either determined by the model reports or inferred by their model release time (*i.e.*, the cut-off date of the training data should at least goes before the model’s first release).

In the following, we list the detailed evidence of each model’s model and data information. ① **StarCoder** [52] explicitly states their training data, **the Stack** [51]. The model was first released on May 2023 according to the commit history [53]. While **the Stack** contains over 6TB of permissively licensed source code files covering 358 programming languages, 220.92M active GitHub repository names were collected from the event archives published between January 1st, 2015, and March 31st, 2022. ② The **StarChat** [54] is a fine-tuned version for assisting coding tasks. The model was released on June 7th, 2023 [55]. It is fine-tuned on OpenAssistant dataset [56] whose upload date is April 12th, 2023 according to the commit time [57]. ③ The **WizardCoder** [58] [59], was fine-tuned based on StarCoder-15.5B. The training source includes the training data for **StarCoder**, *i.e.*, **the Stack** and 20K instruction-following data

TABLE I: Large Language Models and Datasets

Code Large Language Models Statistics							
Model Information				Data Information			
Index	Base Model	Model	Size	1st Release	Source	Time Span	Language
①	StarCoder	StarCoder	15.5B	May 2023	the Stack [51]	Jan 2015~ Mar 2022	358 Programming languages
②		StarChat	15.5B	Jun 2023	GitHub	Up to April 2023	80+ programming languages.
③		WizardCoder	15.5B	Jun 2023	the Stack [51], CodeAlpaca-20k	Jan 2015~ Mar 2023	358 Programming languages and English
④	Llama2	CodeLlama-Instruct	7B	Apr 2023	GitHub + StackOverflow	Up to Jan 2023	Python, C/C++, TypeScript, Java and more
⑤		Phind-CodeLlama-34b	34B	Aug 2023	GitHub + StackOverflow	Up to Aug 2023	Python, C/C++, TypeScript, Java and more
⑥	GLM	ChatGLM2	6B	Oct 2023	–	Up to May 2022	Python, Java, JavaScript
⑦	ChatGPT	ChatGPT 3.5-turbo	–	June 2023	Public Data	Up to Sep 2021	95+ Programming languages
⑧		Github-Copilot	–	Jun 2021	Public Data	Up to April 2023	95+ Programming languages

used for fine-tuning the Code Alpaca model [60] whose upload date is May 13th, 2023 [61]. ④ **CodeLlama-7b-Instruct** [62] was first uploaded on August 24th, 2023 [63]. CodeLlama’s training data is the same as that of Llama 2, which was trained between January 2023 and July 2023 [64]. ⑤ **Phind** [65] was fine-tuned on CodeLlama with an additional 1.5B tokens high-quality programming-related data. It was first released on Aug 29, 2023 [66] and since the cut-off-date of the additional data was disclosed, we assume this time as its training source cut-off-date. ⑥ **ChatGLM** [67], [68], [69] was first released on October 26th, 2023 [70]. Since the training source was not disclosed, we assume the training source was collected before October, 2023. ⑦ **ChatGPT 3.5** has various variants. We use the *ChatGPT3.5-turbo-0613* in the evaluation. Its first release time is June 2023, while its training data was cut off on May 2021 [71]. ⑧ **Github Copilot** [72] was first released on June 29th, 2021 [73]. From November 30th, 2023, Copilot is empowered by GPT-4, [74] whose cut-off date is April 2023 [75].

B. Data Collection

RQ1 requires chronological data over a long period of time to construct annual code groups. As such, we collect the Python code from January 1st, 2018, to December 31st, 2023. Though there are several datasets available [51], they do not cover the most recent data (*e.g.*, the latest six months). So, we first collect data from the Stack v2 [26], which contains code until March 31st, 2022 (Section III-B1), and then collect newer data after March 31st, 2022, by crawling newly created repositories on GitHub (Section III-B2).

1) **Data Collection from the Stack v2**: We take the Stack v2 [26] dataset as the code corpus for its representativeness and comprehensiveness. We use the deduplicated version of the *stack*, which contains over 2.9TB data of permissively licensed source code files covering 358 programming languages, from January 1st, 2015 to March 31st, 2022. Typically, we take Python as our subject programming language for its popularity.

2) **Newer Data Collection**: To cover the code afterward (after March 31st, 2022), we crawl the Python repositories from GitHub with permissive licenses. In particular, we collect the repositories created from April 1st, 2022 to December 31st, 2023 with more than 50 stars. Note that we only collected repositories that are **created after** the time range of the Stack v2, so there should not be an overlap between the

collected new data and the Stack v2. As a result, a total of 15,743 repositories are collected.

3) **Data Processing**: We process the collected data by extracting functions from the Python code. We filter out the functions that depend on other functions (*e.g.*, the functions inside `classes` or inside other functions) or the empty functions (*i.e.*, the function body is simply `pass`).

4) **Statistics of Collected Data**: After data collection and preparation, we collect a total of 12,493,174 Python functions in total. The left subfigure of Fig. 3 shows the total number of functions collected yearly. The number of functions increases over years except for 2023. The reason is that we only collect the data of 2023 from GitHub repositories with at least 50 stars, instead of all GitHub repositories like *the Stack v2*. As a result, the collected functions in 2023 are fewer than in previous years. Nevertheless, since we conduct experiments on a sample of functions in each year and the total number of functions in 2023 is large enough (over 1 million), this difference in annual function numbers should not be a significant threat to our study.

Fig. 3 (right) shows the lines of code (LOCs) in 384 sampled functions each year. The lines of code in different years are essentially stable, with an average of 15.88 LOCs each year.

C. Annual Code Groups Preparation (RQ1)

Since there are over ten thousand code functions per year, traversing them all is unrealistic. Thus, we follow the practice of statistics of random sampling, setting a 95% confidence level and 5% margin of error (*i.e.*, there is a 95% confidence that the true population falls within $\pm 5\%$ of the sample estimate). It results in an average of 384 code functions annually. The sampled 384 functions in each year form the annual code groups in Fig. 1. We refer to annual code groups as Code-2018, Code-2019, Code-2020, Code-2021, Code-2022, Code-2023, respectively.

D. Curated Code Groups Preparation (RQ2)

HumanEval [17] dataset contains a set of Python functions, so we construct the Code-HumE code group directly with all functions in HumanEval. CoderEval [20] dataset contains a set of Python files, so we extract functions from all files like in Section III-B3 and then construct the Code-CodE code group.

E. Refactored Code Groups Preparation (RQ3)

The refactoring operators are expected to apply to contaminated code. According to Table I, the earliest cut-off date

is September 2021, so we chose Code-2021 to apply the refactoring. In particular, we apply the operators introduced in Section II-E to Code-2021. Each operator is applied to every function in Code-2021 where applicable. If an operator is not applicable (e.g., when attempting to flip if-else branches in code that does not contain an if-statement) to a certain code function under refactoring, we skip refactoring that function.

If there are multiple ways to apply a particular operator, we randomly select one way to use it. For example, if one function has multiple if-conditions, we randomly pick one if-condition and apply the IFF (if-condition flipping) operator. The rationale behind employing a single operator at most once without repetition is that we want to analyze the impact of the most granular refactoring operator. Besides, by applying at most one type of operator to one function, we aim to isolate the influence of an individual operator rather than the compounded effects resulting from multiple operators. We use Code-IFF, Code-Loop, Code-Renm, Code-Param, and Code-Deco to refer to the code group constructed by applying the five refactoring operations in Section II-E, respectively.

F. Coding Task Preparation

The code completion task understands the code prefix and generates the subsequent code. Since we want to minimize the human factors (e.g., manual annotations, subjective interpretations) in the code, we keep the code untouched without human curation and annotation. Indeed, without a natural language description, CLMs need to understand the code’s intention from the code as-is and complete the suffix.

To retain the code intent to the greatest extent, we only masked out the last statement in each function. Since we filtered out those functions with less than 3 lines when data preparation, this avoids the situation where the function body is empty after making out the last statement in the function. Furthermore, because only the final statement has been masked, the flexibility of code completion is relatively constrained. In other words, the correct ways to complete it should be limited, confined to completing only a single line of code. This setting also makes it easier to judge whether the models’ output is correct. We use exact string matching to determine the outputs’ correctness unless the last statement allows flexibility (such as a print statement). We ignore the concrete content for such cases and only match whether the print statement is predicted.

IV. EVALUATION

We use nucleus sampling [76] in line with recent works [19], [20], [77], where five solution samples are randomly generated with a temperature of 0.2 [17]. The experiments are conducted on a computational infrastructure comprising two NVIDIA RTX 6000 Ada GPUs, each with 48GB of graphics memory.

A. Evaluation on the Assumption of RQ1-RQ3

Before we start RQ1-RQ3, which uses CLMs’ performance as an indicator for data contamination, there is an assumption that the contaminated/cleansed data are on similar difficulty

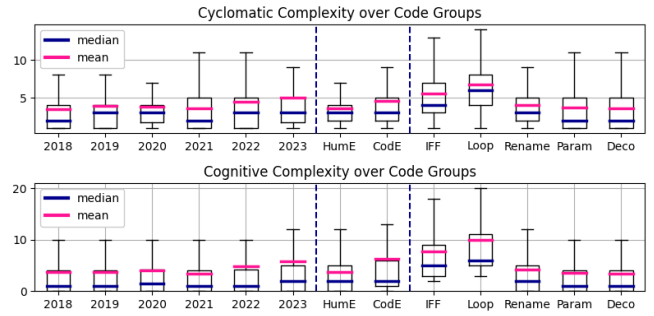


Fig. 4: Complexity Comparison of Various Code Groups

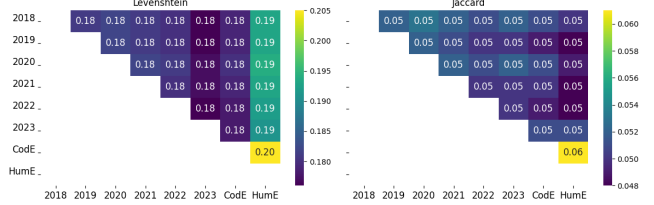


Fig. 5: Code Similarities Over Various Code Groups

levels. Thus, we demonstrate the code complexity of code groups and show the similarities between code groups.

The code complexities over code groups are shown in Fig. 4. The average cyclomatic and cognitive complexities over code groups are generally similar, with slight fluctuation. In particular, among Code-2018 to Code-2023 (groups before the first dashed vertical line), *a slight ascending trend* could be observed on both complexities, reaching the highest mean and medium complexities of Code-2023. For benchmarks (between the first and second vertical dashed lines), the complexities of Code-CodeE are higher than that of Code-HumE in both metrics, meaning that the codes in Code-CodeE are slightly more complex than Code-HumE on average, which is in line with the observation in previous work [20].

Among the refactored code groups (after the second vertical line), the Code-IFF and Code-Loop have the highest scores on both complexities, while other factored groups are similar to Code-2018 to Code-2021 with small fluctuations.

Finding 1: The code groups display *similar code complexity* in terms of cyclomatic and cognitive complexities, validating the assumptions for RQ1-RQ3 despite minor fluctuations.
Finding 2: Code complexity across different years appears to be relatively stable, exhibiting a slight upward trend.
Finding 3: The code complexity of CoderEval is slightly higher than early code groups, while HumanEval is at a similar complexity level as early codes.

We further measure the three code similarities over various code groups. The results of Levenshtein and Jaccard distance are visualized in Fig. 5. The CodeBleu stays at an extremely low level (from 0.000018 to 0.009587), so we omit the visualization of it. The heatmaps show that the similarities over code groups do not differ substantially. Also, clear that the similarities between code groups are relatively low, ranging from 0.18 to

0.20 (Levenshtein similarity) and 0.05 to 0.06 (Jaccard similarity). Note that the existing work [26] uses a threshold of 0.85 Jaccard similarity. The average results over these code groups are significantly lower, indicating the overlapping degree among these code groups are at a relatively minimal level. Similarity metrics among the code groups do not differ substantially.

Finding 4: The various code groups exhibit a *consistent level of similarity* among themselves, which echoes the observation on code complexity.

B. RQ1. How does CLMs’ performance differ on contaminated/recent data?

Table II (from entries “Code-2018” to “Code-2023-oct”) shows Pass@k with nucleus sampling on code groups from different years. Overall, Copilot consistently outperforms other models in all the code groups, followed by Phind. Horizontally, the value of Pass@k *stays stable from Code-2018 to Code-2022*, followed by *a slight decrease in Code-2022* observed in all CLMs. Combining the complexity analysis in Fig. 4, it is reasonable because Code-2022 has higher code complexities than previous code groups.

An interesting situation happens in Code-2023. Though most CLMs were released before 2023 (meaning that these CLMs *are less likely* to see these codes), and the complexity of Code-2023 is higher than the earlier code groups, the *Pass@k in Code-2023 consistently outperforms earlier code groups over all the CLMs*. This observation is counter-intuitive because the general belief is that the more recent data suffer less from data contamination. So, we further sampled a code group after Oct 2023, which is later than *all* CLMs’ cut-off date, to check whether the observation still holds. See the column “Code-2023-oct”; the results *echo* that in Code-2023 again. CLMs even achieve higher performance in Code-2023-oct.

A possible reason behind this is the popularity of AI programming assistants such as Copilot. As reported by Github in June 2022 [78], nearly 40% of the code is being written by GitHub Copilot in popular coding languages like Python in last one year. That means the code created and committed in 2023 may be initially predicted by AI; it stands to reason that CLMs perform better on them.

Finding 5: CLMs do not necessarily perform worse on recent data (*i.e.*, later than their release). On the contrary, the CLMs may perform better on recent codes. It indicates that *using more recent data may not be an ideal strategy* to alleviate data contamination.

Implication: The popularity of AI programming assistants such as Copilot may further *exacerbate data contamination threats*. It also indicates that simply amassing more recent data may not be optimal when creating new benchmarks.

C. RQ2. How does CLMs’ performance differ on contaminated/curated data?

The performance on the curated datasets (*i.e.*, HumanEval [17] and CoderEval [20]) is shown in the last two

entries (*i.e.*, “Code-Code” and “Code-HumE”) in Table II. *In general, CLMs achieve better performance on curated datasets compared to that on the contaminated data in terms of Pass@k scores.* Interestingly, the performance of CLMs is significantly better on CoderEval even though its release time is later than most models’ release time. This may be because these curated benchmarks are carefully screened and processed manually, the code quality is relatively high, and it is more likely to conform to the distribution of instruction data of CLMs that has been supervised and fine-tuned. Such results indicate that the program comprehension and completion ability of CLMs are effectively generalized to datasets curated after the models’ cut-off date.

Finding 6: CLMs perform better on recently curated datasets, compared to the performance on the contaminated code groups. It indicates that *using curated datasets may not be effective* in mitigating data contamination threats.

D. RQ3. How does CLMs’ performance differ on contaminated/refactored data?

Table III shows the Pass@1 scores of different CLMs on the code group Code-2021 and various refactored code groups. **Syntactic Refactoring Operators.** The change in the Pass@1 of CLMs varies. IFF is capable of degrading the Pass@1 score of some CLMs while loop transformation may not be effective. The possible reason behind this may be that such syntactic refactoring may or may not make the code close to the data distribution of the training set. On the contrary, syntactic refactoring could even upgrade the model’s performance (Code-Loop). Hence, the performance may vary. **Semantic Refactoring Operators.** The Pass@1 score of CLMs decreases after refactoring. Renaming identifiers with their synonyms decreases the Pass@1 scores of CLMs in that the new identifier name, although it is the synonym of the original name, may convey different contextual meanings in the code. For instance, in the following code snippet `encoded = base64.b64encode(io.open(mp4, 'r+b').read())`, a video file is read and stored as base64 encoding in variable `encoded`. The refactoring operator renames the variable to `ciphered`. Although `ciphered` is a synonym of `encoded` in natural language, it does not fit in the coding context since there is no cryptographic ciphering procedure. As a result, the renaming may alter the semantics of the program and provide different hints to CLMs to complete the code, so that the CLMs’ performance is influenced.

Special parameter appending and performance measurement decoration prepend additional contexts to CLMs’ prompt but the original functionalities of the program are left unchanged. We observe that such refactoring operators are effective in degrading the CLMs’ performance.

Finding 7: *Syntactic refactoring* may not be useful to alleviate data contamination. On the contrary, refactoring the code structure could *even upgrade* the models’ performance.

TABLE II: **Pass@k with Nucleus Sampling on Contaminated/Recent/Curated Code Groups.** The **greener**, the larger. We **underline** the results on the data collected *after* each model’s release time (*i.e.*, *cleansed*) to facilitate clearer demonstration.

Model	Code-2018			Code-2019			Code-2020			Code-2021			Code-2022			Code-2023			Code-2023-oct			Code-CodE			Code-HumE		
	P@1	P@3	P@5	P@1	P@3	P@5	P@1	P@3	P@5	P@1	P@3	P@5	P@1	P@3	P@5	P@1	P@3	P@5	P@1	P@3	P@5	P@1	P@3	P@5	P@1	P@3	P@5
StarCoder	30.70	32.60	33.60	30.50	32.00	32.80	30.20	31.20	33.10	29.40	30.20	31.20	27.30	29.40	30.50	<u>34.10</u>	<u>35.40</u>	<u>36.50</u>	<u>33.90</u>	<u>36.70</u>	<u>37.80</u>	<u>40.90</u>	<u>41.70</u>	<u>42.60</u>	51.20	54.30	54.90
StarChat	37.20	39.80	40.60	34.40	37.00	38.00	35.70	38.80	39.60	38.50	41.10	42.20	32.60	35.40	37.20	39.30	41.10	41.70	<u>41.90</u>	<u>41.90</u>	<u>41.90</u>	48.70	52.60	53.50	48.80	54.90	54.90
WizardCoder	38.50	39.60	40.90	34.90	37.50	37.80	38.50	40.40	40.90	40.40	43.20	43.80	34.10	36.20	37.20	40.60	41.90	42.40	<u>43.80</u>	<u>47.40</u>	<u>48.40</u>	52.20	56.50	58.30	55.50	57.30	57.90
Codellama	37.00	39.60	40.60	37.20	39.80	40.10	37.00	40.60	41.10	38.50	40.60	41.10	33.60	35.90	37.00	<u>38.00</u>	<u>40.40</u>	<u>41.40</u>	<u>43.20</u>	<u>46.40</u>	<u>46.60</u>	52.20	55.70	57.00	53.00	54.30	54.90
Phind	41.40	42.70	43.80	40.90	43.50	44.00	44.00	45.80	46.40	44.30	46.90	47.40	38.80	41.40	42.40	38.80	39.60	39.60	<u>46.40</u>	<u>47.40</u>	<u>47.40</u>	60.90	60.90	61.30	56.10	57.30	57.90
ChatGLM2	22.70	28.60	32.60	24.20	26.00	27.30	24.20	26.80	27.60	21.40	24.70	25.50	22.10	24.00	25.00	<u>22.70</u>	<u>25.30</u>	<u>26.80</u>	<u>27.30</u>	<u>28.10</u>	<u>30.50</u>	30.90	32.60	<u>34.80</u>	34.80	36.00	37.20
ChatGPT-3.5	32.30	34.10	34.40	30.20	31.80	32.60	30.70	31.00	31.20	31.80	33.30	33.90	<u>28.10</u>	<u>30.20</u>	<u>30.20</u>	<u>34.10</u>	<u>35.20</u>	<u>35.40</u>	<u>37.50</u>	<u>39.10</u>	<u>39.80</u>	<u>47.40</u>	<u>48.70</u>	<u>49.10</u>	51.80	54.30	55.50
Github-Copilot	48.40	49.70	50.50	49.00	49.70	50.30	48.70	50.80	51.60	53.40	55.50	56.20	44.80	46.10	46.60	50.50	52.60	53.10	<u>56.00</u>	<u>58.30</u>	<u>58.90</u>	60.00	63.00	64.80	57.90	61.60	62.20

TABLE III: **Pass@1 Results (Greedy) on Original (Code-2021) and Various Refactored Code Groups.** The entries Δ is the difference between the previous column and the “Origin” column. Foreground color **blue** means decrease and **red** means increase. The background color highlights the value of the Pass@1 result: the **greener**, the larger.

Model	2021	Code-IFF		Code-Loop		Code-Renm		Code-Param		Code-Deco	
	P@1	P@1	Δ	P@1	Δ	P@1	Δ	P@1	Δ	P@1	Δ
Codellama	38.5	35.1	-3.4	41.8	3.3	38.2	-0.3	37.3	-1.2	23.4	-15.1
StarCode	29.4	30.5	1.1	34.4	5.0	30.5	1.1	27.6	-1.8	24.7	-4.7
StarCoder-chat	38.5	38.5	0.0	40.2	1.7	37.1	-1.4	36.7	-1.8	30.5	-8.0
Wizardcode	40.4	39.1	-1.3	47.5	7.1	40.4	0.0	40.2	-0.2	24.5	-15.9
Phind	44.3	48.3	4.0	52.5	8.2	43.3	-1.0	42.5	-1.8	39.1	-5.2
Chatglm2	21.4	20.7	-0.7	23.8	2.4	21.8	0.4	20.2	-1.2	19.5	-1.9
gpt-3.5	31.8	32.8	1.0	36.9	5.1	32.4	0.6	31.8	0.0	29.4	-2.4
Copilot	53.4	58.0	4.6	64.8	11.4	57.1	3.7	46.2	-7.2	46.1	-7.3

Implication: Semantic code refactoring operators (Renm, Param and Deco) could be more effective in evaluating the data contamination threats of CLMs.

E. RQ4. Could existing metrics distinguish contaminated/cleansed data?

Fig. 6 shows six MIA-related scores on different code groups. The figure shows the results calculated on StarCoder; the results on other CLMs are similar, so we omit them due to space limits. Note that these scores indicate how “natural” the models found the inputs are. The lower the score, the closer the input fits the models’ training distribution.

The contaminated/cleansed code groups generally share similar scores in all metrics. The mean perplexity of the code groups ranges from 3.37 (Code-2023) to 5.17 (Code-CodE). The largest score difference is observed in ppl_lower, ranging from 4.06 (Code-2023) to 6.41 (Code-CodE). In addition, the Code-2023, Code-CodE, and the five refactored code groups are unseen for StarCoder whose cut-off date is March 2022, yet there is no apparent difference in these cleansed code groups compared with that on contaminated code groups.

Finding 8: The existing six MIA-related metrics can *hardly distinguish the contaminated/cleansed data*, so they may be hardly used for quantifying data contamination.

Surprisingly, though the recent code, *i.e.*, Code-2023, is later than StarCoder’s cut-off date (March 2022), all six metrics show absolute drops on it compared with other code groups. Code-2023 fits CLMs’ training distribution better than other code groups, so the model found it more natural (lower scores).

This observation may partially explain the counterintuitiveness we observed in RQ1. CLMs find Code-2023 more natural than earlier code groups and closer to the training distribution; their performance will naturally be better.

One reason is the developers’ tendency to duplicate codes through copy-paste. However, this practice has been a long-established norm within the programming community, as discussed in prior literature [79]. We would expect a similar pattern in previous data if it were the primary factor. When we look at the earlier code groups (*e.g.*, Code-2018 to Code-2021), the scores are stable without a significant drop, suggesting that other factors may influence this trend.

Another possible reason is the emergence of AI coding assistants such as Github Copilot. As we mentioned in RQ1, nearly 40% of the code has been written by GitHub Copilot in the last year. As this trend continues, more developers use AI as a programming aid, making the codes generated by AI more naturally fit the training distribution of the models.

In addition, looking at the refactored code groups in all six subplots, code refactoring appears not to have significantly impacted the model’s familiarity with the code. The scores only exhibited minor disturbances. Among these refactoring operators, identifier renaming (Code-Renm) and performance measurement decoration (Code-Deco) have a larger impact, indicating that these two operators are more likely to result in a divergence from the original training data distribution. Such observation also echoes the findings in RQ3.

Finding 9: Identifier renaming (Renm) and performance measurement decoration (Deco) have a greater impact on data, making CLMs unfamiliar with the data.

Implication: Renm and Deco may be more useful for alleviating data contamination.

V. THREATS TO VALIDITY

We acknowledge several threats to the validity of our conclusions. First, **programming language selection bias**. Our study confines itself to the analysis of Python, which, while being one of the most popular programming languages, does not represent other languages. We justify the selection of Python due to its widespread adoption and the belief that insights derived from popular languages may have broader relevance. Second, **coding task selection bias**. The study focuses on code completion tasks with the masking of the last statement may not encapsulate the full range of capabilities required

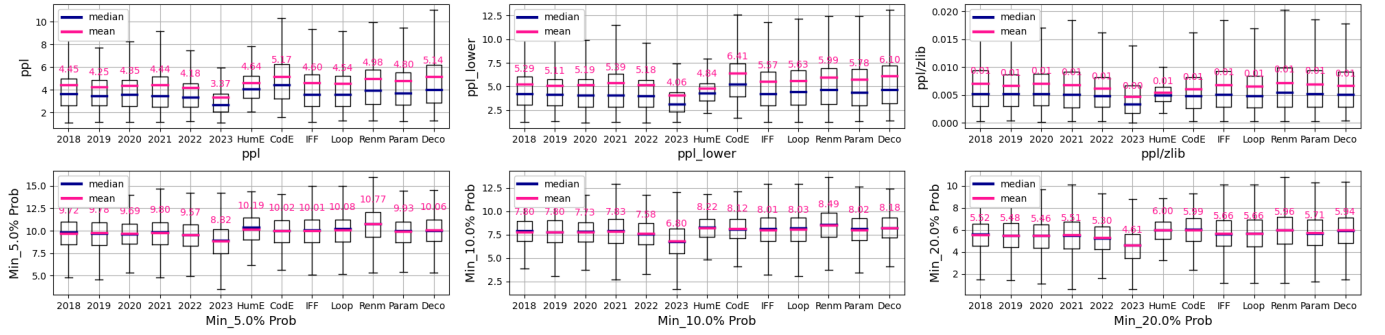


Fig. 6: Contamination-related Scores Over Various Code Groups

for other coding tasks. The performance of CLMs could vary on other tasks. We selected it because of its relevance to real-world coding practices and the advantage of minimizing human intervention, which can introduce additional variability. To enrich our understanding of CLMs’ capabilities, we call for further research to explore a wider array of coding tasks. Third, **possible semantics overlapping between code groups**. While efforts were made to prevent the overlap of code groups, the possibility of semantic overlap cannot be entirely eliminated. To mitigate this threat, three similarity metrics were employed to ascertain the distinctness of code groups, and the results suggest a relatively low level of similarity.

VI. RELATED WORK

A. Data Contamination in Large Language Models

Several recent studies [6], [5], [7], [80], [13], [28], [81], [82] have explored data contamination in large language models (LLMs). Brown *et al.* [15] analyzed the contamination in GPT-3, showing that the performance between contaminated and clean data is not necessarily different. Carlini *et al.* [13] found that the memorization of LLMs grows with model size, training data duplicates, and prompt length. Kandpal *et al.* [28] identified that the success of privacy attacks on LLMs could be attributed to the sequence duplication in the training set. Razeghi *et al.* [81] studied the correlations between LLMs’ performance and the frequency of terms and observed that LLMs perform better on inputs with more frequent terms. Magar *et al.* [82] studied the memorization and exploitation of masked language models in a controlled manner. They fine-tune two models, one with labeled test data and the other without. They identified the number of duplications of the contaminated data and the model size that affected model exploitation. Our work is based on different assumptions. Magar *et al.* pre-trained and fine-tuned models (*i.e.*, BERT) from scratch, with a controlled training set, aiming at identifying factors that affect model exploitation, while we study the efficacy of countermeasures for alleviating data contamination and identifying effective countermeasures.

B. Membership Inference Attacks

Membership inference attacks (MIAs) try to determine whether a particular data is contained in the model’s training

data. Various metrics are proposed to infer data membership, including LOSS [27], reference models [15], perplexity [29], Zlib Entropy [49], Neighborhood attack [83], Min-k% Prob [14], *etc.* Though MIAs are extensively studied in traditional deep learning models, the research on MIA in LLMs is limited. Recently, Duan *et al.* [30] conducted a large-scale MIA on LLMs and found that MIAs barely outperform random guessing across varying LLM sizes and domains. They also identified that different temporal ranges of data may affect model performance. These works aim at extracting sensitive data from LLMs, While our study considers how to evaluate CLMs fairly.

C. Countermeasures Alleviating Data Contamination Threat

Various benchmarks [22], [31], [17], [84], [20], [19], [9] for coding tasks are proposed to evaluate CLMs in diverse abilities and on a different scales. Most of these benchmarks were constructed or processed manually to ensure the quality of source code and other materials (*e.g.*, docstrings, test cases). Yet, it is unclear how CLMs perform on these benchmarks compared with training data and how large the difference is. This study attempts to answer this question.

Several works refactored code to alleviate data contamination. Wu *et al.* [9] adopted identifier renaming and Code Structure Change to the code they collected. They conducted the refactoring manually, and they did not compare the performance change after code refactoring. Besides, several studies [12], [85], [86] leverage code assistants such as GPT-4 to do the code refactoring. Yet, we avoid using AIs to code to get rid of reintroducing variants to the study.

VII. CONCLUSION

We conducted a systematic study on the data contamination threats of CLMs with contaminated and cleansed data. We investigated how CLMs perform on contaminated/recent code, code in curated benchmarks, and refactored code and found that CLMs in general achieve a stable performance between contaminated data and different types of cleansed data. Our study indicates that existing countermeasures in alleviating data contamination threats may not be effective, calling for more discussion on the related topics of CLM-based software engineering research.

REFERENCES

- [1] Y. Deng, C. S. Xia, H. Peng, C. Yang, and L. Zhang, "Large language models are zero-shot fuzzers: Fuzzing deep-learning libraries via large language models," in *Proceedings of the 32nd ACM SIGSOFT international symposium on software testing and analysis*, 2023, pp. 423–435.
- [2] Y. Deng, C. S. Xia, C. Yang, S. D. Zhang, S. Yang, and L. Zhang, "Large language models are edge-case generators: Crafting unusual programs for fuzzing deep learning libraries," in *Proceedings of the 46th IEEE/ACM International Conference on Software Engineering*, 2024, pp. 1–13.
- [3] C. S. Xia, M. Paltenghi, J. Le Tian, M. Pradel, and L. Zhang, "Fuzz4all: Universal fuzzing with large language models," *Proc. IEEE/ACM ICSE*, 2024.
- [4] C. S. Xia, Y. Wei, and L. Zhang, "Automated program repair in the era of large pre-trained language models," in *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*. IEEE, 2023, pp. 1482–1494.
- [5] S. Golchin and M. Surdeanu, "Time travel in llms: Tracing data contamination in large language models," *CoRR*, vol. abs/2308.08493, 2023. [Online]. Available: <https://doi.org/10.48550/arXiv.2308.08493>
- [6] O. Sainz, J. Campos, I. García-Ferrero, J. Etxaniz, O. L. de Lacalle, and E. Agirre, "NLP evaluation in trouble: On the need to measure LLM data contamination for each benchmark," in *Findings of the Association for Computational Linguistics: EMNLP 2023*, H. Bouamor, J. Pino, and K. Bali, Eds. Singapore: Association for Computational Linguistics, Dec. 2023, pp. 10 776–10 787. [Online]. Available: <https://aclanthology.org/2023.findings-emnlp.722>
- [7] C. Li and J. Flanigan, "Task contamination: Language models may not be few-shot anymore," *CoRR*, vol. abs/2312.16337, 2023. [Online]. Available: <https://doi.org/10.48550/arXiv.2312.16337>
- [8] Y. Li, "An open source data contamination report for llama series models," *arXiv preprint arXiv:2310.17589*, 2023.
- [9] Y. Wu, N. Jiang, H. V. Pham, T. Lutellier, J. Davis, L. Tan, P. Babkin, and S. Shah, "How effective are neural networks for fixing security vulnerabilities," in *Proceedings of the 32nd ACM SIGSOFT International Symposium on Software Testing and Analysis*, ser. ISSTA 2023. New York, NY, USA: Association for Computing Machinery, 2023, p. 1282–1294. [Online]. Available: <https://doi.org/10.1145/3597926.3598135>
- [10] T.-O. Li, W. Zong, Y. Wang, H. Tian, Y. Wang, S.-C. Cheung, and J. Kramer, "Nuances are the key: Unlocking chatgpt to find failure-inducing tests with differential prompting," in *2023 38th IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE, 2023, pp. 14–26.
- [11] J. Cao, M. Li, M. Wen, and S.-c. Cheung, "A study on prompt design, advantages and limitations of chatgpt for deep learning program repair," *arXiv preprint arXiv:2304.08191*, 2023.
- [12] A. Fan, B. Gokkaya, M. Harman, M. Lyubarskiy, S. Sengupta, S. Yoo, and J. M. Zhang, "Large language models for software engineering: Survey and open problems," *arXiv preprint arXiv:2310.03533*, 2023.
- [13] K. Tirumala, A. Markosyan, L. Zettlemoyer, and A. Aghajanyan, "Memorization without overfitting: Analyzing the training dynamics of large language models," *Advances in Neural Information Processing Systems*, vol. 35, pp. 38 274–38 290, 2022.
- [14] W. Shi, A. Ajith, M. Xia, Y. Huang, D. Liu, T. Blevins, D. Chen, and L. Zettlemoyer, "Detecting pretraining data from large language models," *arXiv preprint arXiv:2310.16789*, 2023.
- [15] N. Carlini, F. Tramèr, E. Wallace, M. Jagielski, A. Herbert-Voss, K. Lee, A. Roberts, T. B. Brown, D. Song, Ú. Erlingsson, A. Oprea, and C. Raffel, "Extracting training data from large language models," in *30th USENIX Security Symposium, USENIX Security 2021, August 11-13, 2021*, M. D. Bailey and R. Greenstadt, Eds. USENIX Association, 2021, pp. 2633–2650. [Online]. Available: <https://www.usenix.org/conference/usenixsecurity21/presentation/carlini-extracting>
- [16] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale *et al.*, "Llama 2: Open foundation and fine-tuned chat models," *arXiv preprint arXiv:2307.09288*, 2023.
- [17] M. Chen, J. Tworek, H. Jun, Q. Yuan, H. P. de Oliveira Pinto, J. Kaplan, H. Edwards, Y. Burda, N. Joseph, G. Brockman, A. Ray, R. Puri, G. Krueger, M. Petrov, H. Khlaaf, G. Sastry, P. Mishkin, B. Chan, S. Gray, N. Ryder, M. Pavlov, A. Power, L. Kaiser, M. Bavarian, C. Winter, P. Tillet, F. P. Such, D. Cummings, M. Plappert, F. Chantzis, E. Barnes, A. Herbert-Voss, W. H. Guss, A. Nichol, A. Paino, N. Tezak, J. Tang, I. Babuschkin, S. Balaji, S. Jain, W. Saunders, C. Hesse, A. N. Carr, J. Leike, J. Achiam, V. Misra, E. Morikawa, A. Radford, M. Knight, M. Brundage, M. Murati, K. Mayer, P. Welinder, B. McGrew, D. Amodei, S. McCandlish, I. Sutskever, and W. Zaremba, "Evaluating large language models trained on code," 2021.
- [18] J. Austin, A. Odena, M. Nye, M. Bosma, H. Michalewski, D. Dohan, E. Jiang, C. Cai, M. Terry, Q. Le *et al.*, "Program synthesis with large language models," *arXiv preprint arXiv:2108.07732*, 2021.
- [19] X. Du, M. Liu, K. Wang, H. Wang, J. Liu, Y. Chen, J. Feng, C. Sha, X. Peng, and Y. Lou, "Classeval: A manually-crafted benchmark for evaluating llms on class-level code generation," 2023.
- [20] H. Yu, B. Shen, D. Ran, J. Zhang, Q. Zhang, Y. Ma, G. Liang, Y. Li, T. Xie, and Q. Wang, "Codereval: A benchmark of pragmatic code generation with generative pre-trained models," *arXiv preprint arXiv:2302.00288*, 2023.
- [21] Y. Lai, C. Li, Y. Wang, T. Zhang, R. Zhong, L. Zettlemoyer, W.-T. Yih, D. Fried, S. Wang, and T. Yu, "DS-1000: A natural and reliable benchmark for data science code generation," in *Proceedings of the 40th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, A. Krause, E. Brunskill, K. Cho, B. Engelhardt, S. Sabato, and J. Scarlett, Eds., vol. 202. PMLR, 23–29 Jul 2023, pp. 18 319–18 345.
- [22] C. E. Jimenez, J. Yang, A. Wettig, S. Yao, K. Pei, O. Press, and K. R. Narasimhan, "SWE-bench: Can language models resolve real-world github issues?" in *The Twelfth International Conference on Learning Representations*, 2024. [Online]. Available: <https://openreview.net/forum?id=VTF8yNQMM66>
- [23] A. Shirafuji, Y. Oda, J. Suzuki, M. Morishita, and Y. Watanobe, "Refactoring programs using large language models with few-shot examples," *arXiv preprint arXiv:2311.11690*, 2023.
- [24] A. H. Huang, H. Wang, and Y. Yang, "Finbert: A large language model for extracting information from financial text," *Contemporary Accounting Research*, vol. 40, no. 2, pp. 806–841, 2023.
- [25] N. Gruver, M. Finzi, S. Qiu, and A. G. Wilson, "Large language models are zero-shot time series forecasters," *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [26] "bigcode/the-stack-v2-dedup," <https://huggingface.co/datasets/bigcode/the-stack-v2-dedup>, 2023.
- [27] S. Yeom, I. Giacomelli, M. Fredrikson, and S. Jha, "Privacy risk in machine learning: Analyzing the connection to overfitting," in *2018 IEEE 31st Computer Security Foundations Symposium (CSF)*, 2018, pp. 268–282.
- [28] N. Kandpal, E. Wallace, and C. Raffel, "Deduplicating training data mitigates privacy risks in language models," in *International Conference on Machine Learning*. PMLR, 2022, pp. 10 697–10 707.
- [29] F. Jelinek, R. L. Mercer, L. R. Bahl, and J. K. Baker, "Perplexity—a measure of the difficulty of speech recognition tasks," *The Journal of the Acoustical Society of America*, vol. 62, no. S1, pp. S63–S63, 1977.
- [30] M. Duan, A. Suri, N. Mireshghallah, S. Min, W. Shi, L. Zettlemoyer, Y. Tsvetkov, Y. Choi, D. Evans, and H. Hajishirzi, "Do membership inference attacks work on large language models?" *arXiv preprint arXiv:2402.07841*, 2024.
- [31] J. Austin, A. Odena, M. Nye, M. Bosma, H. Michalewski, D. Dohan, E. Jiang, C. Cai, M. Terry, Q. Le *et al.*, "Program synthesis with large language models," *arXiv preprint arXiv:2108.07732*, 2021.
- [32] "openai/human-eval," <https://github.com/openai/human-eval/commits/master/data>, 2021.
- [33] "Codereval/codereval," <https://github.com/CoderEval/CoderEval/commits/main/CoderEval4Python.json>, 2023.
- [34] A. A. B. Baqais and M. Alshayeb, "Automatic software refactoring: a systematic literature review," *Software Quality Journal*, vol. 28, no. 2, pp. 459–502, 2020.
- [35] J. Al Dallal and A. Abdin, "Empirical evaluation of the impact of object-oriented code refactoring on quality attributes: A systematic literature review," *IEEE Transactions on Software Engineering*, vol. 44, no. 1, pp. 44–69, 2017.
- [36] N. Carlini, C. Liu, Ú. Erlingsson, J. Kos, and D. Song, "The secret sharer: Evaluating and testing unintended memorization in neural networks," in *28th USENIX Security Symposium (USENIX Security 19)*, 2019, pp. 267–284.
- [37] Y. Li, "Estimating contamination via perplexity: Quantifying memorisation in language model evaluation," *arXiv preprint arXiv:2309.10677*, 2023.
- [38] V. Raychev, M. Vechev, and E. Yahav, "Code completion with statistical language models," in *Proceedings of the 35th ACM SIGPLAN conference on programming language design and implementation*, 2014, pp. 419–428.

- [39] A. Svyatkovskiy, Y. Zhao, S. Fu, and N. Sundaresan, "Pythia: Ai-assisted code completion system," in *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, 2019, pp. 2727–2735.
- [40] P. Oman and J. Hagemester, "Metrics for assessing a software system's maintainability," in *Proceedings Conference on Software Maintenance 1992*, 1992, pp. 337–344.
- [41] D. Coleman, D. Ash, B. Lowther, and P. Oman, "Using metrics to evaluate software system maintainability," *Computer*, vol. 27, no. 8, pp. 44–49, 1994.
- [42] <https://radon.readthedocs.io/en/latest/intro.html#cyclomatic-complexity>, 2012.
- [43] J. Bieri, "Cognitive complexity-simplicity and predictive behavior," *The Journal of Abnormal and Social Psychology*, vol. 51, no. 2, p. 263, 1955.
- [44] <https://pypi.org/project/cognitive-complexity/>, 2022.
- [45] S. Ren, D. Guo, S. Lu, L. Zhou, S. Liu, D. Tang, N. Sundaresan, M. Zhou, A. Blanco, and S. Ma, "Codebleu: a method for automatic evaluation of code synthesis," *CoRR*, vol. abs/2009.10297, 2020. [Online]. Available: <https://arxiv.org/abs/2009.10297>
- [46] K. Lee, D. Ippolito, A. Nystrom, C. Zhang, D. Eck, C. Callison-Burch, and N. Carlini, "Deduplicating training data makes language models better," *CoRR*, vol. abs/2107.06499, 2021. [Online]. Available: <https://arxiv.org/abs/2107.06499>
- [47] A. Z. Broder, "Identifying and filtering near-duplicate documents," in *Combinatorial Pattern Matching*, R. Giancarlo and D. Sankoff, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2000, pp. 1–10.
- [48] <https://github.com/johnbumgarner/wordhoard>, 2024.
- [49] J.-l. Gailly and M. Adler, "Zlib compression library," 2004.
- [50] Z. Zhang, J. Wen, and M. Huang, "Ethicist: Targeted training data extraction through loss smoothed soft prompting and calibrated confidence estimation," *arXiv preprint arXiv:2307.04401*, 2023.
- [51] D. Kocetkov, R. Li, L. Ben Allal, J. Li, C. Mou, C. Muñoz Ferrandis, Y. Jernite, M. Mitchell, S. Hughes, T. Wolf, D. Bahdanau, L. von Werra, and H. de Vries, "The stack: 3 tb of permissively licensed source code," *Preprint*, 2022.
- [52] "bigcode/starcoder," <https://huggingface.co/bigcode/starcoder>, 2023.
- [53] "bigcode/starcoderbase," <https://huggingface.co/bigcode/starcoderbase>, 2023.
- [54] L. Tunstall, N. Lambert, N. Rajani, E. Beeching, T. Le Scao, L. von Werra, S. Han, P. Schmid, and A. Rush, "Creating a coding assistant with starcoder," *Hugging Face Blog*, 2023, <https://huggingface.co/blog/starchat>.
- [55] "Huggingfaceh4/starchat-beta," <https://huggingface.co/HuggingFaceH4/starchat-beta/commit/4d8424c69643590f193cc97dc7eebf66500ebc6>, 2023.
- [56] "timdettmers/openassistant-guanaco," <https://huggingface.co/datasets/timdettmers/openassistant-guanaco>, 2023.
- [57] "Openassistant/oasst1," <https://huggingface.co/datasets/OpenAssistant/oasst1/tree/main>, 2023.
- [58] Z. Luo, C. Xu, P. Zhao, Q. Sun, X. Geng, W. Hu, C. Tao, J. Ma, Q. Lin, and D. Jiang, "Wizardcoder: Empowering code large language models with evol-instruct," *CoRR*, vol. abs/2306.08568, 2023. [Online]. Available: <https://doi.org/10.48550/arXiv.2306.08568>
- [59] "Wizardlm/wizardcoder," <https://huggingface.co/WizardLM/WizardCoder-15B-V1.0>, 2023.
- [60] S. Chaudhary, "Code alpaca: An instruction-following llama model for code generation," <https://github.com/sahil280114/codealpaca>, 2023.
- [61] "sahil280114/codealpaca," <https://github.com/sahil280114/codealpaca/commit/d269da106a579a623a654529b3cb91b5dfa9c72f>, 2023.
- [62] "codellama/codellama-7b-instruct-hf," <https://huggingface.co/codellama/CodeLlama-7b-Instruct-hf>, 2023.
- [63] "codellama/codellama-7b-instruct-hf commit," <https://huggingface.co/codellama/CodeLlama-7b-Instruct-hf/commit/65db8fcae13921e49c3dd0a2be4757102d0e723f>, 2023.
- [64] "meta-llama/llama-2-7b," <https://huggingface.co/meta-llama/llama-2-7b>, 2023.
- [65] "Phind/phind-codellama-34b-v2," <https://huggingface.co/Phind/Phind-CodeLlama-34B-v2>, 2023.
- [66] "Phind/phind-codellama-34b-v2," <https://huggingface.co/Phind/Phind-CodeLlama-34B-v2/commit/29c3be6006297754f344ba05678c038b0b77f6c0>, 2023.
- [67] Z. Du, Y. Qian, X. Liu, M. Ding, J. Qiu, Z. Yang, and J. Tang, "Glm: General language model pretraining with autoregressive blank infilling," in *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2022, pp. 320–335.
- [68] A. Zeng, X. Liu, Z. Du, Z. Wang, H. Lai, M. Ding, Z. Yang, Y. Xu, W. Zheng, X. Xia *et al.*, "Glm-130b: An open bilingual pre-trained model," *arXiv preprint arXiv:2210.02414*, 2022.
- [69] "Thudm/chatglm3-6b," <https://huggingface.co/THUDM/chatglm3-6b>, 2023.
- [70] "Thudm/chatglm3-6b," <https://huggingface.co/THUDM/chatglm3-6b/commit/62acdaa77a5742d120acf9b6419656b403218c3d>, 2023.
- [71] "Gpt-3.5-turbo model availability," <https://learn.microsoft.com/en-us/azure/ai-services/openai/concepts/models#gpt-35-turbo-model-availability>, 2023.
- [72] "Github copilot," <https://copilot.microsoft.com/>, 2023.
- [73] "Introducing github copilot: Ai pair programmer," <https://github.blog/2021-06-29-introducing-github-copilot-ai-pair-programmer/>, 2023.
- [74] "Github copilot november 30th update," <https://github.blog/changelog/2023-11-30-github-copilot-november-30th-update/>, 2023.
- [75] "Gpt-4 and gpt-4-turbo preview," <https://learn.microsoft.com/en-us/azure/ai-services/openai/concepts/models#gpt-4-and-gpt-4-turbo-preview>, 2023.
- [76] A. Holtzman, J. Buys, L. Du, M. Forbes, and Y. Choi, "The curious case of neural text degeneration," in *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020. [Online]. Available: <https://openreview.net/forum?id=rygGQyFvH>
- [77] S. Ouyang, J. M. Zhang, M. Harman, and M. Wang, "Llm is like a box of chocolates: the non-determinism of chatgpt in code generation," *arXiv preprint arXiv:2308.02828*, 2023.
- [78] "Github copilot is generally available to all developers," <https://github.blog/2022-06-21-github-copilot-is-generally-available-to-all-developers/>, 2022.
- [79] M. Kim, L. Bergman, T. Lau, and D. Notkin, "An ethnographic study of copy and paste programming practices in oopl," in *Proceedings. 2004 International Symposium on Empirical Software Engineering, 2004. ISESE '04.*, 2004, pp. 83–92.
- [80] S. Balloccu, P. Schmidová, M. Lango, and O. Dusek, "Leak, cheat, repeat: Data contamination and evaluation malpractices in closed-source LLMs," in *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, Y. Graham and M. Purver, Eds. St. Julian's, Malta: Association for Computational Linguistics, Mar. 2024, pp. 67–93. [Online]. Available: <https://aclanthology.org/2024.eacl-long.5>
- [81] T. Schick and H. Schütze, "It's not just size that matters: Small language models are also few-shot learners," *arXiv preprint arXiv:2009.07118*, 2020.
- [82] I. Magar and R. Schwartz, "Data contamination: From memorization to exploitation," *arXiv preprint arXiv:2203.08242*, 2022.
- [83] J. Mattern, F. Miresghallah, Z. Jin, B. Schölkopf, M. Sachan, and T. Berg-Kirkpatrick, "Membership inference attacks against language models via neighbourhood comparison," *arXiv preprint arXiv:2305.18462*, 2023.
- [84] J. Liu, C. S. Xia, Y. Wang, and L. Zhang, "Is your code generated by chatGPT really correct? rigorous evaluation of large language models for code generation," in *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. [Online]. Available: <https://openreview.net/forum?id=1qv610Cu7>
- [85] R. A. Poldrack, T. Lu, and G. Beguš, "Ai-assisted coding: Experiments with gpt-4," *arXiv preprint arXiv:2304.13187*, 2023.
- [86] D. Noever and K. Williams, "Chatbots as fluent polyglots: Revisiting breakthrough code snippets," *arXiv preprint arXiv:2301.03373*, 2023.