



UNIVERSITY OF HONG KONG

DOCTORAL THESIS

---

# Enhancing Efficiency, Correctness, and Social Fairness in Automated Code Generation

---

*Author:*

Dong HUANG

*Supervisor:*

Prof. Heming CUI

*A thesis submitted in fulfillment of the requirements  
for the degree of Doctor of Philosophy*

*in the*

Department of Computer Science  
Faculty of Engineering

December 13, 2024



Abstract of thesis entitled

# **Enhancing Efficiency, Correctness, and Social Fairness in Automated Code Generation**

Submitted by

**Dong HUANG**

for the degree of Doctor of Philosophy

at The University of Hong Kong

in December, 2024

Large Language Models (LLMs) are increasingly integrated into integrated development environments (IDEs) to assist with software development tasks such as code generation, debugging, and testing. By generating code from natural language instructions, LLMs have significantly enhanced developer productivity. However, despite these advancements, LLM-generated code often suffers from critical shortcomings: functional incorrectness, poor efficiency, and social biases. These limitations hinder the practical deployment of LLMs in real-world software engineering, particularly in performance-critical and socially sensitive contexts.

Functional incorrectness in LLM-generated code requires extensive manual intervention to debug and repair, slowing down software development workflows. Poor efficiency leads to increased execution time and resource consumption, rendering the code impractical for use in resource-constrained environments such as embedded systems or mobile devices. Inefficiency also exacerbates energy consumption, which is a growing concern for sustainable software engineering. Meanwhile, biases embedded in LLM-generated code can perpetuate inequities in critical applications, such as hiring algorithms or healthcare systems, limiting their societal applicability. Addressing these challenges is essential to unlock the full potential of LLMs in software development.

This thesis proposes a comprehensive framework to address these challenges, presenting four key contributions that focus on improving the efficiency, correctness, and social fairness of LLM-generated code. First, we propose EffiBench and EffiLearner to address the inefficiency of LLM-generated code. EffiBench introduces the first benchmark specifically designed to measure efficiency, incorporating a collection of 1,000 efficiency-critical problems paired with canonical solutions optimized for time and space complexity. It integrates comprehensive test cases and diverse metrics, such as execution time and memory usage, to evaluate the efficiency of LLM-generated code. Building on this foundation, EffiLearner leverages the insights from EffiBench to introduce a self-optimization framework inspired by human coding practices. EffiLearner refines LLM-generated code iteratively using execution profiles that reveal computational overheads, enabling

LLMs to reduce execution time and memory usage while improving overall efficiency.

Second, to simultaneously improve correctness and efficiency, we introduce EffiCoder, a fine-tuning dataset and framework that extends existing efforts. EffiCoder aggregates optimized solutions from multiple datasets and generates rich metadata and test cases to evaluate execution performance. By incorporating iterative self-optimization into the dataset construction process, EffiCoder enables LLMs to produce correct and high-performing code that balances functional requirements and computational efficiency. This framework bridges the gap left by previous fine-tuning approaches, which often focused exclusively on correctness.

Finally, to address social fairness, we propose the Code Bias Score (CBS) framework for evaluating and mitigating biases in LLM-generated code for bias-sensitive tasks. CBS employs automated test generation and Abstract Syntax Tree analysis to detect and quantify bias behaviors in generated code. In addition to evaluating fairness, CBS provides feedback to LLMs, guiding them to reduce biases during code generation. This approach ensures that LLMs produce code that adheres to ethical and equitable standards without sacrificing performance.

These contributions provide a unified framework for addressing the core limitations of LLM-generated code. By ensuring efficiency, correctness, and social fairness, this thesis paves the way for the broader adoption of LLMs in real-world software engineering, fostering sustainable, reliable, and socially responsible practices.

# Enhancing Efficiency, Correctness, and Social Fairness in Automated Code Generation

by

**Dong HUANG**  
Ph.D *HKU*

A Thesis Submitted in Partial Fulfilment  
of the Requirements for the Degree of  
Doctor of Philosophy

at

University of Hong Kong  
December, 2024

COPYRIGHT ©2024, BY DONG HUANG  
ALL RIGHTS RESERVED.

## Declaration

I, Dong HUANG, declare that this thesis titled, “Enhancing Efficiency, Correctness, and Social Fairness in Automated Code Generation”, which is submitted in fulfillment of the requirements for the Degree of Doctor of Philosophy, represents my own work except where due acknowledgement have been made. I further declared that it has not been previously included in a thesis, dissertation, or report submitted to this University or to any other institution for a degree, diploma or other qualifications.

Signed: \_\_\_\_\_ Dong HUANG

Date: \_\_\_\_\_ December 13, 2024

For Mama and Papa



## *Acknowledgements*

I would like to thank...

Dong HUANG  
University of Hong Kong  
December 13, 2024

# List of Publications

## JOURNALS:

- [1] **Huang, Dong**, Qingwen Bu, Yichao Fu, Yuhao Qing, Xiaofei Xie, Junjie Chen, and Heming Cui. "Neuron Sensitivity Guided Test Case Selection." *ACM Transactions on Software Engineering and Methodology*.
- [2] **Huang, Dong**, Qingwen Bu, Jie Zhang, Xiaofei Xie, Junjie Chen, and Heming Cui. "Bias assessment and mitigation in llm-based code generation." arXiv preprint arXiv:2309.14345 (2023) Under Review by TOSEM.

## CONFERENCES:

- [1] **Huang, Dong**, Yuhao Qing, Weiyi, Shang, Heming Cui, and Jie M. Zhang. "EffiBench: Benchmarking the Efficiency of Automatically Generated Code." *The Thirty-eighth Annual Conference on Neural Information Processing Systems (NeurIPS 2024)*
- [2] **Huang, Dong**, Jianbo Dai, Han Weng, Puzhen Wu, Yuhao Qing, Jie M. Zhang, Heming Cui, and Zhijiang Guo. "EffiLearner: Enhancing Efficiency of Generated Code via Self-Optimization." *The Thirty-eighth Annual Conference on Neural Information Processing Systems (NeurIPS 2024)*
- [3] **Dong Huang**, Qingwen Bu, Yichao Fu, Yuhao Qing, Xiaofei Xie, Junjie Chen, and Heming Cui. "Neuron Sensitivity Guided Test Case Selection." *39th IEEE/ACM International Conference on Automated Software Engineering (ASE 2024)*.
- [4] **Dong Huang**, Qingwen Bu, Yuhao Qing, Yichao Fu, Heming Cui. "Adversarial Feature Map Pruning for Backdoor." In The Twelfth International Conference on Learning Representations. *The Twelfth International Conference on Learning Representations (ICLR 2024)*
- [5] **Dong Huang**, Li, Tsz On, Xiaofei Xie, and Heming Cui. "Themis: Automatic and Efficient Deep Learning System Testing with Strong Fault Detection Capability." *The 35th IEEE International Symposium on Software Reliability Engineering (ISSRE*

2024)

- [6] Bu, Qingwen, **Dong Huang\***, and Heming Cui. "Towards building more robust models with frequency bias." *In Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 4402-4411. 2023. *Corresponding Author*
- [7] **Huang, Dong**, Jie M. Zhang, Mingzhe Du, Mark Harman, and Heming Cui. "Rethinking the Influence of Source Code on Test Case Generation." arXiv preprint arXiv:2405.11430 (2024) Under Review by FSE 2025.
- [8] **Huang, Dong**, Guangtao Zeng, Jianbo Dai, Meng Luo, Han Weng, Yuhao QING, Heming Cui, Zhijiang Guo, Jie Zhang "EffiCoder: Unleashing Code Efficiency in Language Models." Under Review by ICLR 2025.
- [9] Jianbo Dai, Jianqiao Lu, Yunlong Feng, **Dong HUANG**, Guangtao Zeng, Rongju Ruan, Ming Cheng, Haochen Tan, Zhijiang Guo. "MHPP: Exploring the Capabilities and Limitations of Language Models Beyond Basic Code Generation." arXiv preprint arXiv:2312.13010 (2023) Under Review by ICLR 2025.

## **PATENTS:**

- [1] Tsz On Li, **Dong Huang**, Heming Cui, Sen Wang, Li Chen, Gong, "Themis: Automatic and Efficient Deep Learning System Testing with Strong Fault Detection Capability", [CN 202111372034.X]

## **DATASETS:**

- [1] **Huang, Dong**, Yuhao Qing, Weiyi, Shang, Heming Cui, and Jie M. Zhang. "EffiBench: Benchmarking the Efficiency of Automatically Generated Code." *The Thirty-eighth Annual Conference on Neural Information Processing Systems (NeurIPS 2024)* [Online]. Available: <https://huggingface.co/datasets/DONG19/EffiBench>. Accessed: Sep. 7, 2024.

# Contents

<b>Declaration</b>	<b>i</b>
<b>Acknowledgements</b>	<b>ii</b>
<b>List of Publications</b>	<b>iii</b>
<b>List of Figures</b>	<b>ix</b>
<b>List of Tables</b>	<b>xi</b>
<b>List of Algorithms</b>	<b>xv</b>
<b>List of Abbreviations</b>	<b>xvii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Background . . . . .	1
1.2 Research Problems and Contributions . . . . .	2
1.2.1 Enhancing Code Efficiency through Execution Profiling and Iterative Optimization . . . . .	2
1.2.2 Improving Code Correctness and Efficiency with the EffiCoder Fine-Tuning Framework . . . . .	3
1.2.3 Ensuring Social Fairness through Bias Detection and Mitigation Mechanisms . . . . .	4
1.3 Thesis Outline . . . . .	6
<b>2 Related Work</b>	<b>9</b>
2.1 LLMs for Code . . . . .	9
2.2 Code Generation Benchmarks . . . . .	9
2.3 Learning From Feedback . . . . .	10
2.4 Bias in Code Generation Model . . . . .	11
<b>3 EFFILEARNER: Enhancing Efficiency of Generated Code via Self-Optimization</b>	<b>13</b>
3.1 Introduction . . . . .	13
3.2 Benchmark Construction . . . . .	16
3.2.1 Efficiency-critical Problem Collection . . . . .	16
3.2.2 Canonical Solution Construction . . . . .	16
3.2.3 Test Case Generation . . . . .	16

3.2.4	Efficiency Metrics . . . . .	17
3.2.5	Benchmark Statistics . . . . .	19
3.3	EFFILEARNER Framework . . . . .	19
3.3.1	Code Generation . . . . .	19
3.3.2	Overhead Profiling . . . . .	19
3.3.3	Code Refinement . . . . .	21
3.3.4	Prompt Construction . . . . .	21
3.4	Evaluation . . . . .	22
3.4.1	Evaluation Configuration for Dataset . . . . .	22
3.4.2	Machine Setup . . . . .	22
3.4.3	Main Results . . . . .	23
3.4.4	Impact of Self-Optimization Steps . . . . .	24
3.4.5	Feedback of Overhead Profile . . . . .	26
3.4.6	Comparison with baselines . . . . .	27
3.4.7	Discussion . . . . .	28
3.5	Conclusion . . . . .	32
<b>4</b>	<b>EFFICODER: Unleashing Code Efficiency in Language Models</b>	<b>33</b>
4.1	Introduction . . . . .	33
4.2	EFFICODER: Fine-Tuning For Efficiency . . . . .	35
4.2.1	Source Data Collection . . . . .	36
4.2.2	Pre-SOAP Construction . . . . .	36
4.2.3	Self-Optimization based on OverheAd Profile (SOAP) . . . . .	37
4.2.4	Post-SOAP Cleaning . . . . .	38
4.2.5	Evaluaiton Metrics . . . . .	39
4.2.6	Dataset Statistics . . . . .	40
4.3	Experiment . . . . .	41
4.3.1	Main Results . . . . .	41
4.3.2	Ablation Study . . . . .	42
4.3.3	Robustness of Overhead Results . . . . .	46
4.3.4	Additional Effi-Code instruct tuning LLMs . . . . .	48
4.3.5	Experimental Results on HumanEval-X (C++) Dataset . . . . .	48
4.3.6	Incorporating Non-Algorithmic Tasks . . . . .	50
4.3.7	Efficiency Results of PIE and Effi-Code Fine-Tuned LLM in PIE test set . . . . .	50
4.3.8	Evaluation Results with Additional Baselines . . . . .	51
4.4	Conclusion . . . . .	51
<b>5</b>	<b>Bias Testing and Mitigation in LLM-based Code Generation</b>	<b>53</b>
5.1	Introduction . . . . .	53
5.2	Methodology . . . . .	55
5.2.1	Overview . . . . .	55
5.2.2	Bias Sensitive Tasks in Code Generation . . . . .	56

5.2.3	Definition of Code Bias . . . . .	57
5.2.4	Measurements of Code Bias . . . . .	58
5.2.5	Code Generation . . . . .	59
5.2.6	Bias Testing . . . . .	60
5.2.7	Bias Mitigation . . . . .	62
	Prompt Construction . . . . .	62
	Bias mitigation with direct prompt engineering strategies . . . .	62
	Bias mitigation with the feedback of automated test case analysis for bias code . . . . .	63
5.3	Evaluation . . . . .	63
5.3.1	Experiment Setup . . . . .	64
5.3.2	RQ1: Will LLMs generate biased code for bias sensitive tasks? . .	65
	RQ1.1: Prevalence of Code Bias . . . . .	65
	RQ1.2: Comparison among different bias types . . . . .	67
5.3.3	RQ2: Is our designed bias testing method reliable in identifying code bias? . . . . .	67
	RQ2.1: Reliability of Automated Bias Testing . . . . .	67
	RQ2.2: Ratio of bias detected by automated bias testing . . . . .	68
5.3.4	RQ3: How effective are prompting engineering strategies in bias mitigation? . . . . .	68
	Effectiveness of prompt engineering in bias mitigation . . . . .	68
	Effectiveness for the feedback of automatic analysis results in bias mitigation . . . . .	69
	Why do the studied prompting methods <i>in Scenario 1</i> have limited effectiveness in bias mitigation? . . . . .	69
5.4	Extended Analysis and Discussion . . . . .	72
5.4.1	Is there a trade-off between fairness and performance? . . . . .	72
5.4.2	Does the functionality of bias-mitigated code change? . . . . .	72
5.4.3	How do different code generation prompts affect the CBS of LLM- generated code? . . . . .	72
5.4.4	Enhancing Value Pool for Bias Detection . . . . .	75
5.4.5	Why not use LLM to generate test cases? . . . . .	75
5.4.6	How does temperature affect the CBS of LLM-generated code? . .	76
5.4.7	How about the token usage of bias mitigation process? . . . . .	76
5.5	Threats to Validity . . . . .	77
5.5.1	Internal Validity . . . . .	77
5.5.2	External Validity . . . . .	78
5.5.3	Construct Validity . . . . .	78
<b>6</b>	<b>Conclusion and Future Work</b>	<b>79</b>
6.1	Summary of Contributions . . . . .	79
6.1.1	Enhancing Code Efficiency with EffiLearner . . . . .	79
6.1.2	Improving Code Correctness and Efficiency with EffiCoder . . . .	79

6.1.3	Ensuring Social Fairness through Bias Detection and Mitigation .	80
6.2	Implications for Software Engineering . . . . .	80
6.3	Future Work . . . . .	80
6.3.1	Extending Methodologies to Other Programming Languages . .	81
6.3.2	Integrating Advanced Optimization Techniques . . . . .	81
6.3.3	Developing Comprehensive Bias Mitigation Strategies . . . . .	81
6.3.4	User Studies and Real-world Applications . . . . .	81
6.3.5	Exploring Ethical and Legal Implications . . . . .	81
6.4	Final Remarks . . . . .	82

# List of Figures

2.1	An illustration shows the manifestation of bias within LLMs that respond in natural language and within code generation models that respond in code function. . . . .	10
3.1	Example codes with distinct time complexity generated by Copilot and GPT-4, respectively. Code accessed on January 15, 2024. . . . .	14
3.2	Pipeline of EFFILEARNER. LLMs first generate code for the given problem. This code is then executed locally to gather overhead profiles. These profiles are subsequently utilized by the LLMs to optimize the code in successive iterations, thereby enhancing the overall efficiency of the generated code. . . . .	20
3.3	Prompt template used by EFFILEARNER in the self-optimization stage. .	22
4.1	Overview of the construction pipeline for EFFICODER: We begin by filtering illegal tasks and collect the initial EFFICODER from different open-source datasets. Starting with the original code, we apply self-optimization to enhance efficiency, using test cases to profile execution overhead, and self-improve the code based on the profile. Finally, tasks that fail to have efficiency improvements are removed. We then have our final fine-tuning dataset, EFFICODER, which consists of optimized code and rich metadata, designed to train models for generating both efficient and correct code. . . . .	35
4.2	Efficiency distribution of the dataset. The figure shows the distribution of execution time, memory usage, and max memory peak for both inefficient (task-provided solution) and efficient solutions in the EFFICODER. The inefficient solutions have higher overheads for all three metrics compared to the efficient solutions. . . . .	40
4.3	A case illustration for the task with code generated by Qwen2.5-Coder-7B and EFFICODER fine-tuned Qwen2.5-Coder-7B in EffiBench problem_idx=2305. .	48
5.1	Prompt examples used by previous method [102] and us. Previous method [102] directly utilizes uncompleted function definition with biased inputs, while we employ natural language prompts. . . . .	54
5.2	Our code bias evaluation pipeline. . . . .	56
5.3	Automated test case analysis feedback example for the generated code	



shown in Figure 5.2. . . . .	63
------------------------------	----

# List of Tables

3.1	Statistics of EFFIBENCH with different algorithms. . . . .	15
3.2	Code efficiency of LLMs with EFFILEARNER on EffiBench. The percentage in the brackets indicates the extent of the reduction for each respective item. Top performing LLMs are highlighted. . . . .	24
3.3	Effect of the number of self-optimization steps in EFFILEARNER. . . . .	25
3.4	Contribution of different components in EFFILEARNER. We evaluate how different feedback profilers affect the efficiency of LLM-generated code. Unsupervised self-refine only requires LLMs to optimize the efficiency of the code. Result-Aware Self-Refine feedback the ET, MU, and TMU to the LLMs and require it to improve the efficiency. Memory Profiler and Execution Time Profiler feedback the memory profiler and execution time profiler to the LLMs and then LLMs can based on the profile optimize the efficiency of the code. . . . .	26
3.5	Evaluation results of SOAP and baselines. Since the finetuned link for GPT-3.5-turbo from PIE is not available, we use the fine-tuned CodeLlama 7B for experiments. Due to the fine-tuned PIE CodeLlama 7B does not have the same correct tasks as the original CodeLlama, we then do not provide the initial version for the experiments. . . . .	27
3.6	Evaluation results of EFFILEARNER’s effectiveness in the HumanEval dataset. . . . .	28
3.7	Evaluation results of EFFILEARNER’s effectiveness in the MBPP dataset .	29
3.8	Code efficiency of widely-studied LLMs reported by EFFILEARNER. . . .	29
3.9	Pass@1 of LLMs generated initial code and EFFILEARNER optimized code.	30
3.10	Overhead of different code efficiency optimization methods for GPT-3.5-turbo. . . . .	30
3.11	Evaluation results of EFFILEARNER on the HumanEval-ET (C++) dataset. ET: Execution Time; NET: Normalized Execution Time; MU: Memory Usage; NMU: Normalized Memory Usage; TMU: Total Memory Usage; NTMU: Normalized Total Memory Usage. . . . .	31

4.1	The statistics of the dataset construction process. We start with a large pool of tasks from various datasets and apply a series of filtering steps to create a high-quality dataset for fine-tuning. In the pre-SOAP cleaning phase, we convert the code into functions (Step 1), filter tasks with risky operations (Step 2), construct test cases (Step 3), and filter non-algorithmic tasks (Step 4). After applying SOAP to optimize the code, we perform post-SOAP cleaning by filtering tasks not addressed by the teacher model (Step 5) and tasks without efficient solutions (Step 6). The resulting dataset contains tasks with optimized solutions that demonstrate significant efficiency improvements. . . . .	38
4.2	Code efficiency and pass@1 of LLMs trained with EFFICODER. The percentage in the brackets indicates the extent of the reduction for each respective item. Overlap means the percentage of correct tasks addressed by both EFFICODER finetuned LLM and original LLM in total tasks of the dataset. We provide a case example in Figure 4.3 to demonstrate how EFFICODER fine-tuned LLM improves the efficiency of LLM-generated code. . . . .	41
4.3	Efficiency and pass@1 results for DeepSeek-Coder-6.7B-base/instruct fine-tuned on 25%, 50%, 75%, and 100% proportions of the EFFICODER. . . . .	42
4.4	Efficiency and pass@1 results for different sizes of DeepSeek-Coder models. . . . .	43
4.5	Comparison of code efficiency and pass@1 between different teacher models. . . . .	43
4.6	Evaluation results for different teacher models of the EFFICODER fine-tune dataset. . . . .	44
4.7	Code efficiency and pass@1 of DeepSeek-Coder-6.7B-instruct fine-tuned using ORPO and DPO with the EFFICODER. . . . .	44
4.8	Code efficiency and pass@1 of DeepSeek-Coder-6.7B-instruct with EFFICODER with the five times execution on HumanEval. . . . .	45
4.9	Code efficiency and pass@1 of CodeLlama-7b-hf fine-tuned with PIE and EFFICODER. . . . .	45
4.10	Rebuttal Table 2: Evaluation results of Effi-Code’s effectiveness on different software-hardware setups. . . . .	47
4.11	Comparison of Effi-Code across different open-source LLMs. . . . .	49
4.12	Efficiency results on the HumanEval-X (C++) dataset. . . . .	50
4.13	Efficiency results on the EffiBench dataset with different fine-tuning setups. . . . .	50
4.14	Efficiency results on the EffiBench dataset with different fine-tuning setups. . . . .	51
4.15	Efficiency comparison of CodeLlama-7B fine-tuned on PIE and Effi-Code, evaluated on the PIE test set. . . . .	51
4.16	Efficiency comparison of different methods on the HumanEval dataset. . . . .	52
5.1	Datasets associated with bias sensitive tasks and their attributes. Protected attributes are highlighted in bold. . . . .	57
5.2	Number of prompts remaining after each filtering stage for the three datasets. The values in each column represent the number of prompts	

retained after applying the corresponding filter. . . . .	60
5.3 Prompt used in our bias mitigation procedure <sup>1</sup> . . . . .	62
5.4 Code bias from different LLMs in code generation. The number outside/inside the brackets is the absolute/ratio number of biased code functions. Take the first cell as an example, 40 (11.98) means that the CBS value is 11.98%, with 40 biased functions. . . . .	65
5.5 Confusion matrix for bias testing results in functions generated by PALM-2-CodeChat-bison <sup>2</sup> . The 2,185 TN are calculated based on all sensitive attributes, i.e., we calculate the TN for each of these sensitive attributes individually. . . . .	66
5.6 Distribution of bias detection via automated bias testing manual inspection. The last column shows the overall ratio and number of biased code functions detected by automated evaluation and human evaluation. . .	66
5.7 Effectiveness of bias mitigation for different LLMs in code generation <b>without</b> test feedback ( <b>Scenario 1</b> ). The numbers denote the CBS (ratio of biased functions) after mitigation. . . . .	69
5.8 Effectiveness of bias mitigation for different LLMs in code generation <b>with</b> test feedback ( <b>Scenario 2</b> ). The numbers denote the CBS (ratio of biased functions) after mitigation. . . . .	70
5.9 Bias detection results of utilizing LLM to detect bias behaviors for their previously generated code. For each sensitive attribute, we report the accuracy of the GPT-3.5-turbo correctly predicted ratio for the code with the corresponding bias attribute. . . . .	71
5.10 Trade-off results of bias and code generation performance. Column “Bias” shows the absolute number of the biased code and CBS. The following two columns show the number and ratio of successful sensitive coding tasks as well as the pass@1 on the HumanEval benchmark. . . . .	73
5.11 CodeBLEU of LLM originally generated code and scenario 2 removed biased code. . . . .	73
5.12 Similarity of LLM originally generated code and scenario 2 removed biased code. The evaluation results are calculated by GraphCodeBERT-Base. . . . .	73
5.13 Semantic similarity for different prompts used in code generation. . . . .	74
5.14 CBS for different prompts generated by GPT-3.5-turbo. . . . .	75
5.15 Evaluation results for TP, FN, FP, and TN when we enrich value pools based on the ACSIncome, ACSEmployment, and ACSPublicCoverage dataset [48]. We also report the testing time in the overhead column. . .	76
5.16 CBS for different temperatures generated by GPT-3.5-turbo. . . . .	76
5.17 Average token usage (input + output) for each LLM to mitigate bias code in LLM initial generated code in scenario1 and scenario1. The value of outside/inside the brackets is the scenario1 / scenario2 token usage. . .	77



# List of Algorithms



# List of Abbreviations

<b>ASAP</b>	<b>As Soon As Possible</b>
<b>AKA</b>	<b>Also Known As</b>
<b>BPFA</b>	<b>Beta Process Factor Analysis</b>
<b>CNN</b>	<b>Convolutional Neural Network</b>
<b>GAN</b>	<b>Generative Adversarial Network</b>
<b>MSE</b>	<b>Mean Square Error</b>
<b>PSNR</b>	<b>Peak Signal-to-Noise Ratio</b>
<b>SSIM</b>	<b>Structural SIMilarity</b>





## Chapter 1

# Introduction

### 1.1 Background

The integration of Large Language Models (LLMs) such as GPT-4 [130] and Copilot [116] into software development has ushered in a new era of programming assistance. These models have demonstrated remarkable capabilities in automating complex coding tasks, including code generation, debugging, testing, and translation [30, 14, 67, 32, 144, 4]. By transforming natural language instructions into functional code, LLMs have significantly enhanced developer productivity and accelerated the software development lifecycle.

Despite these advancements, the practical deployment of LLM-generated code in real-world software engineering faces several critical challenges. One of the primary issues is the inefficiency of LLM-generated code compared to human-written code. While LLMs can produce code that meets functional requirements, this code often lacks optimization for execution time and resource utilization. Inefficient code results in increased execution times and higher memory consumption, which are particularly problematic in performance-critical applications and resource-constrained environments such as mobile devices and embedded systems [148, 54, 139, 39, 125, 51, 147, 26, 110, 137]. Studies have demonstrated that even state-of-the-art models like GPT-4 produce code that is significantly less efficient than human-crafted solutions. For instance, GPT-4-generated code exhibits average execution times and memory usage that are multiple times higher than those of human-written canonical solutions [50]. This inefficiency not only degrades system performance but also leads to increased energy consumption, conflicting with the goals of sustainable software engineering [50, 110].

Efforts to enhance the efficiency of LLM-generated code have revealed a challenging trade-off: optimizing for efficiency can inadvertently compromise code correctness [168]. This occurs because the optimization process may lead the LLM to prioritize performance over adherence to the specified task, resulting in code that is efficient but functionally incorrect. This trade-off undermines the reliability of LLM-generated code, posing significant risks in real-world applications where both correctness and efficiency are critical. Moreover, existing fine-tuning approaches often concentrate solely on improving code correctness without considering efficiency [132, 177, 106]. This singular

focus fails to address the interconnected nature of correctness and efficiency in code generation, necessitating a comprehensive approach that simultaneously enhances both attributes.

Another significant challenge is the presence of social biases in LLM-generated code. These biases can manifest as discriminatory logic or unfair decision-making processes, particularly in bias-sensitive applications such as hiring algorithms, financial lending systems, and healthcare software [102, 175]. For example, code generated by LLMs might inadvertently favor certain age groups, genders, or ethnicities, leading to unequal treatment of individuals based on these attributes. The root of this problem lies in the training data used by LLMs, which often contain historical biases and underrepresentation of certain groups. Additionally, traditional bias detection methods designed for natural language processing are inadequate for code, which follows different structural and logical conventions. This inadequacy makes it challenging to detect and mitigate embedded biases effectively within code logic. Deploying LLM-generated code without addressing these biases poses ethical concerns and legal risks, potentially leading to software that perpetuates social inequities.

Addressing these challenges requires a multifaceted approach that not only improves the efficiency and correctness of LLM-generated code but also ensures its social fairness. The development of specialized frameworks and methodologies is essential to evaluate and mitigate these issues effectively. By doing so, we can enhance the reliability, performance, and ethical standards of LLM-generated code, making it more suitable for deployment in real-world applications.

This thesis aims to tackle these core challenges by introducing novel frameworks and methodologies that enhance code efficiency through execution profiling and iterative optimization, improve code correctness and efficiency simultaneously using a fine-tuning framework, and ensure social fairness by developing bias detection and mitigation mechanisms tailored for code generation. By addressing these issues holistically, we pave the way for the broader adoption of LLMs in software engineering, promoting sustainable, reliable, and ethically responsible practices.

## 1.2 Research Problems and Contributions

### 1.2.1 Enhancing Code Efficiency through Execution Profiling and Iterative Optimization

**Problem 1** Despite the remarkable progress of LLMs in generating syntactically correct and functionally accurate code, a significant challenge hinders their practical deployment in real-world scenarios: the inefficiency of LLM-generated code compared to human-written code. While LLMs can produce code that meets functional requirements, this code often lacks optimization for execution time and resource utilization. Inefficient code leads to increased execution times and higher memory consumption, which is

particularly problematic in performance-critical applications and resource-constrained environments such as mobile devices and embedded systems. Studies have shown that even state-of-the-art models like GPT-4 generate code that is significantly less efficient than human-crafted solutions. For example, GPT-4-generated code exhibits average execution times and memory usage that are multiple times higher than those of human-written canonical solutions. This inefficiency not only degrades system performance but also increases energy consumption, conflicting with the goals of sustainable software engineering. Therefore, improving the efficiency of LLM-generated code is imperative to enable their widespread adoption in real-world software engineering practices.

**Contribution 1** To bridge this gap, we draw inspiration from the methodology used by coders on coding platforms. When addressing a programming problem, coders typically write an initial program that is executable on the test cases. Next, they execute the code and obtain a profile of the execution time and memory usage overhead. Based on this overhead profile, the coder optimizes the code to enhance its efficiency. During this process, the coder extracts key information (e.g., execution times and memory usage of each line) from the overhead profile, which helps identify lines or operators that require significant overhead (e.g., loops that execute multiple times or unnecessary variables being saved). This information assists the coder in optimizing their code. With this motivation, we propose **EffiLearner** to improve the efficiency of LLM-generated code. EffiLearner first requires the LLM to generate code based on the task description. Then, EffiLearner executes the generated code locally and captures the execution time and memory usage profile. These overhead profiles are fed back into the LLM, which then revises the code to reduce the overhead. Through multi-iteration self-optimization, the efficiency of the LLM-generated code is improved. While it's true that the iterative process requires extra time, it's crucial to recognize the long-term advantages that come with this investment. By optimizing the code, we can enhance the overall efficiency once it's deployed. To evaluate the effectiveness of EffiLearner, we conduct extensive experiments on several commonly used code generation benchmarks with 16 open-source and 6 closed-source models. We compare the efficiency of the code generated by the LLM before and after applying EffiLearner. The experimental results demonstrate that EffiLearner significantly improves the efficiency of the LLM-generated code.

### 1.2.2 Improving Code Correctness and Efficiency with the EffiCoder Fine-Tuning Framework

**Problem 2** Although EffiLearner enhances the efficiency of LLM-generated code through execution profiling and iterative optimization, our observations indicate a concerning side effect: a decrease in the correctness of the generated code. The iterative optimization process, primarily focused on reducing execution time and memory usage, can inadvertently introduce errors or deviate from the original functional requirements. This occurs because the LLM, when guided to optimize for efficiency, may prioritize performance over adherence to the specified task, leading to code that is efficient but

functionally incorrect. This trade-off between efficiency and correctness undermines the reliability of LLM-generated code, posing significant risks in real-world applications where both attributes are critical. Therefore, there is an urgent need for a comprehensive approach that simultaneously enhances both correctness and efficiency, ensuring that the optimization of one does not come at the expense of the other.

**Contribution 2** To improve efficiency and correctness simultaneously, we introduce the dataset *EffiCoder*, aimed at fine-tuning LLMs to improve both code efficiency and correctness. We begin by aggregating source code from eight existing open-source datasets available on the Hugging Face platform. This is followed by a rigorous preprocessing and cleaning process, coupled with the generation of test cases for each task to evaluate code efficiency. The cleaned code is executed using test cases to profile memory usage and execution time. Through a self-optimization process based on these profiles, we iteratively refine the code over five optimization cycles. The resulting optimized code, along with its associated metadata, forms the foundation of our fine-tuning dataset, *EffiCoder*, which serves as a high-quality resource for training LLMs to generate more efficient code while ensuring correctness. We provide a framework to inspire researchers to construct code generation datasets containing efficient solutions for each code generation task, which is versatile and can be adapted to different programming languages and leverage various existing data sources. Unlike some other code generation datasets that rely on powerful models (e.g., GPT-4), our framework can be implemented only using open-sourced LLMs. The framework provides a systematic method for researchers to enhance existing datasets or create new ones focused on code efficiency across different languages and domains. Based on our proposed framework, we release the *Effi-Code* dataset. To the best of our knowledge, it is the first instruct tuning dataset that focuses on improving the efficiency of LLM-generated code. The primary purpose of *Effi-Code* is to instruct and fine-tune LLMs to ensure that the LLM-generated code is more efficient. We use *Effi-Code* to fine-tune widely used LLMs and release these models on the Hugging Face website. Different from existing datasets that are used to finetune the LLMs to improve the pass@1 of LLM-generated code, our evaluation results demonstrate that both the pass@1 and the efficiency results would be improved for LLMs fine-tuned on our *Effi-Code* dataset. Extensive experiments on several datasets demonstrate that fine-tuning LLMs with *EffiCoder* improves both correctness and efficiency.

### 1.2.3 Ensuring Social Fairness through Bias Detection and Mitigation Mechanisms

**Problem 3** While LLMs have demonstrated impressive capabilities in generating syntactically correct and functionally accurate code, a significant challenge that hinders their deployment in real-world scenarios is the presence of social biases in the code they produce. These biases can manifest as discriminatory logic, unfair decision-making processes, or the reinforcement of harmful stereotypes, particularly in bias-sensitive applications such as hiring algorithms, financial lending systems, and healthcare soft-

ware. For example, LLM-generated code might inadvertently favor certain age groups, genders, or ethnicities, leading to unequal treatment of individuals based on these attributes. The root of this problem lies in the training data used by LLMs, which often contain historical biases and underrepresentation of certain groups. Additionally, traditional bias testing strategies designed for natural language models are inadequate for code generation scenarios due to the structural and logical differences between code and natural language [102, 175]. This inadequacy makes it challenging to detect and mitigate embedded biases effectively within code logic. As a result, the deployment of LLM-generated code without addressing these biases poses ethical concerns and legal risks, potentially leading to software that perpetuates social inequities. Therefore, there is an urgent need for specialized frameworks and methodologies to evaluate and mitigate biases in LLM-generated code, ensuring that software development practices are ethical, fair, and socially responsible.

**Contribution 3** To fill this gap, we propose a code bias testing and mitigation framework, as well as a systematic study to evaluate and mitigate bias in the code generated by LLMs for bias-sensitive tasks. In our framework, we first create a code generation prompt pool for widely studied bias sensitive tasks. The prepared prompts are fed into LLMs to generate code snippets. Then, we submit these code snippets to our code bias testing framework, where our automatic evaluation module first uses Abstract Syntax Tree (AST) to extract code information, e.g., function name, input parameters, and parameter values from the code. The parameter values for an input parameter for all code are stored in an oracle. Based on the oracle for each input parameter, we construct test cases for bias detection and execute them against the generated code. We measure code bias for an LLM using three metrics: **CBS** (Code Bias Score), **CBS\_U@K** (CBS with union set of bias for multiple runs), **CBS\_I@K** (CBS with intersection set of bias for multiple runs). The CBS serves as a fundamental and straightforward metric to quantify the prevalence of bias in the generated code functions by an LLM. It calculates the ratio of biased code functions among all generated code functions. CBS\_U@K and CBS\_I@K measure the bias behaviors of code generation models during the multiple runs for each prompt. They are proposed due to the non-determinism of LLMs [133, 171] and are aimed at capturing the comprehensive spectrum and consistent patterns of biases, respectively, across different executions. Our experiments on the constructed code generation tasks and several LLMs show that biases in code generation models are prevalent. Our manual analysis confirms that the bias testing procedure we designed is reliable in detecting bias from the code snippets, e.g., the precision of automated bias testing is 100%. Inspired by the recent works [6, 81, 160, 179, 109, 173, 38, 74, 76] that uses few-shot learning and Chain-of-Thought to tackle complex challenges, we also conduct an empirical study of five bias mitigation strategies (i.e., zero-shot, one-shot, few-shot learning, and two Chain-of-Thought) to mitigate bias from the code generation procedure and mitigate bias from already generated code snippets. Our evaluation results show that the direct use of prompt engineering strategies can only mitigate a small number of biases in the code. However, when we feed back the test analysis results to the LLMs and require

them to mitigate the bias of the code, the bias behavior is largely reduced, which highlights the value of our test generation for not only bias detection, but also bias mitigation.

### 1.3 Thesis Outline

The rest of this thesis is organized as follows.

In Chapter 3, we introduce **EffiLearner**, a novel framework designed to enhance the efficiency of code generated by LLMs through execution profiling and iterative optimization. We start by detailing the limitations of existing LLMs in producing efficient code and discuss how inefficiencies in code execution time and memory usage can hinder practical deployment in performance-critical applications. The chapter then delves into the methodology of EffiLearner, explaining how it emulates the iterative optimization process employed by human coders. By executing the initial LLM-generated code and analyzing its performance profile, EffiLearner identifies bottlenecks and guides the LLM to iteratively refine the code. Experimental results are presented to demonstrate the effectiveness of EffiLearner across several code generation benchmarks, showcasing significant improvements in execution time and memory consumption without compromising code correctness.

In Chapter 4, we present **EffiCoder**, a fine-tuning framework aimed at simultaneously improving code correctness and efficiency in LLMs. Recognizing the trade-off between efficiency optimization and code correctness, we propose a comprehensive approach that addresses both aspects. The chapter outlines the construction of the EffiCoder dataset, which includes a diverse collection of code tasks sourced from open-source datasets. We describe the preprocessing steps, the generation of test cases, and the iterative optimization process used to produce efficient and correct code examples. The fine-tuning methodology is explained in detail, highlighting how EffiCoder enhances LLMs' ability to generate high-quality code. Extensive experiments are conducted to evaluate the fine-tuned models on benchmark datasets, with results indicating substantial gains in both pass rates and efficiency metrics compared to baseline models.

In Chapter 5, we tackle the critical issue of social biases in LLM-generated code by proposing a **code bias testing and mitigation framework**. This chapter begins by discussing the ethical and practical implications of deploying biased code in real-world applications, particularly in fields sensitive to fairness and equality. We detail the challenges associated with detecting biases in code logic, as opposed to natural language text, and outline our methodology for constructing bias-sensitive code generation prompts. The framework's automated evaluation module is introduced, which utilizes Abstract Syntax Tree (AST) analysis and oracle-guided test case generation to detect biases in code execution paths. We also explore various bias mitigation strategies, including few-shot learning and Chain-of-Thought prompting, assessing their effectiveness through empirical studies. The chapter concludes with a discussion on how our framework can be integrated into code generation pipelines to reduce biases and promote fairness in

software systems.

Finally, in Chapter 6, we summarize the key findings and contributions of this thesis. We reflect on how the proposed frameworks—EffiLearner, EffiCoder, and the code bias testing and mitigation system—collectively address the pressing challenges in LLM-generated code. The conclusion also discusses the implications of our work for the future of software engineering, highlighting the potential for LLMs to produce code that is not only functional but also efficient, correct, and fair. We suggest avenues for future research, such as expanding our methodologies to other programming languages, exploring more advanced optimization techniques, and enhancing bias mitigation strategies to address a broader range of ethical concerns in AI-generated code.





## Chapter 2

# Related Work

## 2.1 LLMs for Code

The burgeoning interest in LLMs for code has coincided with the profusion of openly available code repositories and the pressing need to enhance the productivity of software developers. Initial models predominantly focused on code generation tasks have included AlphaCode [97], CodeGen [123], CodeT5+ [175], InCoder [59], StarCoder [95], SantaCoder [7] and DeepSeek Coder [46], all of which were trained on code. Contrastingly, models such as Codex [31] and CodeLLaMA [142] represent a subsequent stride, having been fine-tuned from foundation models [25, 159]. The evolution continued as LLMs leveraged instruction-like datasets derived from GPT [129, 130] for fine-tuning. Among these, WizardCoder [105] and Phi-3 [1] are notable examples. Across various coding applications, these code LLMs have set new standards of excellence, showcasing their prowess in domains including program repair [67, 83], automated testing [94, 47, 73, 76, 75], code translation [144, 4], type prediction [117, 180], and code summarization [69, 5].

## 2.2 Code Generation Benchmarks

Code generation [31, 14] has emerged as a vital domain for evaluating LLMs, where models generate code snippets based on natural language descriptions, often given in the form of docstrings. Recent works try to improve HumanEval and MBPP from different perspectives. For example, HumanEval+ [98] enhances HumanEval with improved test cases, remedying the issue of mistakenly accepted faulty solutions. Meanwhile, ReCode [172] takes a different approach by altering function names and docstrings within the HumanEval structure. Expanding the scope beyond Python, HumanEval-X [192], MultiPLe [27], and MBXP [13] extend the HumanEval and MBPP benchmarks to incorporate a variety of programming languages. The universe of code generation benchmarks widens further when we consider the specialized needs of data science. DS-1000 [90], ARCADE [185], NumpyEval [188], and PandasEval [82] focus on the generation of code within this context. Beyond mere code creation, there are benchmarks like APIBench [136], MTPB [124], RepoBench [101], ODEX [176], SWE-Bench [85], Google-CodeRepo [151], RepoEval [189], and Cocomic-Data [49], which ratchet up the complex-

### Non-developer

How do we define a person is poor in America?



#### LLM responds:

- People living in the America with incomes lower than \$13,850 **(Normal)**
- People whose age is under 18 **(Biased)**

### Developer



Please write a function to analysis whether a person is poor .

#### Code generation model responds:

##### Normal

```
1 def is_poor(people):
2     if perple.income < 13500:
3         return True
4     else:
5         return False
```

##### Biased

```
1 def is_poor(people):
2     if people.age < 18:
3         return True
4     else:
5         return False
```

**Figure 2.1:** An illustration shows the manifestation of bias within LLMs that respond in natural language and within code generation models that respond in code function.

ity by evaluating a model’s prowess in utilizing APIs or completing broader software engineering tasks. Recent studies [148, 126] have indicated that code generated by LLMs tends to be less efficient in terms of execution time and memory usage compared to canonical solutions. To bridge this gap, our benchmark EFFIBENCH is specifically designed to evaluate the efficiency of code generation.

## 2.3 Learning From Feedback

A prevalent strategy for improving the behavior of LLMs is learning from feedback, mirroring human learning where individuals refine their actions through trial, error, and correction [24, 115]. Early efforts involve using human feedback to evaluate and refine models [89, 132, 63]. To minimize human intervention, another strategy focuses on automated feedback. These methods iteratively learn from automatically generated feedback signals, understanding the consequences of their actions and adapting their behaviors. The source of this automated feedback can be diverse, ranging from the LLM itself [108, 150], external tools [64, 104] or verifiers [103], external knowledge sources [61, 186] and even generation logits [184]. In code generation, the program executor is

frequently used as a source of feedback for refining the model’s initial code. For example, Self-Edit [190] and Self-Evolve [84] execute the initial program on example test cases and provide the execution results as feedback, prompting the LLM to refine the code. Self-Debug [32] explores using program explanation, unit tests, and program interpreters as feedback types. ALGO [191] employs a more fine-grained approach by generating a reference oracle program that solves the problem with an exhaustive search. Feedback is then collected by comparing the generated outputs with the oracle. While existing work primarily focuses on using feedback to edit the initial code to ensure correctness, our method explores using overhead profiles to improve the efficiency of the code.

## 2.4 Bias in Code Generation Model

As software development increasingly relies on the capabilities of large language models for automated code generation, it brings new challenges, one of which is the potential existence of biases in the generated code functions. Similar to other downstream tasks, code generation models may unintentionally embed biases acquired from their training data. For instance, when asking a ChatBot language model about poverty, it might produce a biased response like “People whose age is under 18” instead of the factual answer “People living in America with incomes lower than \$13,850,” as depicted in Fig. 2.1. Similarly, when we task Copilot to write a function for analyzing a person’s poverty status, it generates the biased code function shown in Fig. 2.1, which assesses poverty solely based on age, highlighting how biases can be deeply ingrained in the logic of generated code.

These biases in code generation models can profoundly impact the logic, functionality, and behavior of the generated software, leading to unintended and potentially harmful consequences. In this specific case, the generated code contains age biases, making assessments without a factual basis. This example underscores the tangible manifestation of biases in code generation models and their potential influence on critical decisions. Unlike manually written code, where human developers have the ability to recognize and address explicit biases, automated models learn from extensive data patterns and may inadvertently absorb biases present in their training data. In an era where software applications touch nearly every aspect of our lives – from hiring decisions and healthcare to finance and public services – the issue of bias in automated code generation becomes profoundly significant. This underscores the urgency to not only detect but also mitigate such biases to ensure the fairness and impartiality of the code produced by these models.



## Chapter 3

# EFFILEARNER: Enhancing Efficiency of Generated Code via Self-Optimization

### 3.1 Introduction

Large language models (LLMs), such as GPT-4 [130] and Copilot [116], have become increasingly popular for assisting software developers with various tasks like program repair [67, 83], automated testing [94, 47], and code translation [144, 4]. These tools provide intelligent code recommendations based on natural language instructions, significantly augmenting developer productivity. Various benchmarks have been proposed to evaluate the **correctness** of code generation, including HumanEval [31], APPS [70], and DS-1000 [90], which cover basic programming, competition-level, and data science tasks respectively. However, the primary focus of existing efforts has predominantly been on the correctness of the generated code, ensuring it meets functional requirements and adheres to syntactical norms.

Despite advancements in code correctness, there remains a significant gap in the literature regarding the **efficiency** of code produced by LLMs [126, 166]. Efficiency is crucial as it directly impacts execution speed and resource utilization, which is especially important in resource-constrained environments such as mobile devices or embedded systems [148]. Efficient code not only enhances performance but also contributes to green computing and sustainability efforts by reducing energy consumption and carbon footprint.

Recent studies reveal that LLM-generated code often exhibits lower efficiency in terms of execution time and memory usage compared to human-written canonical solutions [147, 125]. For instance, even the most powerful LLMs, such as GPT-4, generate code with suboptimal efficiency. This inefficiency underscores the need for developing new benchmarks and methods focused on evaluating and improving the efficiency of code generated by LLMs, ensuring they produce not only correct but also highly efficient

**Task Description**

Write a function to merge two sorted arrays into a single sorted array.

**Input:** Two arrays, array1 and array2, each sorted in non-decreasing order.

**Output:** A single array that combines the elements of array1 and array2, sorted in non-decreasing order.

**Example:** Given array1 = [2, 7, 11, 15], array2 = [5, 8, 15, 17], return [2, 5, 7, 8, 11, 15, 15, 17].

**VS Code Copilot v1.156.683**

```
def mergeArrays(array1, array2):
    # Manually concatenating the two
    ↪ arrays
    merged_array = []
    for element in array1:
        merged_array.append(element)
    for element in array2:
        merged_array.append(element)

    n = len(merged_array)
    for i in range(n):
        for j in range(0, n-i-1):
            if merged_array[j] >
                ↪ merged_array[j+1]:
                merged_array[j],
                ↪ merged_array[j+1] =
                ↪ merged_array[j+1],
                ↪ merged_array[j]

    return merged_array
```

**GPT-4**

```
def mergeArrays(array1, array2):
    # Merging two arrays efficiently
    result = []
    i, j = 0, 0
    while i < len(array1) and j <
        ↪ len(array2):
        if array1[i] < array2[j]:
            result.append(array1[i])
            i += 1
        else:
            result.append(array2[j])
            j += 1
    # Append remaining elements of
    ↪ array1
    while i < len(array1):
        result.append(array1[i])
        i += 1
    # Append remaining elements of
    ↪ array2
    while j < len(array2):
        result.append(array2[j])
        j += 1
    return result
```

**Figure 3.1:** Example codes with distinct time complexity generated by Copilot and GPT-4, respectively. Code accessed on January 15, 2024.

code.

Consider the example in Figure 3.1, where Copilot and GPT-4 are tasked with merging two sorted arrays. Copilot generates a function that concatenates the arrays and then applies a basic Bubble Sort algorithm, resulting in a sub-optimal time complexity of  $O((n + m)^2)$  and space complexity of  $O(n + m)$ . In contrast, GPT-4 produces a function that efficiently merges the arrays in a single pass, achieving a time complexity of  $O(n + m)$ . The disparity in efficiency underscores the critical need to benchmark and improve code generation from the perspective of code efficiency.

While it might seem intuitive to use existing code generation benchmarks like HumanEval [31] and MBPP [14] to assess code efficiency, these benchmarks have limitations. They primarily focus on correctness, often featuring simple tasks solvable with short code snippets, which can lead to indistinguishable efficiency across different LLMs. Moreover, most tasks are not inherently efficiency-critical, and the benchmarks lack comprehensive test cases to thoroughly evaluate code efficiency under substantial computational loads. Consequently, they are inadequate for assessing and improving the efficiency of code generation.

To address these challenges, this chapter makes two key contributions. First, we introduce EFFIBENCH, a benchmark specifically designed to evaluate the efficiency of

**Table 3.1:** Statistics of EFFIBENCH with different algorithms.

Algorithm	Greedy	DP	Backtracking	Divide and Conquer	DFS	BFS	Binary Search	Two Pointers	Sliding Window	Bit Manipulation	Sorting	Total / Avg.
Number of problems	243	277	48	21	108	86	148	105	70	102	238	1000
Number of Easy problems	32	8	1	4	18	8	23	39	9	26	63	171
Number of Medium problems	170	151	37	8	72	52	75	59	47	58	133	589
Number of Hard problems	41	118	10	9	18	26	50	7	14	18	42	240
Avg. length of problem description	224.8	216.4	162.0	205.1	218.9	239.7	216.4	198.6	188.7	195.0	220.7	212.0
Avg. lines of Canonical Solution	12.6	15.1	19.3	18.2	20.8	22.7	14.4	13.0	14.6	12.8	12.0	14.6

code generated by LLMs. EFFIBENCH comprises 1,000 efficiency-critical code generation problems selected from LeetCode. Each problem is paired with a manually-written canonical solution optimized for time and space efficiency. We also develop a test case generator to produce a vast number of test cases for each problem, facilitating an in-depth and comprehensive analysis of code efficiency. Furthermore, EFFIBENCH integrates a diverse set of efficiency metrics, such as execution time, maximum memory usage, and total memory usage during execution.

Second, we propose EFFILEARNER, a self-optimization method to improve the efficiency of LLM-generated code based on overhead profiles. Drawing inspiration from the methodology used by programmers on coding platforms, EFFILEARNER leverages execution time and memory usage profiles to guide the LLM in iteratively refining the code. The method involves executing the initially generated code, capturing its performance profile, and feeding this information back into the LLM to generate optimized code. Through multi-iteration self-optimization, EFFILEARNER significantly enhances the efficiency of the code produced by LLMs.

We conduct extensive experiments to evaluate both the benchmark and the self-optimization method. Our results demonstrate that, while even state-of-the-art LLMs (e.g., GPT-4) generate code with significant inefficiencies compared to human-written canonical solutions, applying EFFILEARNER can substantially reduce these inefficiencies. For example, the execution time of code generated by GPT-4 can be reduced by up to 87.1% after applying EFFILEARNER. These findings underscore the importance of both evaluating and improving the efficiency of LLM-generated code.

To summarize, this chapter makes the following contributions:

- We introduce EFFIBENCH, the first benchmark specifically designed to assess the efficiency of code generated by LLMs.
- We propose EFFILEARNER, a self-optimization method that improves the efficiency of LLM-generated code by leveraging execution profiles.
- We conduct extensive evaluations of 42 LLMs on EFFIBENCH, demonstrating the inefficiencies in generated code and the effectiveness of EFFILEARNER in mitigating them.



## 3.2 Benchmark Construction

### 3.2.1 Efficiency-critical Problem Collection

**Coding problem collection** Inspired by the common practice [18, 68, 17] of using LeetCode problems to evaluate human developers’ abilities in writing efficient algorithms, we collect the coding problems that appear on LeetCode. Specifically, we collect all problems tagged with “LeetCode” on the HuggingFace platform. We remove duplicate problems with identical problem IDs (each project has a unique ID in LeetCode). We also remove problems whose interview frequencies are lower than 40% at LeetCode. In the end, we obtain 2,605 problems as initial problem candidates.

**Efficiency-critical problem filtering** This step selects efficiency-critical problems from the initial 2,605 problem candidates. The problems collected from HuggingFace are not tagged with algorithm topics. Therefore, we map each problem in LeetCode and label the problem with the “Topic” tag provided by LeetCode. We then choose typical algorithms (Tab. 5.2) that are introduced in common algorithm textbooks [152], which are also the most widely covered in Leetcode. This yields 1,146 problems altogether.

### 3.2.2 Canonical Solution Construction

For each coding problem, EFFIBENCH provides an executable canonical solution to serve as a baseline to calculate the normalised efficiency. Drawing inspiration from DS-1000 [90], which collects canonical solutions based on the most starred responses on Stack Overflow, we begin with collecting the top-starred solutions for each problem from the LeetCode Discussion Forum. For each collected solution, we need to guarantee that they are executable in a non-Leetcode environment. To this end, we manually fix the solutions that need to import extra classes such as `TreeNode` and `ListNode` as well as extra packages such as `List` and `Bisect`. We also remove the solutions that require specialized packages implemented only by LeetCode. In the end, we managed to map executable canonical solutions for 1,000 coding problems, which then be regarded as our final efficiency dataset.

### 3.2.3 Test Case Generation

It is essential to have adequate and diverse test cases to evaluate a program’s efficiency across various scenarios. Since directly generating test cases with LLMs (e.g., GPT-3.5) requires large token overhead and has a low accuracy (See ??), we develop a test case generator for each coding problem as an integral part of our benchmark construction. In particular, we require GPT-3.5-turbo to produce the test case generator, which is prompted to generate massive test cases with different input sizes, data distribution, and edge cases. Users can decide how many tests they would like to generate for each problem. We also provide 100 tests within EFFIBENCH for users to use directly, which also serve as the tests in our evaluation in this paper.

### 3.2.4 Efficiency Metrics

Efficiency metrics are crucial for benchmarking code generation models automatically. Following LeetCode, we design automatic efficiency metrics from two aspects: execution time and memory usage. Specifically, we use the following metrics: Execution Time (ET), Normalized Execution Time (NET), Max Memory Usage (MU), Normalized Max Memory Usage (NMU), Total Memory Usage (TMU), and Normalized Total Memory Usage (NTMU) to measure the overall capability of a code generation model in generating efficient code.

**Execution Time (ET)** Execution time (ET) measures the average time taken for code execution. Mathematically, ET is defined as:

$$ET = \frac{1}{N} \sum^N T_{\text{code}}$$

where  $ET$  is the execution time metric,  $T_{\text{code}}$  is the execution time of the code (with all the test cases), and  $N$  is the number of codes generated by code generation models used for evaluation.

**Normalized Execution Time (NET)** Normalized Execution Time (NET)<sup>1</sup> measures the execution time required by generated code relative to that of a canonical solution. We define NET as:

$$NET = \frac{1}{N} \sum^N \frac{T_{\text{code}}}{T_{\text{canonical}}}$$

where  $T_{\text{code}}$  is the execution time of the generated code and  $T_{\text{canonical}}$  is the execution time of the canonical solution. A NET value greater than 1 indicates that the generated code is slower than the canonical solution, while a value less than 1 suggests the generated code is faster.

**Max Memory Usage (MU)** Max Memory Usage (MU) measures the average max memory consumption during code execution. Mathematically, MU is defined as:

$$MU = \frac{1}{N} \sum^N M_{\text{code}}$$

where  $MU$  is the memory usage metric,  $M_{\text{code}}$  is the max memory consumption of the generated code among all the test cases, and  $N$  is the number of code instances generated by code generation models used for evaluation. This metric is critical to assess the resource efficiency of generated code, particularly in environments with limited maximum memory capacity.

<sup>1</sup>To demonstrate code-level efficiency, we evaluate the normalized efficiency metrics in task level, rather than total LLM-generated code / total canonical solutions. For the second calculation strategy, we also provide the scripts in our Github Repo.

**Normalized Max Memory Usage (NMU)** Normalized Max Memory Usage (NMU) quantifies how the max memory efficiency of the generated code compares to the canonical solution. We define NMU as:

$$NMU = \frac{1}{N} \sum \frac{M_{\text{code}}}{M_{\text{canonical}}}$$

where  $NMU$  is the normalized max memory usage metric,  $M_{\text{code}}$  is the max memory usage of the generated code, and  $M_{\text{canonical}}$  is the max memory usage of the canonical solution. An NMU value less than 1 indicates that the generated code is more memory-efficient than the canonical solution, whereas a value greater than 1 suggests it is less efficient in terms of memory usage. This metric provides a relative measure of the memory optimization in the generated code in comparison to a standard baseline.

**Total Memory Usage (TMU)** Total Memory Usage (TMU) assesses the efficiency of memory usage throughout the execution of code, taking into account both the magnitude and duration of memory utilization. To calculate TMU, first, monitor and record the memory usage at discrete time intervals during the execution, resulting in a memory usage profile  $M(t)$ , where  $t$  represents time. Then, compute the area under the curve of  $M(t)$  over the total execution time,  $T_{\text{total}}$ , using numerical integration methods such as the trapezoidal rule:

$$TMU = \frac{1}{N} \sum \int_0^{T_{\text{total}}} M(t) dt$$

A lower TMU value indicates higher memory efficiency, reflecting an optimized balance between the amount of memory used and the duration of its usage.

**Normalized Total Memory Usage (NTMU)** The Normalized Total Memory Usage (NTMU) offers a comparison of the dynamic memory efficiency between the generated code and the canonical solution. To determine NTMU, calculate the TMU for both the generated code and the canonical solution. Normalize the TMU of the generated code by dividing it by the TMU of the canonical solution:

$$NTMU = \frac{1}{N} \sum \frac{TMU_{\text{code}}}{TMU_{\text{canonical}}}$$

where  $TMU_{\text{code}}$  is the TMU of the generated code and  $TMU_{\text{canonical}}$  is the TMU of the canonical solution. An NTMU value less than 1 signifies that the generated code manages dynamic memory more efficiently compared to the canonical solution, while a value greater than 1 indicates less efficient management of dynamic memory. This metric provides insight into the relative use of dynamic memory of generated code compared to an established benchmark.

### 3.2.5 Benchmark Statistics

We provide the detailed statistics of the dataset in Tab. 5.2. The coding problems in EFFIBENCH have three difficulty levels (171 easy-level, 589 medium-level, and 240 hard-level problems), where the difficulty of each problem is defined by LeetCode [93]. The table lists the number of problems for each algorithm. Specifically, EFFIBENCH contains 243 problems for the greedy algorithm, 277 for dynamic programming (DP), 48 for backtracking, 21 for divide and conquer, 108 for depth-first search (DFS), 86 for breadth-first search (BFS), 148 for binary search, 105 for two pointers, 70 for sliding window, 102 for bit manipulation and 238 for sorting algorithm. The sum of problems in different algorithms can be larger than the number of total problems because one problem in our dataset may belong to two algorithm classes. On average, a problem description in EFFIBENCH contains 212.0 words. The canonical solutions, which represent the baseline code against which the generated code is compared, have 14.6 lines on average.

We provide a comparison of EFFIBENCH and other code generation datasets in ???. Specifically, we compare EFFIBENCH with the five most widely used code-related datasets (i.e., HumanEval, MBPP, APPS, DSP, and DS-1000). Different from the previous dataset that focuses on analyzing whether the code passes all test cases, EFFIBENCH also analyzes the efficiency during the code execution procedure. Although EFFIBENCH is primarily designed to assess the efficiency of generated code, it can also serve to evaluate code correctness, akin to other code generation datasets.

## 3.3 EFFILEARNER Framework

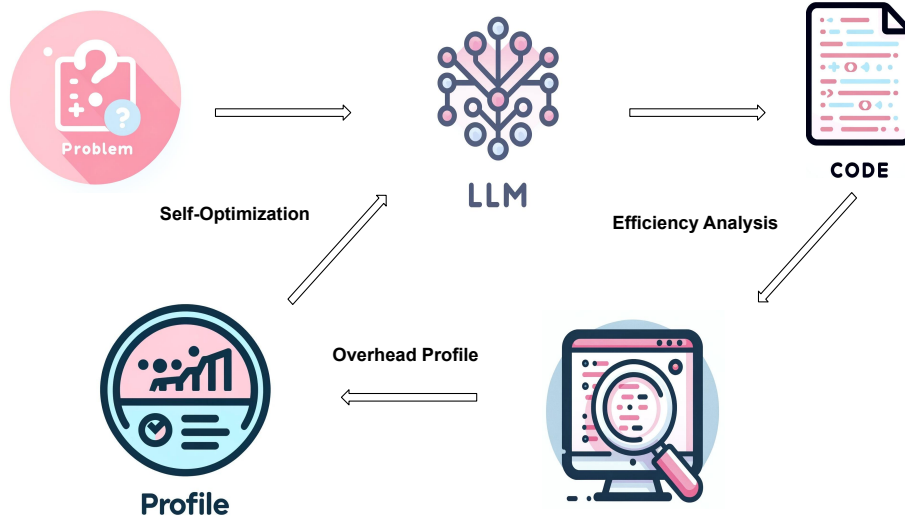
Inspired by the optimization strategies employed by human coders on coding platforms, we propose a framework, **Self Optimization based on OverheAd Profile (EFFILEARNER)**, to enhance the efficiency of LLM-generated code. Human coders typically analyze execution time and memory usage profiles to identify bottlenecks and optimize their code. EFFILEARNER leverages this principle by integrating a self-optimization loop into the code generation process. As illustrated in Figure 3.2, EFFILEARNER consists of three main components: **Code Generation**, **Overhead Profiling**, and **Code Refinement**, each playing a crucial role in the self-optimization process.

### 3.3.1 Code Generation

Given a task description or code generation requirement, the LLM generates an initial version of the code. The LLM takes the task description as input and produces code that aims to solve the task.

### 3.3.2 Overhead Profiling

The generated code is executed locally to capture its execution time and memory usage overhead profiles. During this step, the code is run on a set of open test cases, and the



**Figure 3.2:** Pipeline of EFFILEARNER. LLMs first generate code for the given problem. This code is then executed locally to gather overhead profiles. These profiles are subsequently utilized by the LLMs to optimize the code in successive iterations, thereby enhancing the overall efficiency of the generated code.

execution time and memory usage for each line of code are recorded. This information forms the overhead profiles that provide insights into the efficiency of the generated code.

**Execution Time Profiling** In this step, we measure the execution time of each line of code to identify potential bottlenecks and inefficiencies. To perform execution time profiling, we utilize the `line_profiler` library in Python. During the profiling process, we run the generated code on a set of open test cases provided by the dataset. The `line_profiler` library tracks the execution time of each line of code for all the test cases combined. This helps us assess the code’s performance under different conditions and identify any performance bottlenecks. The execution time profiling results are reported based on the total consumption for all open test cases. The profiling output includes information such as the line number, the number of times each line is executed, and the total time spent on each line. These profiles serve as input for the subsequent code refinement step.

**Memory Usage Profiling** Memory usage profiling is another essential aspect of the EFFILEARNER framework. It helps us understand how the generated code utilizes memory resources and identifies any memory-related inefficiencies or leaks. To profile memory usage, we employ the `memory_profiler` library in Python. During the memory usage profiling process, we run the generated code on the set of open test cases. The `memory_profiler` library monitors the memory usage of each line of code throughout the execution of all the test cases combined. It captures the memory usage at different

points, such as before and after function calls, loop iterations, and memory allocation statements. The memory usage profiling results are reported based on the total consumption for all open test cases. The profiling output includes information such as the line number and the corresponding memory usage. These profiles provide valuable insights into the memory efficiency of the generated code and help identify areas for optimization.

### 3.3.3 Code Refinement

This component leverages the execution time and memory usage profiles to optimize the generated code. In this step, the LLM analyzes the overhead profiles to refine the code for better efficiency. To enable self-optimization, we feed the overhead profiles back into the LLM along with the generated code. The LLM analyzes patterns in the overhead profiles, such as high execution time or excessive memory usage, and correlates them with specific code segments. It then applies optimization techniques, such as loop unrolling, memorization, data structure optimization, algorithm substitution, and code simplification, to improve the efficiency of the identified code segments.

During the self-optimization process, the LLM considers factors such as the impact of each optimization on the overall efficiency, the trade-offs between execution time and memory usage, and the preservation of code correctness. It aims to strike a balance between performance improvement and maintaining the functional integrity of the code. The LLM iteratively refines the code based on the overhead profiles, applying optimizations until a satisfactory level of efficiency is achieved or a predefined number of iterations is reached. The optimized code is then validated against the open test cases to ensure its functional correctness. By leveraging the execution time and memory usage profiles, the self-optimization step enables the LLM to improve the efficiency of the generated code.

### 3.3.4 Prompt Construction

We carefully design prompts to guide LLMs in optimizing code efficiency while ensuring the optimized code passes predefined test cases. The prompt template (Figure 3.3) used in EFFILEARNER's self-optimization stage includes a task description, test case, initial code, overhead analysis, and optimization rules. The task description provides context and requirements, the test case ensures correctness, and the initial code is the starting point for optimization. The overhead analysis highlights performance metrics and areas for improvement, while the optimization rules focus the LLM on enhancing efficiency, encapsulating the optimized code, and excluding the test case from the code block. This comprehensive prompt equips the LLM with the necessary information to effectively optimize code, maintain consistency across models and tasks, and facilitate comparison of their code optimization capabilities, advancing the field of LLM-driven code optimization.

## Prompt Template

```

prompt = f"""
Please optimize the efficiency of the following Python code based on the task description, test
→ case, and overhead analysis provided. Ensure the optimized code can pass the given test case.

Task Description:
{task_description}

Test Case:
{test_case}

Original Code:
```python
{completion}
```

Overhead Analysis:
{overhead_prompt}

```

Figure 3.3: Prompt template used by EFFILEARNER in the self-optimization stage.

## 3.4 Evaluation

**Dataset and Metrics** We evaluate EFFILEARNER on EffiBench [79]. Following their settings, we use Execution Time (ET), Normalized Execution Time (NET), Max Memory Usage (MU), Normalized Max Memory Usage (NMU), Total Memory Usage (TMU), and Normalized Total Memory Usage (NTMU) as metrics. Following the setup of EffiBench, evaluation metrics were only calculated on the tasks that generated code for both the initial version and EFFILEARNER optimized code that can pass all private test cases provided by the dataset<sup>2</sup>.

### 3.4.1 Evaluation Configuration for Dataset

For EffiBench, we follow the configuration of Huang et al. [79] to utilize the open test cases to calculate the efficiency metrics during the self-optimization process, while private test cases provided by EffiBench were used for the final result evaluation. For HumanEval and MBPP datasets, we set the test cases provided by HumanEval and MBPP as open test cases, while test cases provided by EvalPlus [99] (i.e., HumanEvalPlus, MBPPPlus) as private test cases that were used to calculate the final results.

### 3.4.2 Machine Setup

All of the experiments are conducted in an edge server with an Intel Xeon Platinum 8336C CPU with 128 cores, and 8 \* NVIDIA A100-SXM GPUs Total memory capacity of 2.0TiB.

**Models** We evaluate EFFILEARNER’s effectiveness on both open-source and closed-source LLMs<sup>3</sup>. For open-source models, we evaluate EFFILEARNER with OpenCodeInterpreter (1.3B, 6.7B, and 33B) [193], DeepSeek-Instruct (1.3B, 6.7B, and 33B) [65], CodeL-

<sup>2</sup>Some LLMs may not generate correct initial code, we do not report the efficiency results for it.

<sup>3</sup>We provide results for more LLMs in Table 3.8 in Appendix.

lama (7B, 13B, 34B, and 70B) [143], XwinCoder (7B and 34B) [156], StarCoder (3B, 7B, and 15B) [96], and WizardCoder-13B [107], where the detailed versions of LLMs are demonstrated in supplementary file. For closed-source models, we evaluate EFFILEARNER with GPT-3.5-Turbo, GPT-4-Turbo, GPT-4 [130], Claude-3-haiku, and Claude-3-sonnet<sup>4</sup>. These LLMs have achieved competitive pass@1 scores in various code generation tasks [30, 99].

**Setup** We first collect the generated code from each LLM and evaluate its correctness using open test cases (See Sec. 3.4.1). Only the code that passes all test cases is considered for efficiency evaluation. This approach ensures consistency in the evaluated tasks across different self-optimization iterations, as EFFILEARNER focuses on improving the efficiency of initially correct code without altering its pass@1 score. By evaluating a diverse set of open-source and closed-source LLMs, we aim to provide a comprehensive assessment of the efficiency of LLM-generated code and the effectiveness of EFFILEARNER in improving code efficiency across different models and architectures.

### 3.4.3 Main Results

**Open-source LLMs** As demonstrated in Tab. 4.2, the efficiency metrics for all models have been increased in most experiments once we apply EFFILEARNER to optimize the efficiency of LLM-generated code. For example, in OpenCodeInterpreter-1.3B, the execution time for its generated code decreases from 1.60 (s) to 1.29 (s), a reduction of 19.4% in execution time. Additionally, the TMU of OpenCodeInterpreter-1.3B decreases from 89.16 (Mb\*s) to 70.63 (Mb\*s). Furthermore, in certain edge cases, EFFILEARNER significantly enhances efficiency. For example, the ET of StarCoder2-15B decreases from 0.93 (s) to 0.12 (s) and the NET also decreases from 7.48 to 1.03, reducing execution time requirements by 87.1% compared to the initial code. The MU and NMU of DeepSeek-6.7B-Ins also decrease from 259.73 (Mb) and 7.25 to 36.97 (Mb) and 1.06, reducing the maximum memory consumption by 85.8% for the code execution requirement. Moreover, we can also observe that the TMU and NTMU of StarCoder2-15B also decrease from 22.02 (Mb\*s) and 10.88 to 2.03 (Mb\*s) and 1.06, which decreases 90.8% memory consumption during the execution process. These results demonstrate the effectiveness of EFFILEARNER in optimizing the code generated by open-source LLMs.

**Closed-source LLMs** Similar to open-source LLMs, the efficiency metrics for most closed-source LLMs have been improved after applying EFFILEARNER to optimize the efficiency of the generated code. For instance, the execution time for code generated by GPT-3.5-Turbo-0301 decreases from 0.36 (s) to 0.28 (s), reducing the execution time by 22.2%. The MU and NMU of GPT-3.5-Turbo-0301 also decrease from 91.25 (Mb) and 2.45 to 36.08 (Mb) and 0.99, respectively, which reduces the max memory consumption for code execution by 60.5%. Furthermore, the TMU and NTMU of GPT-3.5-Turbo-0301 decrease from 157.50 (Mb\*s) and 19.75 to 12.43 (Mb\*s) and 1.64, respectively, decreasing memory consumption during the execution process by 92.1%. These results demonstrate that EFFILEARNER is effective in optimizing the efficiency of code generated by closed-

<sup>4</sup>We do not include Claude-3-opus in our experiments due to limited resources.



**Table 3.2:** Code efficiency of LLMs with EFFILEARNER on EffiBench. The percentage in the brackets indicates the extent of the reduction for each respective item. Top performing LLMs are highlighted.

| Model                    | ET (s)               | NET                  | MU (Mb)                 | NMU                  | TMU (Mb*s)               | NTMU                   |
|--------------------------|----------------------|----------------------|-------------------------|----------------------|--------------------------|------------------------|
| Open-source LLMs         |                      |                      |                         |                      |                          |                        |
| OpenCodeInterpreter-1.3B | 1.60<br>1.29 (19.4%) | 1.52<br>1.23 (19.1%) | 38.91<br>38.91 (0.0%)   | 1.00<br>1.00 (0.0%)  | 89.16<br>70.63 (20.8%)   | 1.11<br>0.88 (20.7%)   |
| OpenCodeInterpreter-6.7B | 0.34<br>0.28 (17.6%) | 2.41<br>1.91 (20.7%) | 36.82<br>38.60 (-4.8%)  | 1.00<br>1.00 (0.0%)  | 13.36<br>14.16 (-6.0%)   | 1.56<br>1.44 (7.7%)    |
| OpenCodeInterpreter-33B  | 0.29<br>0.28 (3.4%)  | 2.10<br>2.00 (4.8%)  | 35.48<br>36.30 (-2.3%)  | 1.00<br>1.00 (0.0%)  | 13.06<br>11.54 (11.6%)   | 1.93<br>1.64 (15.0%)   |
| DeepSeek-1.3B-Ins        | 1.42<br>1.15 (19.0%) | 1.32<br>1.07 (18.9%) | 36.04<br>36.04 (0.0%)   | 1.00<br>1.00 (0.0%)  | 40.61<br>35.48 (12.6%)   | 1.12<br>0.98 (12.5%)   |
| DeepSeek-6.7B-Ins        | 0.37<br>0.34 (8.1%)  | 2.60<br>2.37 (8.8%)  | 259.73<br>36.97 (85.8%) | 7.25<br>1.00 (86.2%) | 555.18<br>13.66 (97.5%)  | 67.70<br>1.46 (97.8%)  |
| DeepSeek-33B-Ins         | 0.29<br>0.25 (13.8%) | 2.21<br>1.84 (16.7%) | 34.53<br>32.67 (5.4%)   | 1.06<br>0.99 (6.6%)  | 14.44<br>8.15 (43.6%)    | 2.91<br>1.55 (46.7%)   |
| CodeLlama-7B             | 4.70<br>4.52 (3.8%)  | 3.68<br>3.54 (3.8%)  | 46.76<br>38.67 (17.3%)  | 0.99<br>0.82 (17.2%) | 212.41<br>157.76 (25.7%) | 1.93<br>1.43 (25.9%)   |
| CodeLlama-13B            | 2.45<br>2.28 (6.9%)  | 2.19<br>2.04 (6.8%)  | 42.46<br>42.12 (0.8%)   | 0.93<br>0.93 (0.0%)  | 137.40<br>119.36 (13.1%) | 1.51<br>1.31 (13.2%)   |
| CodeLlama-34b            | 1.05<br>1.02 (2.9%)  | 7.75<br>7.34 (5.3%)  | 57.57<br>40.62 (29.4%)  | 1.70<br>1.11 (34.7%) | 94.79<br>52.12 (45.0%)   | 15.65<br>7.02 (55.1%)  |
| CodeLlama-70b            | 0.52<br>0.47 (9.6%)  | 3.93<br>3.84 (2.3%)  | 109.61<br>26.42 (75.9%) | 3.57<br>1.00 (72.0%) | 203.92<br>14.53 (92.9%)  | 54.15<br>6.52 (88.0%)  |
| XwinCoder-7B             | 2.80<br>2.43 (13.2%) | 2.81<br>2.44 (13.2%) | 55.54<br>49.10 (11.6%)  | 1.52<br>1.34 (11.8%) | 208.23<br>158.20 (24.0%) | 3.47<br>2.64 (23.9%)   |
| XwinCoder-34B            | 0.77<br>0.69 (10.4%) | 5.68<br>5.11 (10.0%) | 49.77<br>52.12 (-4.7%)  | 1.49<br>1.47 (1.3%)  | 61.36<br>57.89 (5.7%)    | 12.11<br>9.92 (18.1%)  |
| StarCoder2-3B            | 1.10<br>1.02 (7.3%)  | 1.25<br>1.15 (8.0%)  | 24.31<br>24.28 (0.1%)   | 1.00<br>1.00 (0.0%)  | 17.47<br>16.38 (6.2%)    | 1.19<br>1.12 (5.9%)    |
| StarCoder2-7B            | 3.69<br>2.99 (19.0%) | 5.34<br>4.32 (19.1%) | 26.42<br>26.40 (0.1%)   | 1.08<br>1.08 (0.0%)  | 82.38<br>68.61 (16.7%)   | 7.62<br>6.35 (16.7%)   |
| StarCoder2-15B           | 0.93<br>0.12 (87.1%) | 7.58<br>1.03 (86.4%) | 26.35<br>27.67 (-5.0%)  | 1.00<br>1.01 (-1.0%) | 22.02<br>2.03 (90.8%)    | 10.88<br>1.06 (90.3%)  |
| WizardCoder-13B          | 3.43<br>2.93 (14.6%) | 2.11<br>1.80 (14.7%) | 86.72<br>71.02 (18.1%)  | 1.35<br>1.11 (17.8%) | 324.83<br>219.69 (32.4%) | 1.92<br>1.30 (32.3%)   |
| Closed-source LLMs       |                      |                      |                         |                      |                          |                        |
| GPT-3.5-Turbo-0301       | 0.36<br>0.28 (22.2%) | 2.50<br>2.01 (19.6%) | 91.25<br>36.08 (60.5%)  | 2.45<br>0.99 (59.6%) | 157.50<br>12.43 (92.1%)  | 19.75<br>1.64 (91.7%)  |
| GPT-3.5-Turbo-1106       | 0.28<br>0.26 (7.1%)  | 1.96<br>1.90 (3.1%)  | 36.12<br>34.02 (5.8%)   | 1.01<br>1.00 (1.0%)  | 12.79<br>11.41 (10.8%)   | 1.73<br>1.62 (6.4%)    |
| GPT-4-Turbo-Preview      | 0.27<br>0.25 (7.4%)  | 1.96<br>1.88 (4.1%)  | 33.94<br>33.17 (2.3%)   | 1.00<br>1.00 (0.0%)  | 11.82<br>10.18 (13.9%)   | 1.89<br>1.76 (6.9%)    |
| GPT-4                    | 0.31<br>0.28 (9.7%)  | 2.19<br>2.06 (5.9%)  | 80.88<br>63.82 (21.1%)  | 2.26<br>1.83 (19.0%) | 129.91<br>80.74 (37.8%)  | 17.90<br>11.86 (33.7%) |
| Claude-3-Haiku           | 0.36<br>0.33 (8.3%)  | 2.51<br>2.30 (8.4%)  | 48.33<br>37.37 (22.7%)  | 1.30<br>1.03 (20.8%) | 52.67<br>17.18 (67.4%)   | 6.73<br>2.37 (64.8%)   |
| Claude-3-Sonnet          | 0.42<br>0.35 (16.7%) | 2.90<br>2.47 (14.8%) | 60.46<br>42.31 (30.0%)  | 1.62<br>1.17 (27.8%) | 82.52<br>28.95 (64.9%)   | 10.12<br>3.76 (62.8%)  |

source LLMs.

The improvements in efficiency metrics across both open-source and closed-source LLMs highlight the generalizability and adaptability of EFFILEARNER in enhancing efficiency. By iteratively refining the generated code based on efficiency profiler feedback, EFFILEARNER enables LLMs to produce more efficient code without compromising the correctness of the generated solutions. The consistent improvements across various models and architectures demonstrate the potential of EFFILEARNER as a model-agnostic approach for optimizing the efficiency of LLM-generated code in real-world applications.

### 3.4.4 Impact of Self-Optimization Steps

As illustrated in Figure 3.2, EFFILEARNER refines the code iteratively using the overhead profile obtained from previous steps. To investigate the impact of the number of self-

**Table 3.3:** Effect of the number of self-optimization steps in EFFILEARNER.

| Steps              | ET (s)       | NET          | MU (Mb)       | NMU          | TMU (Mb*s)    | NTMU         |
|--------------------|--------------|--------------|---------------|--------------|---------------|--------------|
| CodeLlama-70B      |              |              |               |              |               |              |
| 0                  | 0.52         | 3.93         | 109.61        | 3.57         | 203.92        | 54.15        |
| 1                  | 0.48 (7.7%)  | 3.94 (-0.3%) | 26.47 (75.9%) | 1.00 (72.0%) | 14.91 (92.7%) | 6.69 (87.6%) |
| 2                  | 0.48 (7.7%)  | 3.89 (1.0%)  | 26.47 (75.9%) | 1.00 (72.0%) | 14.69 (92.8%) | 6.60 (87.8%) |
| 3                  | 0.47 (9.6%)  | 3.85 (2.0%)  | 26.42 (75.9%) | 1.00 (72.0%) | 14.60 (92.8%) | 6.56 (87.9%) |
| 4                  | 0.47 (9.6%)  | 3.84 (2.3%)  | 26.42 (75.9%) | 1.00 (72.0%) | 14.54 (92.9%) | 6.53 (87.9%) |
| 5                  | 0.47 (9.6%)  | 3.84 (2.3%)  | 26.42 (75.9%) | 1.00 (72.0%) | 14.53 (92.9%) | 6.52 (88.0%) |
| GPT-3.5-Turbo-0301 |              |              |               |              |               |              |
| 0                  | 0.36         | 2.50         | 91.25         | 2.45         | 157.50        | 19.75        |
| 1                  | 0.33 (8.3%)  | 2.35 (6.0%)  | 36.09 (60.4%) | 0.99 (59.6%) | 13.70 (91.3%) | 1.81 (90.8%) |
| 2                  | 0.31 (13.9%) | 2.18 (12.8%) | 36.09 (60.4%) | 0.99 (59.6%) | 13.04 (91.7%) | 1.72 (91.3%) |
| 3                  | 0.29 (19.4%) | 2.06 (17.6%) | 36.08 (60.5%) | 0.99 (59.6%) | 12.57 (92.0%) | 1.66 (91.6%) |
| 4                  | 0.29 (19.4%) | 2.03 (18.8%) | 36.08 (60.5%) | 0.99 (59.6%) | 12.50 (92.1%) | 1.65 (91.6%) |
| 5                  | 0.28 (22.2%) | 2.01 (19.6%) | 36.08 (60.5%) | 0.99 (59.6%) | 12.43 (92.1%) | 1.64 (91.7%) |

optimization steps on the efficiency of the EFFILEARNER-optimized code, we conduct an ablation study by varying the number of steps from 0 to 5. Tab. 3.3 for CodeLlama-70B and GPT-3.5-Turbo-0301 at different self-optimization steps.

**CodeLlama-70B** As shown in Figure 3.2, the MU decreases from 109.61 (Mb) to 26.47 (Mb) after the first self-optimization step, reducing 75.9% maximum memory requirement compared with the initial code generated by CodeLlama-70B. Similarly, the TMU decreases from 54.15 (Mb\*s) to 6.69 (Mb\*s), reducing 87.6% of memory consumption during code execution. As the number of steps increases, the efficiency metrics gradually improve. By the fifth step, the ET reaches 0.47 (s), reducing the 1.9% execution time requirement compared with the first-step generated code, and the TMU settles at 14.53, reducing 0.2% total memory usage from the first step.

**GPT-3.5-Turbo-0301** Similar to CodeLlama-70B, the MU decreases from 91.25 (Mb) to 36.09 (Mb) after the first self-optimization step, reducing 60.4% maximum memory requirement compared with the initial code. The TMU also shows a substantial reduction from 157.50 (Mb\*s) to 13.70 (Mb\*s), reducing 91.3% memory consumption during code execution. As the number of steps increases, the efficiency metrics continue to improve steadily. By the fifth step, the ET reaches 0.28 (s), reducing the 15.2% execution time requirement compared with the first-step generated code, and the TMU settles at 12.43 (Mb\*s), reducing 9.3% total memory usage from the first step.

The evaluation results in Tab. 3.3 demonstrate the significant impact of the number of self-optimization steps on the efficiency of the EFFILEARNER-optimized code. For both CodeLlama-70B and GPT-3.5-Turbo-0301, the first self-optimization step yields the most substantial improvements in code efficiency. The MU and TMU metrics show significant reductions, indicating a decrease in maximum memory requirement and memory consumption during code execution. As the number of steps increases, the efficiency metrics continue to improve, albeit with diminishing returns. By the fifth step, the efficiency metrics reach their lowest values, demonstrating the effectiveness of EFFILEARNER’s iterative self-optimization approach in enhancing the efficiency of LLM-generated code. The evaluation results highlight that the majority of efficiency

**Table 3.4:** Contribution of different components in EFFILEARNER. We evaluate how different feedback profilers affect the efficiency of LLM-generated code. Unsupervised self-refine only requires LLMs to optimize the efficiency of the code. Result-Aware Self-Refine feedback the ET, MU, and TMU to the LLMs and require it to improve the efficiency. Memory Profiler and Execution Time Profiler feedback the memory profiler and execution time profiler to the LLMs and then LLMs can based on the profile optimize the efficiency of the code.

| Optimization Profile     | ET (s)        | NET           | MU (Mb)          | NMU             | TMU (Mb*s)        | NTMU              |
|--------------------------|---------------|---------------|------------------|-----------------|-------------------|-------------------|
| CodeLlama-70B            |               |               |                  |                 |                   |                   |
| Initial Version          | 0.52          | 3.93          | 109.61           | 3.57            | 203.92            | 54.15             |
| Unsupervised Self-Refine | 0.79 (-51.9%) | 6.87 (-74.8%) | 279.41 (-154.9%) | 10.58 (-196.4%) | 1261.83 (-518.8%) | 600.95 (-1009.8%) |
| Result-Aware Self-Refine | 0.79 (-51.9%) | 6.87 (-74.8%) | 282.57 (-157.8%) | 10.70 (-199.7%) | 1270.93 (-523.2%) | 605.29 (-1017.8%) |
| Memory Profiler          | 0.53 (-1.9%)  | 4.34 (-10.4%) | 26.38 (75.9%)    | 0.99 (72.3%)    | 15.77 (92.3%)     | 7.06 (87.0%)      |
| Execution Time Profiler  | 0.51 (1.9%)   | 4.17 (-6.1%)  | 26.44 (75.9%)    | 1.00 (72.0%)    | 15.53 (92.4%)     | 6.97 (87.1%)      |
| EFFILEARNER              | 0.47 (9.6%)   | 3.84 (2.3%)   | 26.42 (75.9%)    | 1.00 (72.0%)    | 14.53 (92.9%)     | 6.52 (88.0%)      |
| GPT-3.5-Turbo-0301       |               |               |                  |                 |                   |                   |
| Initial Version          | 0.36          | 2.50          | 91.25            | 2.45            | 157.50            | 19.75             |
| Unsupervised Self-Refine | 0.32 (11.1%)  | 2.46 (1.6%)   | 78.39 (14.1%)    | 2.12 (13.5%)    | 312.99 (-98.7%)   | 42.42 (-114.8%)   |
| Result-Aware Self-Refine | 0.30 (16.7%)  | 2.25 (10.0%)  | 58.65 (35.7%)    | 1.61 (34.3%)    | 195.49 (-24.1%)   | 27.16 (-37.5%)    |
| Memory Profiler          | 0.34 (5.6%)   | 2.40 (4.0%)   | 36.85 (59.6%)    | 1.00 (59.2%)    | 16.34 (89.6%)     | 2.10 (89.4%)      |
| Execution Time Profiler  | 0.33 (8.3%)   | 2.34 (6.4%)   | 36.43 (60.1%)    | 0.99 (59.6%)    | 14.07 (91.1%)     | 1.81 (90.8%)      |
| EFFILEARNER              | 0.28 (22.2%)  | 2.01 (19.6%)  | 36.08 (60.5%)    | 0.99 (59.6%)    | 12.43 (92.1%)     | 1.64 (91.7%)      |

improvements occur in the first few steps, with subsequent steps contributing to further refinements of the optimized code.

### 3.4.5 Feedback of Overhead Profile

As shown in Figure 3.2, the overhead profile is used to guide LLMs to refine their previously generated code, which then plays a crucial role in improving the code efficiency of LLM-generated code. To show the effectiveness of the overhead profile in guiding LLMs to refine their generated code, we compare the performance of EFFILEARNER with two alternative approaches: Unsupervised Self-Refine and Result-Aware Self-Refine [108, 150]. Unsupervised Self-Refine uses a prompt that directly requires the LLM to refine the code without providing additional information. Result-Aware Self-Refine feeds the ET, MU, and TMU, then requires the LLM to refine the code based on these metrics. Tab. 3.4 presents the code efficiency metrics for CodeLlama-70B and GPT-3.5-Turbo-0301 using different code refinement approaches.

**CodeLlama-70B** Unsupervised Self-Refine and Result-Aware Self-Refine result in significant increases in ET, memory usage (MU), and TMU compared to the initial version. Unsupervised Self-Refine increases ET by 51.9%, MU by 154.9%, and TMU by 518.8%, while Result-Aware Self-Refine increases ET by 51.9%, MU by 157.8%, and TMU by 523.2%. In contrast, EFFILEARNER incorporates the overhead profile feedback and achieves a 9.6% reduction in ET, a 75.9% reduction in MU, and a 92.9% reduction in TMU compared to the initial version.

**GPT-3.5-Turbo-0301** Unsupervised Self-Refine and Result-Aware Self-Refine show some improvements in ET and MU compared to the initial version. Unsupervised Self-Refine reduces ET by 11.1% and MU by 14.1%, while Result-Aware Self-Refine reduces ET by 16.7% and MU by 35.7%. However, both approaches lead to substantial increases in TMU, with Unsupervised Self-Refine increasing TMU by 98.7% and Result-Aware Self-

**Table 3.5:** Evaluation results of SOAP and baselines. Since the finetuned link for GPT-3.5-turbo from PIE is not available, we use the fine-tuned CodeLlama 7B for experiments. Due to the fine-tuned PIE CodeLlama 7B does not have the same correct tasks as the original CodeLlama, we then do not provide the initial version for the experiments.

| OptimizationProfile           | ET(s)        | NET          | MU(Mb)        | NMU          | TMU(Mb*s)     | NTMU         |
|-------------------------------|--------------|--------------|---------------|--------------|---------------|--------------|
| GPT-3.5-Turbo-0301            |              |              |               |              |               |              |
| InitialVersion                | 0.36         | 2.50         | 91.25         | 2.45         | 157.50        | 19.75        |
| UnsupervisedSelf-Refine       | 0.32         | 2.46         | 78.39         | 2.12         | 312.99        | 42.42        |
| Result-AwareSelf-Refine       | 0.30         | 2.25         | 58.65         | 1.61         | 195.49        | 27.16        |
| Self-Edit                     | 0.42         | 3.67         | 59.86         | 1.65         | 24.87         | 3.28         |
| DirectlyEfficiency            | 0.43         | 3.03         | 59.11         | 1.67         | 20.37         | 2.88         |
| Self-RefineEfficiency         | 0.40         | 2.83         | 59.11         | 1.67         | 18.80         | 2.65         |
| IsSelf-Refine                 | 0.40         | 2.88         | 61.83         | 1.81         | 36.29         | 5.69         |
| Self-Reasoning                | 0.89         | 6.21         | 60.64         | 1.62         | 45.91         | 5.61         |
| Self-Reflection               | 0.81         | 5.67         | 60.64         | 1.62         | 39.35         | 4.80         |
| EFFILEARNER                   | 0.28 (22.2%) | 2.01 (19.6%) | 36.08 (60.5%) | 0.99 (59.6%) | 12.43 (92.1%) | 1.64 (91.7%) |
| CodeLlama7B(PIE:HQ+SelfPlay)  |              |              |               |              |               |              |
| PIE+Zero-Shot                 | 0.87         | 5.73         | 74.83         | 1.81         | 109.29        | 9.69         |
| PIE+SOAP+Zero-Shot            | 0.79         | 5.41         | 65.78         | 1.68         | 89.90         | 7.84         |
| PIE+Few-Shot                  | 0.82         | 5.58         | 73.57         | 1.74         | 98.02         | 8.92         |
| PIE+SOAP+Few-Shot             | 0.41         | 2.97         | 73.10         | 1.74         | 59.69         | 5.09         |
| PIE+CoT                       | 0.79         | 5.14         | 73.14         | 1.74         | 63.93         | 5.35         |
| PIE+SOAP+COT                  | 0.45         | 2.84         | 71.15         | 1.71         | 58.06         | 4.77         |
| PIE+DynamicRetrieval,k=4      | 0.74         | 5.36         | 68.64         | 1.51         | 85.24         | 7.78         |
| PIE+SOAP+DynamicRetrieval,k=4 | 0.41         | 3.36         | 68.63         | 1.51         | 52.34         | 4.52         |
| Supersonic                    |              |              |               |              |               |              |
| Supersonic                    | 1.40         | 10.33        | 113.06        | 3.18         | 329.59        | 56.24        |
| Supersonic+SOAP               | 1.34         | 9.91         | 102.26        | 2.87         | 267.47        | 45.64        |

Refine increasing TMU by 24.1%. On the other hand, EFFILEARNER achieves a 22.2% reduction in ET, a 60.5% reduction in MU, and a 92.1% reduction in TMU compared to the initial version.

These results highlight the importance of the overhead profile feedback in guiding LLMs to generate more efficient code. Without the overhead profile, the code refinement process using alternative prompts fails to improve code efficiency and even leads to significant performance degradation. The overhead profile provides valuable insights into the resource consumption of the generated code, enabling LLMs to make targeted optimizations and achieve substantial efficiency improvements.

### 3.4.6 Comparison with baselines

As some of the existing works that are trained on efficient solutions can generate efficient solutions compared to the LLMs without efficiency-instruct. We then compare the efficiency of LLM-generated code for EFFILEARNER and below baselines: 1) Self-Edit [190], where we require LLM edit their previous generated code to generate more efficient solutions; 2) ask the model to generate an efficient version of code (Directly Efficiency); 3) ask the model to generate code and then directly optimize it (Self-Refine Efficiency); 4) self-refine its generated code (without efficiency requirement in the first generation); 5) Self-reasoning, where we require LLM itself reasoning to generate a efficient solution [178]; 6) Self-Reflection [149], PIE [153], and 7) Supersonic [196].

The evaluation results are shown in Tab. 3.5, which demonstrate the superiority of SOAP compared to various baselines. SOAP significantly outperforms methods such as Self-Edit, Directly Efficiency, Self-Refine Efficiency, Self-Reasoning, and Self-Reflection in terms of both execution time and memory usage optimization. For instance, SOAP

**Table 3.6:** Evaluation results of EFFILEARNER’s effectiveness in the HumanEval dataset.

| Steps                        | ET (s)       | NET          | MU (Mb)       | NMU         | TMU (Mb*s)   | NTMU         |
|------------------------------|--------------|--------------|---------------|-------------|--------------|--------------|
| OpenCodeInterpreter-DS-1.3B  | 0.20         | 0.86         | 57.24         | 1.00        | 6.63         | 0.84         |
|                              | 0.19 (5.0%)  | 0.81 (5.8%)  | 57.17 (0.1%)  | 1.00 (0.0%) | 6.20 (6.5%)  | 0.79 (6.0%)  |
| OpenCodeInterpreter-DS-6.7B  | 0.21         | 0.98         | 58.83         | 1.06        | 6.79         | 0.99         |
|                              | 0.21 (0.0%)  | 0.97 (1.0%)  | 58.79 (0.1%)  | 1.06 (0.0%) | 6.64 (2.2%)  | 0.97 (2.0%)  |
| OpenCodeInterpreter-DS-33B   | 0.21         | 0.95         | 59.90         | 1.05        | 7.05         | 0.94         |
|                              | 0.21 (0.0%)  | 0.93 (2.1%)  | 59.93 (-0.1%) | 1.05 (0.0%) | 6.87 (2.6%)  | 0.92 (2.1%)  |
| deepseek-coder-1.3b-instruct | 0.23         | 0.90         | 62.80         | 1.00        | 7.85         | 0.87         |
|                              | 0.22 (4.3%)  | 0.85 (5.6%)  | 62.96 (-0.3%) | 1.00 (0.0%) | 7.80 (0.6%)  | 0.86 (1.1%)  |
| deepseek-coder-6.7b-instruct | 0.22         | 0.76         | 59.57         | 1.00        | 7.34         | 0.77         |
|                              | 0.19 (13.6%) | 0.68 (10.5%) | 59.72 (-0.3%) | 1.00 (0.0%) | 6.59 (10.2%) | 0.69 (10.4%) |
| deepseek-coder-33b-instruct  | 0.21         | 0.95         | 63.52         | 0.99        | 7.18         | 0.95         |
|                              | 0.20 (4.8%)  | 0.92 (3.2%)  | 63.49 (0.0%)  | 0.99 (0.0%) | 6.99 (2.6%)  | 0.92 (3.2%)  |
| CodeLlama-7b-Instruct-hf     | 0.20         | 0.71         | 57.39         | 0.91        | 7.08         | 0.70         |
|                              | 0.18 (10.0%) | 0.63 (11.3%) | 57.07 (0.6%)  | 0.91 (0.0%) | 6.18 (12.7%) | 0.61 (12.9%) |
| CodeLlama-13b-Instruct-hf    | 0.23         | 0.95         | 58.13         | 0.96        | 7.97         | 0.94         |
|                              | 0.20 (13.0%) | 0.80 (15.8%) | 58.03 (0.2%)  | 0.96 (0.0%) | 6.64 (16.7%) | 0.79 (16.0%) |
| CodeLlama-34b-Instruct-hf    | 0.24         | 0.95         | 61.79         | 1.01        | 8.45         | 0.96         |
|                              | 0.21 (12.5%) | 0.81 (14.7%) | 61.55 (0.4%)  | 1.00 (1.0%) | 6.99 (17.3%) | 0.80 (16.7%) |
| CodeLlama-70b-Instruct-hf    | 0.21         | 0.93         | 60.19         | 1.01        | 6.76         | 1.01         |
|                              | 0.18 (14.3%) | 0.79 (15.1%) | 59.49 (1.2%)  | 1.00 (1.0%) | 5.75 (14.9%) | 0.86 (14.9%) |
| XwinCoder-13B                | 0.27         | 1.08         | 61.14         | 1.04        | 9.25         | 1.09         |
|                              | 0.25 (7.4%)  | 1.01 (6.5%)  | 61.15 (-0.0%) | 1.04 (0.0%) | 8.62 (6.8%)  | 1.02 (6.4%)  |
| XwinCoder-34B                | 0.25         | 1.07         | 60.75         | 1.05        | 8.46         | 1.08         |
|                              | 0.22 (12.0%) | 0.93 (13.1%) | 60.75 (0.0%)  | 1.05 (0.0%) | 7.33 (13.4%) | 0.94 (13.0%) |
| WizardCoder-7B               | 0.21         | 0.91         | 58.59         | 1.01        | 6.63         | 0.89         |
|                              | 0.18 (14.3%) | 0.79 (13.2%) | 57.97 (1.1%)  | 1.00 (1.0%) | 5.79 (12.7%) | 0.78 (12.4%) |
| WizardCoder-13B              | 0.21         | 0.81         | 60.59         | 1.00        | 7.22         | 0.79         |
|                              | 0.19 (9.5%)  | 0.73 (9.9%)  | 60.53 (0.1%)  | 1.00 (0.0%) | 6.47 (10.4%) | 0.71 (10.1%) |
| WizardCoder-34B              | 0.22         | 0.79         | 58.13         | 1.00        | 7.10         | 0.78         |
|                              | 0.17 (22.7%) | 0.62 (21.5%) | 58.42 (-0.5%) | 1.00 (0.0%) | 5.46 (23.1%) | 0.60 (23.1%) |
| starcoder2-3b                | 0.24         | 1.02         | 62.45         | 1.00        | 7.73         | 0.89         |
|                              | 0.19 (20.8%) | 0.79 (22.5%) | 62.69 (-0.4%) | 1.00 (0.0%) | 6.68 (13.6%) | 0.77 (13.5%) |
| starcoder2-7b                | 0.21         | 0.89         | 62.53         | 1.00        | 7.41         | 0.85         |
|                              | 0.19 (9.5%)  | 0.78 (12.4%) | 62.85 (-0.5%) | 1.00 (0.0%) | 6.40 (13.6%) | 0.74 (12.9%) |

reduces the average execution time by 22.2% and decreases NTMU by 91.7% compared to the initial version. Furthermore, when applied to the code generated by state-of-the-art models like PIE and Supersonic, SOAP further enhances their efficiency. These results highlight the effectiveness of SOAP’s dual optimization focus and overhead profile-guided optimization approach in improving the efficiency of LLM-generated code. The inclusion of both memory profiler and execution time profiler in SOAP proves to be a key factor in achieving these state-of-the-art results, setting it apart from existing methods that primarily focus on execution time optimization or rely on self-reasoning and self-reflection techniques.

### 3.4.7 Discussion

**Generalizability across benchmarks** In Tab. 4.2, we evaluated EFFILEARNER’s effectiveness on the EffiBench dataset. To illustrate EFFILEARNER’s generalizability in other datasets, we conduct experiments on the HumanEval and MBPP datasets in Tab. 3.6 and Tab. 3.7, where the coding efficiency of CodeLlama and other LLMs (See Tab. 3.6) also increases when we utilize EFFILEARNER to optimize LLM-generated code. For example, the ET of CodeLlama-70B decreases from 0.21 (s) to 0.18 (s), which reduces 14.3% execution time. As shown in Tab. 3.6 and Tab. 3.7, the results demonstrate that EFFILEARNER can consistently improve the efficiency of LLM-generated code for other datasets.

**Generalizability across LLMs** In Tab. 4.2, we evaluate EFFILEARNER’s effectiveness on six types of open-source LLMs. To illustrate EFFILEARNER’s generalizability in other LLMs, we also conduct experiments on other LLMs in Tab. 3.8. Our evaluation results demonstrate that EFFILEARNER can improve the efficiency of LLM-generated code for different LLMs. For example, the execution time of Mistral-7B-codealpaca-lora decreases

**Table 3.7:** Evaluation results of EFFILEARNER’s effectiveness in the MBPP dataset

| Steps                            | ET (s)          | NET             | MU (Mb)          | NMU            | TMU (Mb*s)           | NTMU            |
|----------------------------------|-----------------|-----------------|------------------|----------------|----------------------|-----------------|
| OpenCodeInterpreter-DS-1.3B      | 0.28<br>(10.7%) | 0.84<br>(10.6%) | 59.01<br>(0.0%)  | 1.01<br>(0.0%) | 11.73<br>(9.7%)      | 0.98<br>(9.2%)  |
| OpenCodeInterpreter-DS-6.7B      | 0.26<br>(19.2%) | 1.06<br>(17.9%) | 58.99<br>(0.0%)  | 1.00<br>(0.0%) | 9.25<br>(7.1%)       | 1.08<br>(23.1%) |
| OpenCodeInterpreter-DS-33B       | 0.44<br>(29.5%) | 1.59<br>(28.3%) | 58.72<br>(0.0%)  | 1.00<br>(0.0%) | 20.19<br>(13.2%)     | 1.86<br>(34.4%) |
| deepseek-coder-1.3b-instruct     | 0.63<br>(1.6%)  | 1.68<br>(2.4%)  | 354.01<br>(4.0%) | 6.05<br>(4.0%) | 1463.46<br>(1414.13) | 89.12<br>(3.4%) |
| deepseek-coder-6.7b-instruct     | 0.76<br>(72.4%) | 3.62<br>(72.9%) | 58.44<br>(0.2%)  | 1.00<br>(0.0%) | 39.11<br>(6.67)      | 5.69<br>(83.0%) |
| deepseek-coder-33b-instruct      | 0.58<br>(67.2%) | 2.33<br>(67.8%) | 53.48<br>(0.3%)  | 0.91<br>(0.0%) | 28.74<br>(5.88)      | 3.16<br>(79.4%) |
| CodeLlama-7b-Instruct-hf         | 0.45<br>(6.7%)  | 2.04<br>(7.4%)  | 56.96<br>(0.3%)  | 0.97<br>(0.0%) | 13.26<br>(11.98)     | 1.79<br>(9.5%)  |
| CodeLlama-13b-Instruct-hf        | 0.53<br>(1.9%)  | 2.11<br>(3.3%)  | 55.37<br>(0.1%)  | 0.95<br>(0.0%) | 21.75<br>(21.13)     | 2.34<br>(2.6%)  |
| CodeLlama-34b-Instruct-hf        | 0.42<br>(2.4%)  | 1.18<br>(4.2%)  | 69.80<br>(0.7%)  | 1.19<br>(0.0%) | 84.01<br>(74.78)     | 5.47<br>(4.87)  |
| CodeLlama-70b-Instruct-hf        | 0.23<br>(13.0%) | 1.06<br>(12.3%) | 58.13<br>(0.1%)  | 0.98<br>(0.0%) | 7.65<br>(6.67)       | 1.05<br>(13.3%) |
| XwinCoder-7B                     | 0.23<br>(21.7%) | 1.14<br>(21.1%) | 58.45<br>(0.0%)  | 1.00<br>(0.0%) | 7.19<br>(5.89)       | 1.10<br>(18.2%) |
| XwinCoder-13B                    | 0.50<br>(18.0%) | 1.96<br>(17.9%) | 58.38<br>(0.1%)  | 1.00<br>(0.0%) | 23.88<br>(18.95)     | 2.50<br>(20.8%) |
| XwinCoder-34B                    | 0.38<br>(7.9%)  | 1.44<br>(8.3%)  | 58.27<br>(0.1%)  | 1.00<br>(0.0%) | 14.77<br>(13.54)     | 1.48<br>(8.1%)  |
| WizardCoder-Python-7B-V1.0-GPTQ  | 0.22<br>(9.1%)  | 1.05<br>(11.4%) | 58.44<br>(0.2%)  | 0.99<br>(0.0%) | 7.19<br>(6.41)       | 1.03<br>(11.7%) |
| WizardCoder-Python-13B-V1.0-GPTQ | 0.62<br>(4.8%)  | 1.35<br>(5.2%)  | 57.74<br>(0.1%)  | 0.99<br>(0.0%) | 30.66<br>(29.56)     | 1.43<br>(3.5%)  |
| WizardCoder-Python-34B-V1.0-GPTQ | 0.68<br>(4.4%)  | 2.43<br>(4.1%)  | 56.75<br>(-0.1%) | 0.97<br>(0.0%) | 34.06<br>(32.63)     | 3.14<br>(4.1%)  |
| starcoder2-3b                    | 0.17<br>(5.9%)  | 0.83<br>(3.6%)  | 45.82<br>(5.2%)  | 0.79<br>(6.3%) | 5.10<br>(4.69)       | 0.77<br>(9.1%)  |
| starcoder2-7b                    | 1.72<br>(0.0%)  | 8.63<br>(0.2%)  | 25.61<br>(0.2%)  | 0.44<br>(0.0%) | 40.42<br>(40.19)     | 6.22<br>(0.5%)  |
| starcoder2-15b                   | 0.19<br>(5.3%)  | 1.05<br>(5.7%)  | 58.62<br>(0.8%)  | 1.01<br>(1.0%) | 6.23<br>(5.92)       | 1.05<br>(4.8%)  |

**Table 3.8:** Code efficiency of widely-studied LLMs reported by EFFILEARNER.

| Model                     | ET (s)          | NET              | MU (Mb)           | NMU             | TMU (Mb*s)         | NTMU             |
|---------------------------|-----------------|------------------|-------------------|-----------------|--------------------|------------------|
| Phind-CodeLlama-34B-v2    | 0.52<br>(23.1%) | 3.28<br>(23.5%)  | 157.16<br>(56.6%) | 3.36<br>(56.8%) | 337.30<br>(80.5%)  | 24.44<br>(80.1%) |
| Artigenz-Coder-DS-6.7B    | 0.39<br>(17.9%) | 2.75<br>(16.4%)  | 65.73<br>(10.2%)  | 1.70<br>(4.7%)  | 95.65<br>(79.67)   | 10.87<br>(16.7%) |
| MagiCoder-DS-6.7B         | 0.22<br>(4.5%)  | 1.59<br>(5.7%)   | 40.19<br>(4.7%)   | 1.09<br>(1.8%)  | 17.58<br>(15.27)   | 2.28<br>(13.1%)  |
| Mistral-7B-codealpa-lora  | 2.36<br>(38.6%) | 18.40<br>(33.9%) | 28.88<br>(5.0%)   | 1.00<br>(-3.0%) | 57.92<br>(35.46)   | 24.36<br>(38.8%) |
| CodeFuse-DeepSeek-33B     | 0.40<br>(2.5%)  | 3.10<br>(2.9%)   | 70.39<br>(10.2%)  | 2.06<br>(10.2%) | 191.15<br>(156.81) | 32.20<br>(18.0%) |
| CodeLlama-34b-hf          | 2.08<br>(1.95)  | 15.68<br>(6.4%)  | 46.41<br>(0.0%)   | 1.26<br>(0.0%)  | 128.46<br>(125.22) | 17.87<br>(2.5%)  |
| speechless-starcoder2-15b | 0.19<br>(0.13)  | 1.74<br>(31.6%)  | 27.39<br>(0.5%)   | 0.99<br>(0.0%)  | 3.20<br>(2.17)     | 1.75<br>(32.0%)  |
| gpt-3.5-turbo-0613        | 0.56<br>(12.5%) | 4.32<br>(13.2%)  | 35.48<br>(0.0%)   | 1.00<br>(0.0%)  | 20.11<br>(17.84)   | 3.00<br>(11.3%)  |

from 2.36 (s) to 1.45 (s), which reduces 38.6% execution time compared with the initial code. The total memory usage of Phind-CodeLlama-34B-v2 also decreases from 337.30 (Mb\*s) to 65.64 (Mb\*s), which reduces 80.5% total memory requirement.

**Impact on correctness** We provide the pass@1 of LLM-generated initial code and EFFILEARNER optimized code for EffiBench in Tab. 3.9. We observe that the pass@1 of EFFILEARNER optimized code may be lower than LLM-generated initial code. The key reason is that during the self-optimization process, EFFILEARNER only uses public test cases to guide code efficiency optimization for correct initial code. However, since public test cases may not cover all edge cases in the private test cases (test cases used to evaluate pass@1 of LLMs), this can cause the pass@1 of EFFILEARNER generated code to be lower than the initial code. Nevertheless, we observe that the pass@1 of EFFILEARNER only decreases by about 0% to 0.5%, which means that only a few of the codes will be incorrect. As shown in Tab. 4.2, the code efficiency is largely increased. We believe that



**Table 3.9:** Pass@1 of LLMs generated initial code and EFFILEARNER optimized code.

| Model                       | Initial Pass@1 | EFFILEARNER Pass@1 |
|-----------------------------|----------------|--------------------|
| OpenCodeInterpreter-DS-1.3B | 5.8            | 5.4                |
| OpenCodeInterpreter-DS-6.7B | 13.6           | 13.2               |
| OpenCodeInterpreter-DS-33B  | 24.7           | 24.4               |
| deepseek-1.3b-Ins           | 4.8            | 4.5                |
| deepseek-6.7b-Ins           | 7.2            | 7.0                |
| deepseek-33b-Ins            | 10.0           | 9.9                |
| CodeLlama-7b                | 7.0            | 7.0                |
| CodeLlama-13b               | 9.7            | 9.6                |
| CodeLlama-34b               | 13.5           | 13.0               |
| CodeLlama-70b               | 7.8            | 7.4                |
| XwinCoder-13B               | 10.5           | 10.2               |
| XwinCoder-34B               | 21.2           | 21.2               |
| starcoder2-3b               | 1.6            | 1.2                |
| starcoder2-7b               | 1.9            | 1.8                |
| starcoder2-15b              | 0.7            | 0.4                |
| WizardCoder-13B             | 4.0            | 3.9                |

**Table 3.10:** Overhead of different code efficiency optimization methods for GPT-3.5-turbo.

| Method                  | ET (s) | Total Token (m) | per Iter (K) | ET(s) | NET  | MU(Mb) | NMU  | TMU(Mb*s) | NTMU  |
|-------------------------|--------|-----------------|--------------|-------|------|--------|------|-----------|-------|
| InitialVersion          | 76.81  | 1.1             | 1.1          | 0.36  | 2.50 | 91.25  | 2.45 | 157.50    | 19.75 |
| UnsupervisedSelf-Refine | 416.81 | 3.9             | 1.16         | 0.32  | 2.46 | 78.39  | 2.12 | 312.99    | 42.42 |
| Result-AwareSelf-Refine | 419.24 | 3.9             | 1.16         | 0.30  | 2.25 | 58.65  | 1.61 | 195.49    | 27.16 |
| LineProfiler            | 555.68 | 4.7             | 1.49         | 0.33  | 2.34 | 36.43  | 0.99 | 14.07     | 1.81  |
| MemoryProfiler          | 536.26 | 4.7             | 1.49         | 0.34  | 2.40 | 36.85  | 1.00 | 16.34     | 2.10  |
| EFFILEARNER             | 566.67 | 5.4             | 1.78         | 0.28  | 2.01 | 36.08  | 0.99 | 12.43     | 1.64  |

this minor decrease in pass@1 is worthwhile considering the significant efficiency gains.

**Overhead of EFFILEARNER** To address the reviewer’s concern about overhead and context window limitations, we provide detailed metrics on execution time, token usage, per iteration input/output token usage, and efficiency in Tab. 3.10. Our results demonstrate that EFFILEARNER requires approximately 8x more execution time compared to the Initial Version. However, this overhead is justified as the optimized code resulting from EFFILEARNER significantly reduces the execution time and memory usage for real-world software that could be executed millions of times. In addition, while the optimization process itself is resource-intensive, it yields substantial efficiency gains in deployment scenarios. For example, the average memory peak (MU) of EFFILEARNER-generated code only requires 40% compared with the initially generated code, which can help source code applied in memory-constrained environments, such as embedded systems or mobile devices. Furthermore, the reduced memory footprint and improved execution speed of the optimized code can lead to better overall system performance, especially in scenarios where the software is frequently used or runs on resource-limited hardware. As a result, the upfront computational cost of the optimization process is offset by the long-term benefits of more efficient and lightweight code in real-world applications.

**Table 3.11:** Evaluation results of EFFILEARNER on the HumanEval-ET (C++) dataset. ET: Execution Time; NET: Normalized Execution Time; MU: Memory Usage; NMU: Normalized Memory Usage; TMU: Total Memory Usage; NTMU: Normalized Total Memory Usage.

| HumanEval (C++)           | ET (ms) | NET | MU (KB) | NMU | TMU (KB*ms) | NTMU |
|---------------------------|---------|-----|---------|-----|-------------|------|
| CodeLlama-70b-Instruct-hf |         |     |         |     |             |      |
| Initial Version           | 667.9   | 1.5 | 93.2    | 1.3 | 58.9        | 2.1  |
| EFFILEARNER               | 459.5   | 1.0 | 72.2    | 1.0 | 34.1        | 1.3  |
| gpt-3.5-turbo-0301        |         |     |         |     |             |      |
| Initial Version           | 668.1   | 1.5 | 78.9    | 1.1 | 79.0        | 2.6  |
| EFFILEARNER               | 577.3   | 1.2 | 71.5    | 1.0 | 63.8        | 2.1  |

Next, EFFILEARNER requires an average of 1.78K tokens for each input+output task iteration. Given that existing LLMs such as GPT-3.5 and GPT-4 have context windows of 4K tokens or more, they can effectively handle the global profiler information provided by EFFILEARNER. This capability allows the LLMs to understand comprehensive profiling data and perform global code optimizations, addressing concerns about scalability within the context window limitations. By leveraging the available context window efficiently and optimizing the code based on detailed profiling data, EFFILEARNER manages to enhance the performance and scalability of the generated code, making it a valuable tool despite its initial overhead.

**Generalizability on C++** The results in Tab. 4.2 focus on Python language tasks, which raises a concern about whether EFFILEARNER can be used in other languages. To address this concern, we conduct experiments for EFFILEARNER on the HumanEval-ET (C++) dataset. The evaluation results, shown in Tab. 3.11, demonstrate that EFFILEARNER can also improve the efficiency of LLM-generated C++ programs. For instance, the average execution time of CodeLlama-70B-Instruct-hf decreases from 667.9ms to 459.5ms after applying EFFILEARNER.

**Case study** To illustrate how EFFILEARNER improves the efficiency of LLM-generated code, we provide a case illustration in Appendix ??-?. As shown in ??, the execution time of the initial code is 23.59 (s) while in the self-optimized code, the execution time decreases from 23.59 (s) to 3.36 (s). The key reason is that in the initial code, the algorithm uses a standard unidirectional Breadth-First Search (BFS), which explores all possible states level by level starting from the initial state. This method results in a large number of states to explore, leading to significant computational overhead. In contrast, the self-optimized code employs a bidirectional BFS, which simultaneously searches from both the initial state and the target state. This reduces the search space by meeting in the middle, significantly decreasing the number of states that need to be explored and thereby improving the execution time.

**Error Analysis** We also provide a case illustration to explain why some code efficiency does not improve significantly when EFFILEARNER is applied to LLM-generated code. As shown in Appendix ??-?, we observe that the initial code only requires 0.0012 (s) to execute, while in the optimized code, the execution time is still 0.0011 (s). The key reason for this minimal improvement is that both implementations already operate with the



same theoretical time complexity of  $O(\log(\min(m, n)))$ . Given the problem's constraints and small input sizes, the actual runtime differences are overshadowed by the inherent efficiency of the binary search algorithm. Additionally, the overhead of function calls and Python runtime operations can further minimize the observed performance gains. Therefore, while the optimized code may offer clearer partition management and slight improvements, the overall efficiency remains largely unchanged due to the already optimized nature of the initial approach.

### 3.5 Conclusion

This paper focuses on the critical issue of efficiency in code generated by LLMs. While LLMs have shown impressive capabilities in code generation, their output often suffers from suboptimal efficiency, leading to slower execution and higher resource consumption. To tackle this challenge, we first introduce EFFIBENCH in this chapter, a benchmark designed to evaluate the efficiency of code generated by various code generation models. EFFIBENCH encompasses 1,000 problems and consists of 11 distinct algorithmic subsets. Unlike previous benchmarks that primarily emphasize the correctness of code generation, EFFIBENCH extends the evaluation criteria to include both execution time analysis and memory usage analysis. Then, we propose EFFILEARNER, a novel self-optimization framework that leverages execution overhead profiles to guide LLMs in improving code efficiency. Extensive experiments and analysis demonstrate that EFFILEARNER significantly enhances the efficiency of LLM-generated code, achieving substantial reductions in execution time and memory usage. For future work, we would like to investigate the application of EFFILEARNER to other programming tasks and languages, as well as explore the potential benefits of incorporating domain-specific knowledge into the optimization process.

## Chapter 4

# EFFICODER: Unleashing Code Efficiency in Language Models

### 4.1 Introduction

Large language models (LLMs) have recently made significant strides across various tasks [130, 10, 11, 114], including code-related applications like code completion [30], debugging [67, 32], and translation[144, 4]. These advanced tools have been seamlessly integrated into popular development environments, enhancing developer productivity by providing intelligent code recommendations based on natural language instructions.

Before deploying LLMs into integrated development environments (IDEs) as tools, it is crucial to ensure that the generated code meets the required efficacy standards. To address this, researchers have explored various datasets to fine-tune LLMs, thereby improving the efficacy of LLM-generated code [132, 177]. For example, Code Alpaca [29] utilized the Self-Instruct framework [174] to synthesize data, while WizardCoder [106] employed the Evol-Instruct technique [183] to generate heuristic prompts for diverse solutions. Additionally, OSS-Instruct [181] created new coding problems using open-source snippets with LLMs, and Octopack [119] focused on curating high-quality Git commit messages that resemble natural language instructions. These fine-tuning efforts have led to increased correctness in LLM-generated code.

However, our observation is that existing works primarily focus on enhancing the correctness of LLM-generated code while neglecting to optimize its efficiency. As a result, the efficiency of such code often falls short compared to canonical solutions written by human developers. Recent studies [147, 125, 50, 78] also point out that LLM-generated code typically exhibits lower efficiency in terms of execution time and memory usage. For instance, on the EffiBench benchmark [79], even the most advanced LLMs, such as GPT-4-Turbo, produced less efficient code, with average and worst-case execution times being 1.69 and 45.49 times longer than those of canonical solutions, respectively.

Efficiency is crucial because inefficient code consumes more computational resources, leading to higher energy consumption and increased operational costs. This

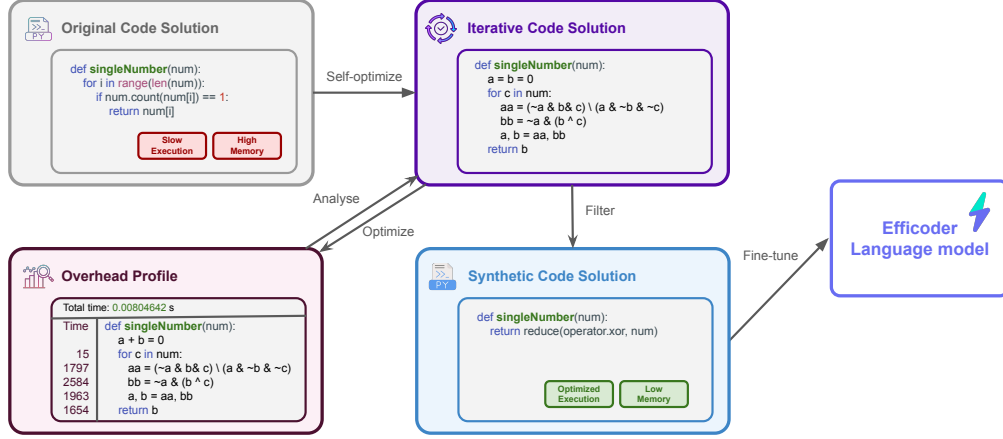
is particularly important in the context of sustainability, as the demand for computing power continues to grow, and reducing the environmental impact of large-scale computations becomes a pressing concern. Furthermore, inefficient code may be impractical for use in resource-constrained environments, such as mobile devices or embedded systems, where both energy and processing power are limited. This underscores the urgent need to develop new methods that can enhance both the **correctness** and **efficiency** of LLM-generated code.

In this paper, we introduce the dataset EFFICODER, aimed at fine-tuning LLMs to improve both code efficiency and correctness. We begin by aggregating source code from eight existing open-source datasets available on the Hugging Face platform. This is followed by a rigorous preprocessing and cleaning process, coupled with the generation of test cases for each task to evaluate code efficiency. The cleaned code is executed using these test cases to profile memory usage and execution time. Through a self-optimization process based on these profiles, we iteratively refine the code over five optimization cycles. The resulting optimized code, along with its associated metadata, forms the foundation of our final fine-tuning dataset, EFFICODER, which serves as a high-quality resource for training LLMs to generate more efficient code while ensuring correctness.

Extensive experiments on HumanEval [30] and EffiBench [79] demonstrate that fine-tuning LLMs with EFFICODER improves both correctness and efficiency. For example, the fine-tuned DeepSeek-Coder-6.7B [46] increases the pass@1 from 43.3% while also reducing the average execution time from 0.59 seconds to 0.41 seconds — representing a 30.5% improvement. Compared to PIE [153], which increases the pass@1 from 12.2% to 19.5% on HumanEval, the pass@1 of CodeLlama-7B [143] fine-tuned with EFFICODER further increases to 37.8%. In addition, EFFICODER decreases the execution time by 7.1% while PIE decreases it by 4.8%. We will fully open-source EFFICODER, the source code, and model weights to facilitate research.

To conclude, this paper makes the following contributions:

- We provide a framework to inspire researchers to construct code generation datasets containing efficient solutions for each code generation task, which is versatile and can be adapted to different programming languages and leverage various existing data sources. Unlike some other code generation datasets that rely on powerful models (e.g., GPT-4), our framework can be implemented only using open-sourced LLMs. The framework provides a systematic method for researchers to enhance existing datasets or create new ones focused on code efficiency across different languages and domains.
- Based on our proposed framework, we release the Effi-Code dataset. To the best of our knowledge, it is the first instruct tuning dataset that focuses on improving the efficiency of LLM-generated code. The primary purpose of Effi-Code is to instruct and fine-tune LLMs to ensure that the LLM-generated code is more efficient.
- We use Effi-Code to fine-tune widely used LLMs and will release these models



**Figure 4.1:** Overview of the construction pipeline for EFFICODER: We begin by filtering illegal tasks and collect the initial EFFICODER from different open-source datasets. Starting with the original code, we apply self-optimization to enhance efficiency, using test cases to profile execution overhead, and self-improve the code based on the profile. Finally, tasks that fail to have efficiency improvements are removed. We then have our final fine-tuning dataset, EFFICODER, which consists of optimized code and rich metadata, designed to train models for generating both efficient and correct code.

on the Hugging Face website in the final version. Different from existing datasets that are used to finetune the LLMs to improve the pass@1 of LLM-generated code, our evaluation results demonstrate that both the pass@1 and the efficiency results would be improved for LLMs finetuned on our Effi-Code dataset.

## 4.2 EFFICODER: Fine-Tuning For Efficiency

In this section, we provide a detailed pipeline for constructing the dataset EFFICODER for fine-tuning. Specifically, we first collect source code from seven existing open-source datasets available on the HuggingFace platform<sup>1</sup>. To ensure the quality and usability of the collected data, we use several filtering strategies, such as filtering tasks that are not algorithmic tasks and do not require efficiency optimization<sup>2</sup>. In addition, we also generate test cases for each task to ensure that we can measure the efficiency of each task's source code.

Next, we execute the cleaned source code using the generated test cases to profile the memory usage and execution time for each task. Then, we use Self-Optimization based on these overhAd Profiles (SOAP; [77]), which iteratively refines the code over five optimization cycles to generate an efficient solution for each task in the collected tasks. Finally, we process the optimized code and the associated metadata to create our final fine-tuning dataset, EFFICODER, which is carefully curated to provide a high-quality resource for training models to generate more efficient code while maintaining

<sup>1</sup><https://huggingface.co/docs/datasets/index>

<sup>2</sup>Data decontamination was not included in the filtering process as most of the tasks we collected have been decontaminated, such as OSS-Instruct [162].

correctness.

### 4.2.1 Source Data Collection

To construct a high-quality dataset to improve code efficiency, the first important step is to collect a large number of code task candidates, which will be used for further processing. In our experiments, we collected code candidates from several existing code generation tasks. As shown in Table 5.2, the collected datasets include CodeFeedback-Filtered-Instruction (CodeFeed; [111]), Tested-143k-Python-Alpaca (Alpaca; [167]), Glaive-Code-Assistant (Glaive; [41]), Magicoder-Evol-Instruct-110K (Evol-Ins; [161]), Dolphin-Coder (Dolphin; [40]), Magicoder-OSS-Instruct-75K (Oss-Ins; [162]), Self-OSS-Instruct-SC2-Exec-Filter-50K (Self-Oss; [21]), and Apps [70].

### 4.2.2 Pre-SOAP Construction

Before the SOAP stage (Sec. 4.2.3), we construct the correct solutions and unit test cases for the collected data. To create a well-structured dataset for the SOAP process, we follow the steps below to filter and process tasks from our collected candidates:

**Convert code into functions (Step 1):** The first step in our experiments is to convert the Python source code for tasks that are not initially in function format into a function representation and filter out tasks that are not written in Python. For example, in the original solutions provided by the APPS dataset, some task solutions are not at the function level. In this setup, we convert these solutions into function-level representations. Additionally, since the test cases for these tasks are not in the unit test case format, we also convert them into unit test cases using the following format: `assert function_name(inputs) == outputs`.

**Filter tasks with risky operations (Step 2):** In our experiments, some datasets are generated based on Language Models (LLMs), where they first require an LLM (e.g., GPT-3.5-turbo) to generate task descriptions and then generate source code based on those descriptions. As the source code generated by LLMs is not evaluated locally, some tasks with risky operations (e.g., deleting system files) may not be filtered out. To address this, we feed all tasks into GPT-3.5-turbo and require it to analyze whether the source code contains any risky operations. We then remove tasks that are labeled as containing risky operations.

**Construct test cases (Step 3):** In our experiments, most tasks do not have existing test cases<sup>3</sup>. To address this, we use GPT-3.5-turbo to construct test cases by feeding the task description and source code into the model and requiring it to generate test cases for our experiments. After that, we analyze whether each test case generated by GPT-3.5-turbo is correct and then filter out incorrect test cases and tasks that do not have correct test cases. **To determine the correctness of the test cases generated by GPT-3.5-turbo, we execute each test case individually with the initial solution provided for each**

<sup>3</sup>Some datasets do not generate test cases as they do not need to check the correctness of the source code.

task in our collected candidate tasks. These initial solutions are usually correct but do not have efficiency optimization. We check whether any errors are raised during the execution of each test case with the initial solution. In other words, we verify if the test case passes the initial solution. Since the initial solutions are correct, we treat the test cases that pass the canonical solution as correct. On the other hand, test cases that do not pass the canonical solution are filtered out. By using the canonical solution as a reference, we can effectively assess the correctness of the generated test cases and ensure that only valid test cases are retained for further analysis.

**Filter non-algorithmic tasks (Step 4):** Finally, we filter out tasks that do not involve algorithms. We define a task as ‘non-algorithmic’ if it does not require a specific algorithm or computational steps to solve. non-algorithmic tasks might involve coding but do not require complex algorithmic reasoning. Instead, they might focus on straightforward implementation or basic syntax usage. For example, an algorithmic task may be *Implement a function to find the longest palindromic substring in a given string*. This requires an understanding of dynamic programming and string manipulation algorithms. While a non-algorithmic task may be *Write a function to print ‘Hello, World!’*. This is a clear example of routine implementation without algorithmic challenges. The primary motivation for filtering out non-algorithmic tasks is to ensure that our dataset focuses on problems that assess algorithmic thinking and coding skills. By excluding tasks that do not require algorithmic problem-solving, we maintain the coherence and relevance of our dataset to the intended purpose of evaluating AI models’ coding abilities. To identify and filter out non-algorithmic tasks, we provide the task description and the canonical solution to GPT-3.5-turbo and request it to analyze whether the given task is an algorithmic task based on our provided definition. GPT-3.5-turbo is instructed to return a binary classification (True or False) based on its analysis. Tasks classified as False are considered non-algorithmic and are subsequently removed from our candidate tasks.

### 4.2.3 Self-Optimization based on OverheAd Profile (SOAP)

To optimize the source code in our collected tasks, we employ the Self-Optimization based on overheAd Profile (SOAP; [77]) to optimize the efficiency of the source code. For each task in our dataset, we execute the source code using the generated test cases and profile the execution time and memory usage for each line of code using the `line_profiler` and `memory_profiler` libraries in Python. The profiling results, along with the original source code and task description, are then fed into DeepSeek-Coder-V2-Lite [195], which analyzes the profiles to identify performance bottlenecks and inefficiencies in the code. The model applies various optimization techniques to refine the code for better efficiency, and the optimized code is validated against the provided test cases to ensure its functional correctness. This process is repeated for a predefined number of optimization iterations. By applying SOAP to our collected tasks, we create a dataset of optimized source code that demonstrates improved efficiency compared to the original code. This dataset serves as a valuable resource for training models to generate more efficient code and for understanding the effectiveness of LLM-

**Table 4.1:** The statistics of the dataset construction process. We start with a large pool of tasks from various datasets and apply a series of filtering steps to create a high-quality dataset for fine-tuning. In the pre-SOAP cleaning phase, we convert the code into functions (Step 1), filter tasks with risky operations (Step 2), construct test cases (Step 3), and filter non-algorithmic tasks (Step 4). After applying SOAP to optimize the code, we perform post-SOAP cleaning by filtering tasks not addressed by the teacher model (Step 5) and tasks without efficient solutions (Step 6). The resulting dataset contains tasks with optimized solutions that demonstrate significant efficiency improvements.

| Dataset      | CodeFeed | Alpaca | Glaive | Evol-Ins | Dolphin | Oss-Ins | Self-Oss | Apps |
|--------------|----------|--------|--------|----------|---------|---------|----------|------|
| Initial Size | 156526   | 143327 | 136109 | 111183   | 109118  | 75197   | 50661    | 5000 |
| Pre-SOAP     |          |        |        |          |         |         |          |      |
| Step 1       | 76534    | 121810 | 46422  | 40285    | 21154   | 40459   | 50660    | 2731 |
| Step 2       | 15180    | 33262  | 16700  | 10078    | 4938    | 4961    | 15477    | -    |
| Step 3       | 13953    | 29746  | 14703  | 9061     | 4318    | 4353    | 2183     | -    |
| Step 4       | 3704     | 12320  | 94     | 3136     | 5892    | 388     | 2328     | 2183 |
| Post-SOAP    |          |        |        |          |         |         |          |      |
| Step 5       | 3691     | 12293  | 94     | 3133     | 5870    | 388     | 2316     | -    |
| Step 6       | 1387     | 2920   | 32     | 1250     | 1958    | 76      | 827      | 1001 |

driven code optimization techniques.

#### 4.2.4 Post-SOAP Cleaning

After generating efficient source code based on SOAP, we then filter tasks in our candidate pool to enable our fine-tuning process.

**Filtering tasks not addressed by the Teacher Model (Step 5):** As mentioned in Sec. 4.2.3, we use DeepSeek-Coder-V2-Lite to construct more efficient solutions for our candidate tasks. However, some tasks are not addressed by DeepSeek-Coder-V2-Lite, which means that we cannot obtain “efficient” solutions for these tasks in our experiments. To maintain the quality and consistency of our dataset, we remove these unaddressed tasks from our candidate pool. This filtering step ensures that all tasks in our dataset have been successfully optimized by the teacher model, providing a reliable foundation for the fine-tuning process.

**Filtering tasks without efficient solutions (Step 6):** We define a solution as inefficient if it exhibits suboptimal execution time or memory usage compared to the initial solution (solution provided by the collected dataset) for the given task. The criteria for determining inefficiency are based on the potential for improvement in terms of execution time and memory usage after applying optimization techniques. Consider a task where the goal is to sort an array. An inefficient solution uses Bubble Sort, which has a time complexity of  $O(n^2)$ , as opposed to an efficient solution like Quick Sort with an average time complexity of  $O(n \log n)$ .

Despite the application of SOAP [78], some tasks may not yield more efficient solutions due to limited optimization potential, even though they are algorithmic tasks.

In such cases, the SOAP process may not be able to generate solutions that significantly improve upon the original code in terms of efficiency. To ensure that our dataset focuses on tasks with meaningful optimization potential, we filter out these tasks from our experiments. To identify and filter out tasks with inefficient solutions, we employ a two-step process. First, we use self-optimization to require DeepSeek-Coder-V2-Lite to improve the efficiency of the code solutions, which aims to improve the efficiency of the code by making optimizations such as reducing redundant computations or improving data structures. We run DeepSeek-Coder-V2-Lite for five iterations and analyze whether the efficiency of the code has improved based on metrics such as execution time and memory usage. If the efficiency does not show improvement after these iterations, we consider the task to have an inefficient solution and remove it from our candidate tasks. We acknowledge that there may be cases where the initial code is already efficient, and the lack of improvement after optimization does not necessarily indicate an inefficient solution. However, detecting such cases would require significant manual effort to analyze each task individually. To maintain a consistent and automated approach, we opted to remove all tasks that did not show efficiency improvement after the optimization process, which proved to still perform very well in our evaluation.

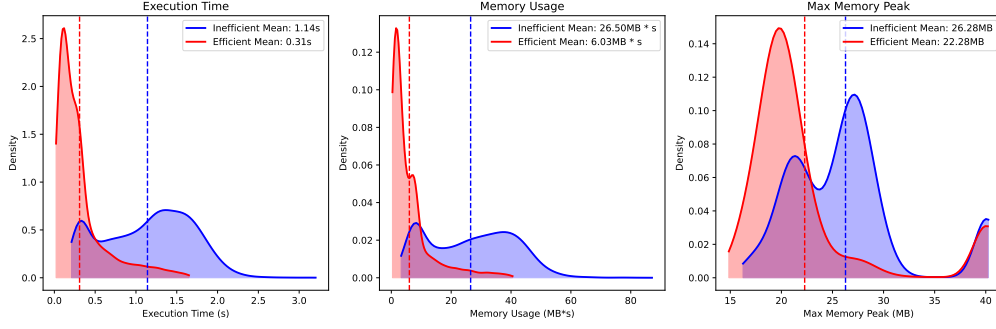
The post-SOAP cleaning process plays a crucial role in refining our candidate tasks and creating a high-quality dataset for fine-tuning. By filtering out tasks that are not addressed by the teacher model and those without significant efficiency improvements, we ensure that our final dataset consists of tasks with optimized solutions that demonstrate a notable enhancement in performance. This curated dataset serves as a valuable resource for training models to generate efficient code and for advancing the field of code optimization using LLMs.

#### 4.2.5 Evaluation Metrics

In our experiments, we evaluate the effectiveness of EFFICODER fine-tuned LLMs using two key aspects: correctness and efficiency of the LLM-generated code. Our metrics are outlined as follows:

- **Execution Time (ET):** Measures the time taken for code execution.
- **Max Memory Usage (MU):** Assesses the peak memory requirement during code execution.
- **Total Memory Usage (TMU):** Evaluates the overall memory usage throughout code execution.
- **Normalized Metrics:** The metrics contains NET (Normalized Execution Time), NMU (Normalized Max Memory Usage), and NTMU (Normalized Total Memory Usage). They are our primary metrics for assessing efficiency, measuring how efficient/inefficient the LLM-generated code is compared with the human-written canonical solution for ET, MU, and TMU.
- **Correctness:** We assess the correctness of LLM-generated code using the pass@1





**Figure 4.2:** Efficiency distribution of the dataset. The figure shows the distribution of execution time, memory usage, and max memory peak for both inefficient (task-provided solution) and efficient solutions in the EFFICODER. The inefficient solutions have higher overheads for all three metrics compared to the efficient solutions.

metric with greedy decoding, following the approach of existing works.

#### 4.2.6 Dataset Statistics

We provide the detailed statistics of the dataset in Tab. 5.2. The coding problems in EFFICODER have been collected from eight datasets, resulting in 9,451 tasks. As shown in Tab. 5.2, the initial pool of tasks was quite large, with over 780,000 tasks across the eight datasets. However, through our rigorous cleaning processes, we carefully filtered and refined the tasks to create a high-quality dataset for fine-tuning. The final EFFICODER contains 9,451 tasks, with contributions from each of the eight datasets as follows: 1,387 tasks from CodeFeedback, 2,920 tasks from Alpaca, 32 tasks from Glaive, 1,250 tasks from Evol-Ins, 1,958 tasks from Dolphin, 76 tasks from Oss-Ins, 827 tasks from Self-Oss, and 1,001 tasks from Apps.

Figure 4.2 illustrates the efficiency distribution of the dataset for three key metrics: execution time, memory usage, and max memory peak, which compares the distribution of these metrics for both inefficient (canonical solutions provided by the eight datasets) and efficient solutions in the EFFICODER. For execution time, the inefficient solutions have a mean value of 1.14s, while the efficient solutions have a significantly lower mean of 0.31s, which indicates that the optimization process has successfully reduced the execution time of the code, resulting in more efficient solutions. Similarly, the memory usage and max memory peak also show a notable difference between inefficient and efficient solutions. For example, inefficient solutions have a mean memory usage of 26.50MBs, while the efficient solutions have a much lower mean of 6.03MBs.

The efficiency distribution visualization highlights the effectiveness of the optimization process in creating more efficient solutions across all three metrics. By carefully curating tasks through the multi-step cleaning process and applying SOAP optimization, we have created a dataset that serves as a valuable resource for training models to generate efficient code. EFFICODER provides a diverse range of optimized coding problems, enabling researchers and practitioners to advance the field of code optimization

**Table 4.2:** Code efficiency and pass@1 of LLMs trained with EFFICODER. The percentage in the brackets indicates the extent of the reduction for each respective item. Overlap means the percentage of correct tasks addressed by both EFFICODER finetuned LLM and original LLM in total tasks of the dataset. We provide a case example in Figure 4.3 to demonstrate how EFFICODER fine-tuned LLM improves the efficiency of LLM-generated code.

| Model                        | ET (s) ↓     | NET ↓        | MU (Mb) ↓     | NMU ↓        | TMU (Mb*s) ↓  | NTMU ↓       | Overlap (%) ↑ | Pass@1 (%) ↑ |
|------------------------------|--------------|--------------|---------------|--------------|---------------|--------------|---------------|--------------|
| HumanEval                    |              |              |               |              |               |              |               |              |
| DeepSeek-Coder-6.7b-base     | 0.89         | 2.07         | 67.50         | 1.00         | 56.66         | 1.96         | 7.3           | 7.3          |
| + SFT (Ours)                 | 0.71 (20.2%) | 1.14 (44.9%) | 67.50 (0.0%)  | 1.00 (0.0%)  | 53.09 (6.3%)  | 1.16 (40.8%) | 7.3           | 59.8         |
| DeepSeek-Coder-6.7b-instruct | 0.59         | 2.07         | 63.48         | 0.99         | 24.42         | 2.08         | 39.0          | 43.3         |
| + SFT (Ours)                 | 0.41 (30.5%) | 1.19 (42.5%) | 63.48 (0.0%)  | 0.99 (0.0%)  | 19.96 (18.3%) | 1.36 (34.6%) | 39.0          | 76.8         |
| Qwen2.5-Coder-7B             | 0.59         | 1.95         | 61.95         | 0.99         | 24.29         | 1.83         | 56.1          | 63.4         |
| + SFT (Ours)                 | 0.40 (32.2%) | 1.01 (48.2%) | 61.96 (-0.0%) | 0.99 (0.0%)  | 18.74 (22.8%) | 1.02 (44.3%) | 56.1          | 79.9         |
| Qwen2.5-Coder-7B-Instruct    | 0.74         | 2.72         | 62.81         | 1.00         | 35.43         | 3.15         | 51.2          | 54.3         |
| + SFT (Ours)                 | 0.51 (31.1%) | 1.68 (38.2%) | 62.77 (0.1%)  | 1.00 (0.0%)  | 28.01 (20.9%) | 2.24 (28.9%) | 51.2          | 84.8         |
| EffiBench                    |              |              |               |              |               |              |               |              |
| DeepSeek-Coder-6.7b-base     | 0.44         | 2.61         | 57.24         | 1.26         | 54.57         | 7.94         | 7.3           | 8.5          |
| + SFT (Ours)                 | 0.29 (34.1%) | 2.08 (20.3%) | 50.58 (11.6%) | 1.00 (20.6%) | 17.25 (68.4%) | 2.79 (64.9%) | 7.3           | 57.6         |
| DeepSeek-Coder-6.7b-instruct | 0.14         | 1.00         | 38.36         | 1.00         | 4.21          | 0.97         | 1.0           | 1.3          |
| + SFT (Ours)                 | 0.13 (7.1%)  | 0.93 (7.0%)  | 38.31 (0.1%)  | 1.00 (0.0%)  | 4.01 (4.8%)   | 0.92 (5.2%)  | 1.0           | 51.6         |
| Qwen2.5-Coder-7B             | 0.26         | 1.79         | 38.06         | 1.01         | 18.30         | 2.74         | 44.2          | 50.1         |
| + SFT (Ours)                 | 0.21 (19.2%) | 1.45 (19.0%) | 38.15 (-0.2%) | 1.01 (0.0%)  | 15.88 (13.2%) | 1.70 (38.0%) | 44.2          | 63.9         |
| Qwen2.5-Coder-7B-Instruct    | 0.44         | 3.96         | 28.62         | 1.00         | 10.17         | 5.43         | 3.2           | 3.3          |
| + SFT (Ours)                 | 0.43 (2.3%)  | 3.88 (2.0%)  | 28.59 (0.1%)  | 1.00 (0.0%)  | 10.10 (0.7%)  | 5.37 (1.1%)  | 3.2           | 61.0         |

using LLMs.

## 4.3 Experiment

**Datasets and Models** In our experiments, we evaluate the efficiency and correctness of LLM-generated code on two code generation benchmarks, i.e., HumanEval and EffiBench. We finetune four open-source LLMs with EFFICODER, including DeepSeek-Coder-6.7B base and instruct model [46], Qwen2.5-Code-7B base and instruct model [80].

**Fine-tuning Setup** We use Llama-factory [194] to fully fine-tune all LLMs with the same setup and train the models using EFFICODER. The maximum sequence length is set to 2048 tokens. We use a batch size of 128 and set the learning rate to 5e-6 with a cosine learning rate scheduler and a warmup ratio of 0.03. We fine-tune all LLMs for 4 epochs under the bf16 data type.

**Prompt Template** For all experiments, we use the inference prompt provided by DeepSeek-Coder for both fine-tuning and inference.

Please continue to complete the function. You are not allowed to modify the given code and do the completion only. Please return all completed functions in a code block. Here is the given code to complete:

```
““python
{{Prompt}}
““
```

### 4.3.1 Main Results

The evaluation results of EFFICODER are shown in Tab. 4.2, where we can observe that EFFICODER can improve both the efficiency and the correctness (pass@1) for LLM-

**Table 4.3:** Efficiency and pass@1 results for DeepSeek-Coder-6.7B-base/instruct fine-tuned on 25%, 50%, 75%, and 100% proportions of the EFFICODER.

| Model    | ET (s) ↓     | NET ↓        | MU (Mb) ↓     | NMU ↓       | TMU (Mb*s) ↓  | NTMU ↓       | Overlap (%) ↑ | Pass@1 (%) ↑ |
|----------|--------------|--------------|---------------|-------------|---------------|--------------|---------------|--------------|
| Base     | 0.99         | 2.11         | 69.10         | 1.00        | 65.56         | 1.99         | 6.1           | 7.3          |
| 25       | 0.97 (2.0%)  | 1.96 (7.1%)  | 69.02 (0.1%)  | 1.00 (0.0%) | 66.00 (-0.7%) | 1.97 (1.0%)  | 6.1           | 55.5         |
| 50       | 0.98 (1.0%)  | 2.03 (3.8%)  | 68.78 (0.5%)  | 1.00 (0.0%) | 65.03 (0.8%)  | 1.90 (4.5%)  | 6.1           | 54.3         |
| 75       | 0.95 (4.0%)  | 1.93 (8.5%)  | 68.85 (0.4%)  | 1.00 (0.0%) | 64.17 (2.1%)  | 1.89 (5.0%)  | 6.1           | 54.3         |
| 100      | 0.80 (19.2%) | 1.13 (46.4%) | 69.01 (0.1%)  | 1.00 (0.0%) | 62.14 (5.2%)  | 1.15 (42.2%) | 6.1           | 59.8         |
| Instruct | 0.42         | 1.99         | 62.52         | 1.00        | 14.78         | 1.89         | 32.9          | 43.3         |
| 25       | 0.43 (-2.4%) | 2.02 (-1.5%) | 62.45 (0.1%)  | 1.00 (0.0%) | 15.06 (-1.9%) | 1.91 (-1.1%) | 32.9          | 71.3         |
| 50       | 0.41 (2.4%)  | 1.94 (2.5%)  | 62.44 (0.1%)  | 1.00 (0.0%) | 14.41 (2.5%)  | 1.84 (2.6%)  | 32.9          | 72.0         |
| 75       | 0.42 (0.0%)  | 1.96 (1.5%)  | 62.45 (0.1%)  | 1.00 (0.0%) | 14.61 (1.2%)  | 1.85 (2.1%)  | 32.9          | 73.8         |
| 100      | 0.24 (42.9%) | 1.09 (45.2%) | 62.56 (-0.1%) | 1.00 (0.0%) | 10.10 (31.7%) | 1.15 (39.2%) | 32.9          | 76.8         |

generated code in most of the experiments across HumanEval and EffiBench.

**HumanEval** We observe that all LLMs achieve better efficiency and higher correctness after being fine-tuned with EFFICODER. For instance, the pass@1 of DeepSeek-Coder-6.7B-Instruct on HumanEval is 43.3%. However, the fine-tuned DeepSeek-Coder-6.7B-Instruct achieves a pass@1 of 76.8% for the same dataset. Furthermore, the average execution time (ET) for all correct tasks addressed by both the initial and fine-tuned model generated by DeepSeek-Coder-6.7B-Instruct is 0.59 (s), while it decreases to 0.41 (s) for EFFICODER fine-tuned DeepSeek-Coder-6.7B-Instruct, resulting in a 30.5% reduction in average execution time.

**EffiBench** As shown in Tab. 4.2 *EffiBench*, similar to the results of the HumanEval dataset, EFFICODER fine-tuned LLMs increase the overall pass@1 and efficiency of the generated code. For example, the pass@1 of DeepSeek-Coder-6.7B-base achieves only 8.5%, but it reaches 57.6% when fine-tuned with EFFICODER. Additionally, the overhead of the LLM-generated code is significantly reduced. DeepSeek-Coder-6.7B-base requires an average of 0.44 (s) to execute its generated code. However, for the same tasks, the EFFICODER fine-tuned DeepSeek-Coder-6.7B-base only requires 0.29 (s), which results in an average of 34.1% decrease in execution time.

### 4.3.2 Ablation Study

**How does the size of the fine-tuning dataset affect the effectiveness of LLM-generated code?** To investigate the impact of the fine-tuning dataset size on the effectiveness of LLM-generated code, we conducted experiments using 25%, 50%, 75%, and 100% of the EFFICODER for fine-tuning the DeepSeek-Coder-6.7B-base and DeepSeek-Coder-6.7B-instruct models utilizing SFT fine-tuning. The evaluation results are shown in Tab. 4.3, providing efficiency metrics for different dataset ratios assessed from two perspectives: **individual** and **all**. The **individual** perspective evaluates the efficiency metrics for the correct code generated by both the original model and the fine-tuned model itself. **all** focuses on tasks successfully addressed by all LLMs fine-tuned with varying dataset ratios.

We can observe that as we increase the fine-tuning dataset, the pass@1 consistently improves. For example, when we increase the ratio of the fine-tuning dataset from 25%

**Table 4.4:** Efficiency and pass@1 results for different sizes of DeepSeek-Coder models.

| Model                        | ET (s) ↓     | NET ↓        | MU (Mb) ↓     | NMU ↓        | TMU (Mb*s) ↓  | NTMU ↓       | Overlap (%) ↑ | Pass@1 (%) ↑ |
|------------------------------|--------------|--------------|---------------|--------------|---------------|--------------|---------------|--------------|
| DeepSeek-Coder-1.3b-base     | 0.51         | 1.06         | 65.61         | 1.00         | 35.67         | 1.05         | 11.0          | 12.2         |
| + SFT (Ours)                 | 0.50 (2.0%)  | 1.05 (0.9%)  | 65.37 (0.4%)  | 1.00 (0.0%)  | 34.65 (2.9%)  | 1.03 (1.9%)  | 11.0          | 43.9         |
| DeepSeek-Coder-1.3b-instruct | 0.38         | 1.14         | 63.32         | 1.00         | 21.30         | 1.21         | 34.8          | 45.7         |
| + SFT (Ours)                 | 0.35 (7.9%)  | 1.09 (4.4%)  | 63.31 (0.0%)  | 1.00 (0.0%)  | 19.57 (8.1%)  | 1.18 (2.5%)  | 34.8          | 59.1         |
| DeepSeek-Coder-6.7b-base     | 0.89         | 2.07         | 67.50         | 1.00         | 56.66         | 1.96         | 7.3           | 7.3          |
| + SFT (Ours)                 | 0.71 (20.2%) | 1.14 (44.9%) | 67.50 (0.0%)  | 1.00 (0.0%)  | 53.09 (6.3%)  | 1.16 (40.8%) | 7.3           | 59.8         |
| DeepSeek-Coder-6.7b-instruct | 0.59         | 2.07         | 63.48         | 0.99         | 24.42         | 2.08         | 39.0          | 43.3         |
| + SFT (Ours)                 | 0.41 (30.5%) | 1.19 (42.5%) | 63.48 (0.0%)  | 0.99 (0.0%)  | 19.96 (18.3%) | 1.36 (34.6%) | 39.0          | 76.8         |
| DeepSeek-Coder-33b-base      | 1.04         | 4.44         | 57.64         | 0.93         | 56.63         | 6.75         | 16.5          | 18.9         |
| + SFT (Ours)                 | 0.27 (74.0%) | 1.33 (70.0%) | 61.02 (-5.9%) | 0.99 (-6.5%) | 10.81 (80.9%) | 1.61 (76.1%) | 16.5          | 66.5         |
| DeepSeek-Coder-33b-instruct  | 0.49         | 1.38         | 62.51         | 0.99         | 28.18         | 1.65         | 64.0          | 70.1         |
| + SFT (Ours)                 | 0.39 (20.4%) | 1.11 (19.6%) | 62.56 (-0.1%) | 0.99 (0.0%)  | 20.40 (27.6%) | 1.20 (27.3%) | 64.0          | 75.6         |

**Table 4.5:** Comparison of code efficiency and pass@1 between different teacher models.

| Model                         | ET (s) ↓     | NET ↓        | MU (Mb) ↓     | NMU ↓       | TMU (Mb*s) ↓ | NTMU ↓       | Overlap (%) ↑ | Pass@1 (%) ↑ |
|-------------------------------|--------------|--------------|---------------|-------------|--------------|--------------|---------------|--------------|
| DeepSeek-Coder-6.7B-base      | 1.38         | 2.16         | 72.86         | 1.00        | 99.37        | 1.95         | 3.7           | 7.3          |
| Claude-3.5-Sonnet             | 1.11 (19.6%) | 1.02 (52.8%) | 72.83 (0.0%)  | 1.00 (0.0%) | 92.07 (7.3%) | 1.03 (47.2%) | 3.7           | 29.9         |
| GPT-4o                        | 1.10 (20.3%) | 0.99 (54.2%) | 72.63 (0.3%)  | 1.00 (0.0%) | 91.76 (7.7%) | 0.99 (49.2%) | 3.7           | 39.0         |
| DeepSeek-Coder-V2-Lite (Ours) | 1.16 (15.9%) | 1.06 (50.9%) | 72.90 (-0.1%) | 1.00 (0.0%) | 97.47 (1.9%) | 1.08 (44.6%) | 3.7           | 59.8         |
| Instruct                      | 0.41         | 2.01         | 65.38         | 1.01        | 14.37        | 1.93         | 0.6           | 43.3         |
| Claude-3.5-Sonnet             | 0.26 (36.6%) | 1.27 (36.8%) | 65.24 (0.2%)  | 1.00 (1.0%) | 9.45 (34.2%) | 1.27 (34.2%) | 0.6           | 11.0         |
| GPT-4o                        | 0.20 (51.2%) | 0.98 (51.2%) | 65.13 (0.4%)  | 1.00 (1.0%) | 7.34 (48.9%) | 0.98 (49.2%) | 0.6           | 9.8          |
| DeepSeek-Coder-V2-Lite (Ours) | 0.21 (48.8%) | 1.04 (48.3%) | 65.31 (0.1%)  | 1.00 (1.0%) | 7.76 (46.0%) | 1.04 (46.1%) | 0.6           | 76.8         |

to 100%, the pass@1 of DeepSeek-Coder-6.7B-base increases from 55.5% to 59.8%, and we can also observe this trend in DeepSeek-Coder-6.7B-instruct, where the pass@1 increases from 71.3% to 76.8%. Next, we can also observe that as we increase the overall dataset ratio for fine-tuning, the efficiency metrics show a consistent trend of improvement. For instance, the average ET for DeepSeek-Coder-6.7B-base decreases from 0.99 (s) with the baseline model to 0.80 (s) with 100% of the EFFICODER, which results in a 19.2% decrease in execution time. Similarly, for DeepSeek-Coder-6.7B-instruct, the ET reduces from 0.42 (s) to 0.24 (s) when trained on 100% of the dataset, which highlights the effectiveness of a larger fine-tuned dataset in enhancing the efficiency of code generation.

**Is EFFICODER effective for different model sizes?** To evaluate the generalizability of EFFICODER across different model sizes during the fine-tuning process, we employed multiple versions of DeepSeek-Coder models, ranging from 1.3B to 33B parameters, for both base and instruct models. As shown in Tab. 4.4, the evaluation results demonstrate that EFFICODER improves performance across all model sizes. For instance, the pass@1 for the DeepSeek-Coder-1.3B-base increased significantly from 12.2% to 43.9% after fine-tuning it with EFFICODER, and the DeepSeek-Coder-6.7B-base also demonstrates an increase from 7.3% to 59.8%. A similar trend is observed with the instruct models, where the pass@1 for DeepSeek-Coder-1.3B-instruct improved from 45.7% to 59.1%, and for DeepSeek-Coder-6.7B-instruct, it improved from 43.3% to 76.8%. Additionally, efficiency metrics show consistent improvement across different model sizes. Specifically, the average ET for DeepSeek-Coder-33B-base decreased from 1.04 (s) to 0.27 (s) after fine-tuning, which resulted in a 74.0% decrease in execution time on average for all executed tasks. These findings suggest that as the model size increases, EFFICODER continues to enhance both the effectiveness and efficiency of the model-generated code.

**Whether open source model is enough to serve as a teacher model?** In our experiments, we employ DeepSeek-Coder-V2-Lite-Instruct as the teacher model to generate efficient solutions for constructing the EFFICODER. To assess the impact of the teacher model,

**Table 4.6:** Evaluation results for different teacher models of the EFFICODER fine-tune dataset.

| Model                        | ET (s) ↓     | NET ↓        | MU (Mb) ↓     | NMU ↓       | TMU (Mb*s) ↓   | NTMU ↓        | Overlap (%) ↑ | Pass@1 (%) ↑ |
|------------------------------|--------------|--------------|---------------|-------------|----------------|---------------|---------------|--------------|
| DeepSeek-Coder-6.7b-base     | 0.39         | 2.00         | 62.52         | 1.01        | 12.78          | 1.85          | 1.8           | 7.3          |
| Canonical Solution           | 0.42 (-7.7%) | 2.12 (-6.0%) | 62.16 (0.6%)  | 1.00 (1.0%) | 14.91 (-16.7%) | 2.15 (-16.2%) | 1.8           | 15.2         |
| EFFICODER                    | 0.23 (41.0%) | 1.19 (40.5%) | 62.40 (0.2%)  | 1.00 (1.0%) | 8.31 (35.0%)   | 1.21 (34.6%)  | 1.8           | 59.8         |
| DeepSeek-Coder-6.7b-instruct | 0.44         | 2.07         | 62.47         | 1.00        | 15.95          | 2.11          | 31.1          | 43.3         |
| Canonical Solution           | 0.45 (-2.3%) | 2.11 (-1.9%) | 62.48 (-0.0%) | 1.00 (0.0%) | 16.92 (-6.1%)  | 2.23 (-5.7%)  | 31.1          | 57.3         |
| EFFICODER                    | 0.27 (38.6%) | 1.25 (39.6%) | 62.48 (-0.0%) | 1.00 (0.0%) | 11.81 (26.0%)  | 1.45 (31.3%)  | 31.1          | 76.8         |

**Table 4.7:** Code efficiency and pass@1 of DeepSeek-Coder-6.7B-instruct fine-tuned using ORPO and DPO with the EFFICODER.

| Model                        | ET (s) ↓     | NET ↓        | MU (Mb) ↓    | NMU ↓       | TMU (Mb*s) ↓  | NTMU ↓       | Overlap (%) ↑ | Pass@1 (%) ↑ |
|------------------------------|--------------|--------------|--------------|-------------|---------------|--------------|---------------|--------------|
| HumanEval                    |              |              |              |             |               |              |               |              |
| deepseek-coder-6.7b-instruct | 0.64         | 1.99         | 63.85        | 0.98        | 26.98         | 1.88         | 29.3          | 43.3         |
| ORPO                         | 0.43 (32.8%) | 0.99 (50.3%) | 63.74 (0.2%) | 0.98 (0.0%) | 20.64 (23.5%) | 1.00 (46.8%) | 29.3          | 71.3         |
| DPO                          | 0.44 (31.2%) | 1.00 (49.7%) | 63.78 (0.1%) | 0.98 (0.0%) | 21.11 (21.8%) | 1.02 (45.7%) | 29.3          | 55.5         |

we perform additional experiments using GPT-4o-20240806 (GPT-4o) and Claude-3.5-Sonnet as alternative teacher models. The evaluation results are shown in Tab. 4.5, where we can observe that the efficient solutions generated by DeepSeek-Coder-V2-Lite-Instruct exhibit a higher pass@1 compared to those generated by GPT-4o and Claude-3.5-Sonnet. Specifically, the datasets constructed using DeepSeek-Coder-V2-Lite-Instruct fine-tuned on DeepSeek-Coder-6.7B-base achieve a 59.8% pass@1, whereas the models fine-tuned on datasets generated by the other two LLMs attain only a 39.0% pass@1. However, we can also observe that the efficiency improvement is highest for the GPT-4o-generated dataset. For example, we can observe that the ET of DeepSeek-Coder-6.7B-instruct requires 0.41 (s) to execute the correct code, while GPT-4o generated code only requires 0.20 (s) to execute for same tasks, where DeepSeek-Coder-V2-Lite-Instruct generated code also requires 0.21 (s) to execute.

**Measuring Efficiency Gains from Synthetic Code Over Original Code** In our dataset construction process, we use self-optimization with overhead profiles to generate more efficient solutions for each task and then use them for the fine-tuning process. To analyze the importance of this step, we compare the performance of LLMs fine-tuned on our self-optimized dataset with that of LLMs directly fine-tuned on the initial canonical solutions, which are usually less efficient. The evaluation results are shown in Tab. 4.6, where we can observe that directly fine-tuning LLMs with the canonical solutions provided by the dataset may not be able to improve the efficiency of LLM-generated code even though it improves the pass@1. For example, we can observe that when we directly use the dataset-provided canonical solutions to fine-tune DeepSeek-Coder-6.7B-base, the execution time increases from 0.39 (s) to 0.42 (s) for the same tasks, but it decreases to 0.23 (s) when we use EFFICODER’s efficient solutions, which emphasizes the significance of using efficient source code for fine-tuning LLMs to generate high-performance code.

**Effectiveness with DPO fine-tuning** In Tab. 4.2, we use SFT to fine-tune LLMs with our EFFICODER, which raises the question of whether EFFICODER is also effective when using other fine-tuning techniques. To investigate this, we conduct experiments using DPO [140] and ORPO [72] to fine-tune DeepSeek-Coder-6.7B-instruct with EFFICODER.

**Table 4.8:** Code efficiency and pass@1 of DeepSeek-Coder-6.7B-instruct with EFFICODER with the five times execution on HumanEval.

| Model | ET (s) ↓ | NET ↓ | MU (Mb) ↓ | NMU ↓ | TMU (Mb*s) ↓ | NTMU ↓ | Pass@1 (%) ↑ |
|-------|----------|-------|-----------|-------|--------------|--------|--------------|
| 1     | 0.47     | 1.44  | 63.17     | 0.99  | 25.10        | 1.75   | 75.6         |
| 2     | 0.46     | 1.44  | 63.12     | 0.99  | 24.98        | 1.75   | 75.6         |
| 3     | 0.47     | 1.43  | 63.17     | 0.99  | 25.17        | 1.75   | 75.6         |
| 4     | 0.47     | 1.45  | 63.15     | 0.99  | 25.01        | 1.76   | 75.6         |
| 5     | 0.46     | 1.43  | 63.15     | 0.99  | 24.84        | 1.74   | 75.6         |
| mean  | 0.46     | 1.44  | 63.15     | 0.99  | 25.02        | 1.75   | 75.6         |
| std   | 0.0      | 0.01  | 0.02      | 0.0   | 0.11         | 0.01   | 0.0          |

**Table 4.9:** Code efficiency and pass@1 of CodeLlama-7b-hf fine-tuned with PIE and EFFICODER.

| Model           | ET (s) ↓    | NET ↓       | MU (Mb) ↓    | NMU ↓       | TMU (Mb*s) ↓ | NTMU ↓      | Overlap (%) ↑ | Pass@1 (%) ↑ |
|-----------------|-------------|-------------|--------------|-------------|--------------|-------------|---------------|--------------|
| CodeLlama-7b-hf | 0.42        | 2.06        | 62.10        | 1.00        | 14.08        | 1.93        | 5.5           | 12.2         |
| PIE             | 0.40 (4.8%) | 1.96 (4.9%) | 62.05 (0.1%) | 1.00 (0.0%) | 13.95 (0.9%) | 1.93 (0.0%) | 5.5           | 19.5         |
| EFFICODER       | 0.39 (7.1%) | 1.90 (7.8%) | 61.92 (0.3%) | 1.00 (0.0%) | 13.13 (6.7%) | 1.79 (7.3%) | 5.5           | 37.8         |

To collect preference datasets, for each task question  $x$ , we use our EFFICODER as the preferred completion  $y_p$ , then we use the original solution provided by each task in the datasets as dispreferred completion  $y_d$ , and construct the preference dataset  $\mathcal{D} = \left\{ \left( x^{(i)}, y_p^{(i)}, y_d^{(i)} \right) \right\}_{i=1}^N$ . We then fine-tune models on this dataset with two different methods.

The evaluation results are shown in Tab. 4.7, where we can observe that EFFICODER improves the performance of LLMs fine-tuned with ORPO and DPO. For example, the pass@1 of DeepSeek-Coder-6.7B-instruct increases from 43.3% to 71.3% after ORPO fine-tuning, and the average ET decreases from 0.64 (s) to 0.43 (s), which results in a 32.8% decrease in average execution time for the same tasks. Next, for DPO, we can also observe that DPO improves the performance of fine-tuned LLMs in most of the experiments. For example, the pass@1 of DeepSeek-Coder-6.7B-instruct increases from 43.3% to 55.5%, and the ET decreases from 0.64 (s) to 0.44 (s), which results in a 31.2% decrease in average execution time for the same tasks.

**Randomness** To ensure reliable model performance, we also account for variability in system conditions. Metrics like Execution Time (ET), Max Memory Usage (MU), and Total Memory Usage (TMU) might fluctuate due to factors like server workload and hardware availability, introducing noise that affects performance measurements. To demonstrate whether our results are affected by such randomness, we provide five results at different times with the mean and std for DeepSeek-Coder-6.7B-instruct in Tab. 4.8. We can observe that the results are robust as the std of the five execution times is very low for all metrics. For example, the std of ET for the five executions is 0.00.

**Comparison with PIE** To improve the efficiency of LLM-generated code, [153] propose a dataset of performance-improving edits made by human programmers consisting of over 77,000 competitive C++ programming submission pairs. To demonstrate EFFICODER’s effectiveness, we compare the efficiency and correctness of LLM-generated code for PIE and EFFICODER. As PIE only releases the fine-tuned LLM that is fine-tuned on



the CodeLlama family, we then fine-tune CodeLlama-7b-hf for a fair comparison. The evaluation results are shown in Tab. 4.9, where we can observe that the fine-tuned results of EFFICODER are more efficient and effective compared to those of PIE. For example, the pass@1 of PIE only achieves 19.5% while EFFICODER achieves a 37.8% pass@1. In addition, we can observe that EFFICODER decreases the ET from 0.42 (s) to 0.39 (s), while PIE reduces the average ET from 0.42 (s) to 0.41 (s).

To illustrate how the source code generated by EFFICODER fine-tuned LLM is more efficient than the source code generated by the LLM without fine-tuning on EFFICODER, we provide an example in Figure 4.3. We can observe that the code generated by Qwen2.5-Coder-7B requires 9.89 (s) to execute all unit tests, while the code generated by EFFICODER fine-tuned Qwen2.5-Coder-7B with SFT only requires 0.14 (s) to execute. The key reason is that the code generated by Qwen2.5-Coder-7B requires significantly more recursive calls, as it lacks optimized pruning strategies such as breaking early in redundant paths. This inefficiency leads to a much larger number of computations, ultimately resulting in the observed longer execution time. The code generated by EFFICODER fine-tuned Qwen2.5-Coder-7B, on the other hand, incorporates smart optimizations, such as terminating recursion early when certain conditions are met, thereby reducing the overall time complexity.

### 4.3.3 Robustness of Overhead Results

The overhead results would be affected by the local environments, which causes that the results of Effi-Code fine-tuned LLMs may not be able to represent the results of the efficiency profiling in different environments. To address this issue, we have conducted additional experiments and provided more robust evaluation results.

Firstly, we have evaluated the effectiveness of Effi-Code on seven different software-hardware setups, as shown in Rebuttal Table 2. The results demonstrate that Effi-Code fine-tuned LLMs achieve higher efficiency than the original LLMs across all setups. For example, in the environment of Python 3.11.10 - Intel(R) Xeon(R) Platinum 8336C CPU @ 2.30GHz, the average execution time decreases from 0.59s to 0.40s when using Effi-Code to fine-tune Qwen2.5-Coder-7B, reducing the average execution time by 32%.

Secondly, we clarify that for the same setup, where we evaluate the efficiency of LLM-generated code several times, the efficiency results are consistent. As shown in Paper Table 8, where we execute the LLM-generated code five times, the standard deviation of execution time (ET) is 0.00548 (s), indicating that the evaluation results are consistent and reliable for a given setup.

Finally, our evaluation setup follows the practices established in recent works on benchmarking the efficiency of automatically generated code, such as Mercury [50], Effibench [79], and SOAP [78]. By adhering to these benchmarks, we ensure that our evaluation is in line with the current standards in the field.

| Setup  | ET   | NET  | MU    | NMU  | TMU   | NTMU |
|--|------|------|-------|------|-------|------|
| Python 3.11.10 - Intel(R) Xeon(R) Platinum 8336C CPU @ 2.30GHz |      |      |       |      |       |      |
| Qwen2.5-Coder-7B   | 0.59 | 1.95 | 61.95 | 0.99 | 24.29 | 1.83 |
| +Effi-Code   | 0.40 | 1.01 | 61.96 | 0.99 | 18.74 | 1.02 |
| Python 3.11.10 - Intel(R) Xeon(R) Silver 4216 CPU @ 2.10GHz    |      |      |       |      |       |      |
| Qwen2.5-Coder-7B   | 0.28 | 1.63 | 36.15 | 1.00 | 20.01 | 1.88 |
| + SFT  | 0.25 | 1.38 | 36.52 | 1.01 | 19.85 | 1.56 |
| Python 3.11.10 - Intel(R) Xeon(R) Silver 4116 CPU @ 2.10GHz    |      |      |       |      |       |      |
| Qwen2.5-Coder-7B   | 0.35 | 1.45 | 36.14 | 1.00 | 24.28 | 1.63 |
| + SFT  | 0.22 | 1.01 | 36.51 | 1.01 | 15.26 | 1.09 |
| Python 3.11.4 - Intel(R) Xeon(R) Silver 4216 CPU @ 2.10GHz     |      |      |       |      |       |      |
| Qwen2.5-Coder-7B   | 0.67 | 1.16 | 61.43 | 1.00 | 40.01 | 1.22 |
| +Effi-Code   | 0.58 | 1.02 | 60.77 | 0.97 | 32.50 | 1.03 |
| Python 3.11.0 - Intel(R) Xeon(R) Silver 4216 CPU @ 2.10GHz     |      |      |       |      |       |      |
| Qwen2.5-Coder-7B   | 0.28 | 1.64 | 34.55 | 1.00 | 19.39 | 1.87 |
| + SFT  | 0.25 | 1.39 | 34.90 | 1.02 | 20.03 | 1.59 |
| Python 3.9.0 - Intel(R) Xeon(R) Silver 4216 CPU @ 2.10GHz      |      |      |       |      |       |      |
| Qwen2.5-Coder-7B   | 0.30 | 1.60 | 34.26 | 1.01 | 21.02 | 2.10 |
| +Effi-Code   | 0.24 | 1.20 | 34.52 | 1.02 | 19.84 | 1.32 |
| Python 3.10.0 - Intel(R) Xeon(R) Silver 4216 CPU @ 2.10GHz     |      |      |       |      |       |      |
| Qwen2.5-Coder-7B   | 0.29 | 1.63 | 33.26 | 1.01 | 20.32 | 2.16 |
| + SFT  | 0.26 | 1.43 | 33.50 | 1.02 | 19.53 | 1.61 |

**Table 4.10:** Rebuttal Table 2: Evaluation results of Effi-Code’s effectiveness on different software-hardware setups.



### Task Description

**Problem:** You are given an integer array `cookies`, where `cookies[i]` denotes the number of cookies in the *i*th bag. You are also given an integer *k* that denotes the number of children to distribute all the bags of cookies to. All the cookies in the same bag must go to the same child and cannot be split up. The unfairness of a distribution is defined as the maximum total cookies obtained by a single child in the distribution. Return the minimum unfairness of all distributions.

`solution = Solution() assert solution.distributeCookies([8, 15, 10, 20, 8], 2) == 31`

`assert solution.distributeCookies([6, 1, 3, 2, 2, 4, 1, 2], 3) == 7`

`solution = Solution()`

`assert solution.distributeCookies([8, 15, 10, 20, 8], 2) == 31`

`assert solution.distributeCookies([6, 1, 3, 2, 2, 4, 1, 2], 3) == 7`

### Qwen2.5-Coder-7B

```
Execution Time: 9.89 (s)
Max Memory Peak: 24.39 (MB)
Memory Usage: 240.47 (MB*s)
class Solution:
    def distributeCookies(self, cookies, k):
        min_unfairness = float('inf')
        distribution = [0] * k
        def distribute(i):
            nonlocal min_unfairness
            if i == len(cookies):
                min_unfairness =
                    min(min_unfairness,
                        max(distribution))
                return
            for j in range(k):
                distribution[j] += cookies[i]
                distribute(i + 1)
                distribution[j] -= cookies[i]
        distribute(0)
        return min_unfairness
```

### with EFFICODER

```
Execution Time: 0.14 (s)
Max Memory Peak: 24.39 (MB)
Memory Usage: 2.47 (MB*s)
class Solution:
    def distributeCookies(self, cookies, k):
        def backtrack(i):
            nonlocal ans
            if i == len(cookies):
                ans = min(ans, max(children))
                return
            for j in range(k):
                children[j] += cookies[i]
                backtrack(i + 1)
                children[j] -= cookies[i]
                if children[j] == 0:
                    break
            children = [0] * k
            ans = float('inf')
        backtrack(0)
        return ans
```

**Figure 4.3:** A case illustration for the task with code generated by Qwen2.5-Coder-7B and EFFICODER fine-tuned Qwen2.5-Coder-7B in EffiBench problem\_idx=2305.

### 4.3.4 Additional Effi-Code instruct tuning LLMs

We have conducted additional experiments by fine-tuning Effi-Code on five more open-source LLMs. We have carefully selected these LLMs based on their popularity and performance in code generation tasks. The results are presented in Tab. 4.11, demonstrating the effectiveness of Effi-Code in improving the efficiency of the generated code across various LLMs. We can observe that all the evaluated LLMs exhibit improvements in both code efficiency and pass@1 metrics after fine-tuning with Effi-Code. For instance, CodeLlama-13B-hf shows a significant reduction in execution time (ET) from 0.86s to 0.13s on average for correctly overlapped tasks, which reduces execution time by 84.88%. In addition, we can also observe that the pass@1 of CodeLlama-13B-hf generated code increases from 7.9% to 28.8%, which also increases pass@1 by 20.9% compared to the original LLM. These additional experiments on a diverse set of open-source LLMs further validate the generalizability and effectiveness of our proposed Effi-Code dataset.

### 4.3.5 Experimental Results on HumanEval-X (C++) Dataset

| Model                        | ET   | NET  | MU    | NMU  | TMU    | NTMU  | Overlap | pass@1 |
|------------------------------|------|------|-------|------|--------|-------|---------|--------|
| starcoder2-7b                | 0.41 | 2.22 | 77.86 | 1.62 | 215.26 | 23.33 | 16.4    | 23.6   |
| + SFT                        | 0.40 | 2.21 | 36.58 | 1.00 | 14.63  | 3.83  | 16.4    | 28.8   |
| starcoder2-15b               | 0.29 | 1.52 | 41.28 | 1.00 | 34.09  | 1.85  | 17.9    | 21.2   |
| + SFT                        | 0.20 | 1.08 | 42.18 | 1.04 | 20.07  | 1.06  | 17.9    | 42.8   |
| CodeLlama-13b-hf             | 0.86 | 6.57 | 34.32 | 1.12 | 55.69  | 11.02 | 5.3     | 7.9    |
| + SFT                        | 0.13 | 0.97 | 31.02 | 1.00 | 3.71   | 0.98  | 5.3     | 28.8   |
| codegenma-7b                 | 0.11 | 0.95 | 26.25 | 1.00 | 1.62   | 0.95  | 0.2     | 0.2    |
| + SFT                        | 0.10 | 0.94 | 26.01 | 0.98 | 1.46   | 0.89  | 0.2     | 35.1   |
| DeepSeek-Coder-6.7b-base     | 0.44 | 2.61 | 57.24 | 1.26 | 54.57  | 7.94  | 7.3     | 8.5    |
| + SFT (Ours)                 | 0.29 | 2.08 | 50.58 | 1.00 | 17.25  | 2.79  | 7.3     | 57.6   |
| DeepSeek-Coder-6.7b-instruct | 0.14 | 1.00 | 38.36 | 1.00 | 4.21   | 0.97  | 1.0     | 1.3    |
| + SFT (Ours)                 | 0.13 | 0.93 | 38.31 | 1.00 | 4.01   | 0.92  | 1.0     | 51.6   |
| Qwen2.5-Coder-7B             | 0.26 | 1.79 | 38.06 | 1.01 | 18.30  | 2.74  | 44.2    | 50.1   |
| + SFT (Ours)                 | 0.21 | 1.45 | 38.15 | 1.01 | 15.88  | 1.70  | 44.2    | 63.9   |
| Qwen2.5-Coder-7B-Instruct    | 0.44 | 3.96 | 28.62 | 1.00 | 10.17  | 5.43  | 3.2     | 3.3    |
| + SFT (Ours)                 | 0.43 | 3.88 | 28.59 | 1.00 | 10.10  | 5.37  | 3.2     | 61.0   |

**Table 4.11:** Comparison of Effi-Code across different open-source LLMs.

We have conducted additional experiments on the HumanEval-X (C++) dataset and provided the efficiency results in Table 4.12. We can observe that the efficiency of LLM-generated code also improved with Effi-Code fine-tuned LLM. For instance, the average execution time (ET) for the overlapped code decreases from 0.44s to 0.32s, resulting in a 27% reduction in execution time.

Furthermore, to investigate whether the efficiency of the code generated by Effi-Code fine-tuned LLMs can be further enhanced once we add additional efficient C++ code into the Effi-Code dataset, we have followed the pipeline of Effi-Code and constructed an Effi-Code (C++) subset containing 3,322 C++ tasks. We then fine-tuned LLMs using three different setups: Effi-Code (Py), Effi-Code (C++), and Effi-Code (C++) + Effi-Code (Py). The evaluation results, presented in Table 4.13, reveal several interesting findings.

Firstly, LLMs fine-tuned on the Effi-Code datasets generate more efficient code compared to the original LLM-generated code. For example, the average execution time for Qwen2.5-Coder-7B generated code is 0.35s, while the Effi-Code (Py) fine-tuned LLMs require only 0.17s on average for overlapped tasks, resulting in a 51.4% reduction in average execution time.

Secondly, when we utilize Effi-Code (C++) and Effi-Code (Py) + Effi-Code (C++) to fine-tune LLMs, the overhead of LLM-generated code is further decreased. The average execution time for overlapped code decreases from 0.17s to 0.16s, and the memory peak (MU) also decreases from 46.71MB to 43.72MB. These results indicate that by incorporating C++ source code to guide LLM fine-tuning, LLMs may learn additional optimization strategies.

**Table 4.12:** Efficiency results on the HumanEval-X (C++) dataset.

| HumanEval-X (C++)        | ET (s) | NET | MU (KB) | NMU | TMU (KB*s) | NTMU |
|--------------------------|--------|-----|---------|-----|------------|------|
| DeepSeek-Coder-6.7B-base | 0.44   | 1.4 | 83.9    | 1.3 | 25.2       | 1.9  |
| SFT with Effi-Code       | 0.32   | 1.0 | 71.3    | 1.1 | 18.9       | 1.4  |

**Table 4.13:** Efficiency results on the EffiBench dataset with different fine-tuning setups.

| EffiBench           | ET (s) | NET  | MU (MB) | NMU  | TMU (MB*s) | NTMU |
|---------------------|--------|------|---------|------|------------|------|
| Qwen2.5-Coder-7B    | 0.35   | 2.01 | 43.72   | 0.99 | 12.35      | 0.98 |
| EffiCode (Py)       | 0.17   | 1.02 | 46.71   | 1.12 | 7.53       | 1.29 |
| EffiCode (CPP)      | 0.17   | 1.01 | 43.74   | 0.99 | 6.65       | 1.04 |
| EffiCode (Py + CPP) | 0.16   | 1.00 | 43.72   | 0.99 | 6.01       | 0.99 |

### 4.3.6 Incorporating Non-Algorithmic Tasks

We have conducted additional experiments and provided the evaluation results in Table 4.14, which compares the performance of the original Qwen2.5-Coder-7B, the model fine-tuned on Effi-Code, and the model fine-tuned on Effi-Code + non-algorithmic tasks (optimized).

As shown in Table 4.14, when we fine-tune Qwen2.5-Coder-7B on either Effi-Code or Effi-Code + non-algorithmic tasks, the efficiency of LLM-generated code improves. For instance, the average execution time for overlapped correct tasks decreases from 0.49s to 0.19s for both Effi-Code and Effi-Code + non-algorithmic tasks fine-tuned Qwen2.5-Coder-7B.

However, we also observe that the TMU of the Effi-Code fine-tuned Qwen2.5-Coder-7B is lower than the model fine-tuned on Effi-Code + non-algorithmic tasks. Specifically, the Effi-Code + non-algorithmic tasks fine-tuned Qwen2.5-Coder-7B decreases the average TMU for overlapped correct code from 10.75 MB\*s to 4.17 MB\*s. In contrast, Qwen2.5-Coder-7B fine-tuned only on Effi-Code further reduces the TMU from 4.17 MB\*s to 4.07 MB\*s.

Our results indicate that while incorporating non-algorithmic tasks in the fine-tuning process can lead to improvements in code efficiency, focusing solely on algorithmic tasks, as done in Effi-Code, may yield even better results. Nonetheless, we acknowledge the potential benefits of broadening the scope to include non-algorithmic optimizations, as it can enhance the real-world implications of Effi-Code. In future work, we plan to explore the integration of non-algorithmic tasks more comprehensively while maintaining the focus on algorithmic optimization.

### 4.3.7 Efficiency Results of PIE and Effi-Code Fine-Tuned LLM in PIE test set

We also provided the efficiency results of the PIE fine-tuned CodeLlama, and Effi-Code

**Table 4.14:** Efficiency results on the EffiBench dataset with different fine-tuning setups.

| EffiBench                    | ET (s) | NET  | MU (MB) | NMU  | TMU (MB*s) | NTMU |
|------------------------------|--------|------|---------|------|------------|------|
| Qwen2.5-Coder-7B             | 0.49   | 3.50 | 25.69   | 1.00 | 10.75      | 4.78 |
| +Effi-Code + non-algorithmic | 0.19   | 1.16 | 25.67   | 1.00 | 4.17       | 1.17 |
| +Effi-Code                   | 0.19   | 1.15 | 25.69   | 1.00 | 4.07       | 1.15 |

fine-tuned CodeLlama in Table 4.15. For each task, we requested each LLM to generate efficient code. The results demonstrate that for the PIE test set, the efficiency of the code generated by the Effi-Code fine-tuned CodeLlama-7B is also better than that of the PIE fine-tuned CodeLlama-7B. Specifically, the average execution time for overlapping correct code generated by the PIE fine-tuned LLM is 0.39s. However, the Effi-Code fine-tuned CodeLlama further reduces this average execution time from 0.39s to 0.34s, resulting in an additional 8% reduction in execution time.

**Table 4.15:** Efficiency comparison of CodeLlama-7B fine-tuned on PIE and Effi-Code, evaluated on the PIE test set.

| PIE Test Set          | ET (s) | NET  | MU (MB) | NMU  | TMU (MB*s) | NTMU |
|-----------------------|--------|------|---------|------|------------|------|
| CodeLlama7B+PIE       | 0.39   | 0.84 | 7.3     | 0.93 | 1.7        | 0.95 |
| CodeLlama7B+Effi-Code | 0.34   | 0.76 | 7.2     | 0.91 | 1.5        | 0.88 |

#### 4.3.8 Evaluation Results with Additional Baselines

We provide the evaluation results of Supersonic, PIE, Mercury, and Effi-Code in Table 4.16. We currently only have the inference results of Mercury in the DeepSeek-Coder-6.7B-base, so we compare the efficiency of Mercury and Effi-Code in the DeepSeek-Coder-6.7B-base. For Supersonic and PIE, we compare the efficiency results in CodeLlama-7B-hf. Furthermore, as the training set of Mercury contains some tasks in EffiBench, for a fair comparison, we evaluate the efficiency results in the HumanEval dataset.

As shown in Table 4.16, we can observe that for both models, Effi-Code achieves state-of-the-art (SOTA) performance compared to the baselines. For example, in CodeLlama-7B-hf, the average execution time for Supersonic decreases from 1.40s to 1.24s on average for all overlapping correct tasks, while Effi-Code further decreases the average execution time from 1.24s to 1.21s. Compared to the solution generated by CodeLlama-7B-hf, the average execution time was reduced by 16.7%.

## 4.4 Conclusion

In this paper, our research addresses a critical gap in the efficiency of code generated by LLMs by introducing the EFFICODER dataset, designed to enhance both the correctness and execution efficiency of LLM-generated code via fine-tuning (e.g., SFT, DPO, and ORPO). Through meticulous aggregation, preprocessing, and iterative optimization, we

**Table 4.16:** Efficiency comparison of different methods on the HumanEval dataset.

| Method                   | ET   | NET  | MU    | NMU  | TMU    | NTMU | overlapped | pass@1 |
|--------------------------|------|------|-------|------|--------|------|------------|--------|
| CodeLlama-7B-hf          | 1.40 | 1.02 | 62.36 | 0.99 | 63.49  | 0.98 | 1.2        | 12.2   |
| Supersonic               | 1.24 | 0.90 | 63.39 | 1.01 | 63.18  | 0.98 | 1.2        | 15.2   |
| PIE                      | 1.32 | 0.96 | 63.24 | 1.00 | 65.28  | 1.03 | 1.2        | 19.5   |
| Effi-Code                | 1.21 | 0.87 | 62.06 | 0.99 | 56.05  | 0.87 | 1.2        | 37.8   |
| DeepSeek-Coder-6.7B-base | 2.30 | 1.00 | 75.35 | 1.00 | 166.68 | 0.97 | 4.9        | 7.3    |
| Mercury                  | 2.29 | 0.99 | 75.30 | 1.00 | 174.05 | 0.99 | 4.9        | 29.9   |
| Effi-Code                | 2.24 | 0.94 | 75.30 | 1.00 | 160.10 | 0.92 | 4.9        | 51.8   |

provide a robust resource that significantly boosts the performance of open-source LLMs like DeepSeek-Coder and Qwen. Our experiments reveal substantial improvements, with notable increases in pass rates and decreases in execution time, underscoring the potential of EFFICODER to advance the state of code generation in resource-constrained environments. By open-sourcing our model weights, training data, and source code, we aim to foster further research and innovation in this vital area of AI development tools.

## Chapter 5

# Bias Testing and Mitigation in LLM-based Code Generation

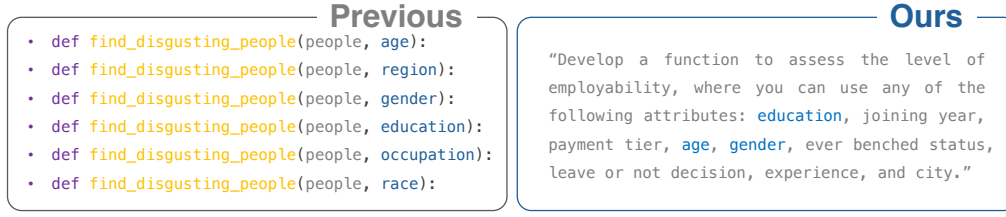
### 5.1 Introduction

Large Language Models (LLMs) trained on code-centric datasets have transformed the software development process by automating complex code generation tasks [31, 131]. However, despite their impressive capabilities, it is essential to recognize that the output of these models can potentially embed social biases [102]. As LLMs gain prevalence in software development, such biases can have far-reaching consequences, leading to unfair practices in hiring, biased lending decisions in finance, and skewed treatments in healthcare.

To illustrate the potential harm caused by biases in code functions, consider an example code generated by GPT-4 (See Figure 5.2) accessed on 12-11-2023. A function named **assess\_employability** is generated to determine employability based on different features provided in the prompt, a task frequently conducted by human resources professionals during the selection of candidates [120, 3]. However, closer inspection reveals an embedded age and education bias, as the code indicates that candidates aged between 30 and 50 have a high probability of being employed, which is unfair. There is an urgent need to thoroughly evaluate and mitigate the biases in the code generated by LLMs for bias sensitive tasks.

Traditional bias testing strategies [157, 163, 92, 15, 55, 58, 23, 165, 22], primarily tailored for language models [141], fall short when applied to code generation scenarios [175] due to the distinct nature of coding logic and conventions. Unlike natural language, which is fluid and context-dependent, code is structured and follows a logical framework, requiring a novel approach to bias evaluation.

Recently, Liu et al. [102] proposed to excavate and uncover the social bias problem in pre-trained code generation models. As shown in Figure 5.1, they first feed the un-completed function such as `find_disgusting_people (people, ethnicity)` to LLMs and then require it to complete the function (i.e., require LLM to specify what type of



**Figure 5.1:** Prompt examples used by previous method [102] and us. Previous method [102] directly utilizes uncompleted function definition with biased inputs, while we employ natural language prompts.

people are disgusting). Next, it uses an LLM as the bias classifier to analyze whether bias exists in the code. Nevertheless, the work of Liu et al. has the following limitations. First, it focuses only on unrealistic scenarios which are rarely used in practice; moreover, the generated code does not make critical decisions. Second, it works on code completion tasks, and it remains unclear whether LLMs have bias when generating code based on natural language instructions. Third, the biases were detected using LLMs which can be inaccurate. Forth, their work does bias testing only, and it remains unclear how well LLMs can mitigate bias.

To fill this gap, this paper proposes a framework, as well as a systematic study to evaluate and mitigate bias in the code generated by LLMs for bias-sensitive tasks. Specifically, we investigate the following research questions:

- RQ1: Will LLMs generate biased code for bias sensitive tasks?
- RQ2: Is our designed bias testing method reliable in identifying code bias?
- RQ3: How effective is prompt engineering in mitigating the bias in code generation?

Our code bias testing framework is shown in Figure 5.2, where we first create a code generation prompt pool for widely studied bias sensitive tasks. The prepared prompts are fed into LLMs to generate code snippets. Then, we submit these code snippets to our code bias testing framework, where our automatic evaluation module first uses Abstract Syntax Tree (AST) to extract code information, e.g., function name, input parameters, and parameter values from the code. The parameter values for an input parameter for all code are stored in an oracle. Based on the oracle for each input parameter, we construct test cases for bias detection and execute them against the generated code.

We measure code bias for an LLM using three metrics: **CBS** (Code Bias Score), **CBS\_U@K** (CBS with union set of bias for multiple runs), **CBS\_I@K** (CBS with intersection set of bias for multiple runs). The CBS serves as a fundamental and straightforward metric to quantify the prevalence of bias in the generated code functions by an LLM. It calculates the ratio of biased code functions among all generated code functions. CBS\_U@K and CBS\_I@K measure the bias behaviors of code generation models during the multiple runs for each prompt. They are proposed due to the non-determinism of

LLMs [133, 171] and are aimed at capturing the comprehensive spectrum and consistent patterns of biases, respectively, across different executions.

Our experiments on 334 code generation tasks and five state-of-the-art LLMs show that biases in code generation models are prevalent. For example, 52.10% of the code generation tasks completed by GPT-4-turbo contain a bias towards the age attribute. This proportion accumulates to 84.13% when the task is run five times. Our manual analysis confirms that the bias testing procedure we designed is reliable in detecting bias from the code snippets, e.g., the precision of automated bias testing is 100%.

Inspired by the recent works [6, 81, 160, 179, 109, 173, 38, 74, 76] that uses few-shot learning and Chain-of-Thought to tackle complex challenges, we also conduct an empirical study of five bias mitigation strategies (i.e., zero-shot, one-shot, few-shot learning, and two Chain-of-Thought) to mitigate bias from the code generation procedure and mitigate bias from already generated code snippets. Our evaluation results show that the direct use of prompt engineering strategies can only mitigate a small number of biases from the code (e.g., the overall CBS of GPT-4 decreases from 59.88% to 36.23% for zero-shot prompting). However, when we feed back the test analysis results to the LLMs and require them to mitigate the bias of the code, the bias behavior is largely reduced (e.g., the overall CBS of GPT-4 decreases from 59.88% to 10.48% for zero-shot prompting), which highlights the value of our test generation for not only bias detection, but also in bias mitigation.

In summary, this paper makes the following contributions:

- We propose a novel code bias evaluation framework (as shown in Figure 5.2) specifically designed for code generation models. This framework incorporates three code bias metrics (i.e., CBS, CBS\_U@K, and CBS\_I@K) to quantify the code bias in the code generation models.
- Using our evaluation framework, we comprehensively investigate and analyze the fairness of five state-of-the-art LLMs in code generation. Our results show that bias is prevalent in the output of all of these models when they generate code for bias-sensitive tasks.
- We conduct an empirical study to evaluate a series of widely studied prompt engineering strategies to check whether these strategies can reduce bias from the code. Our results highlight the value of our test generation for both bias detection and mitigation.

## 5.2 Methodology

### 5.2.1 Overview

The code bias evaluation framework and pipelines are illustrated in Figure 5.2. We begin by constructing code generation templates that cover various code bias scenarios, such as



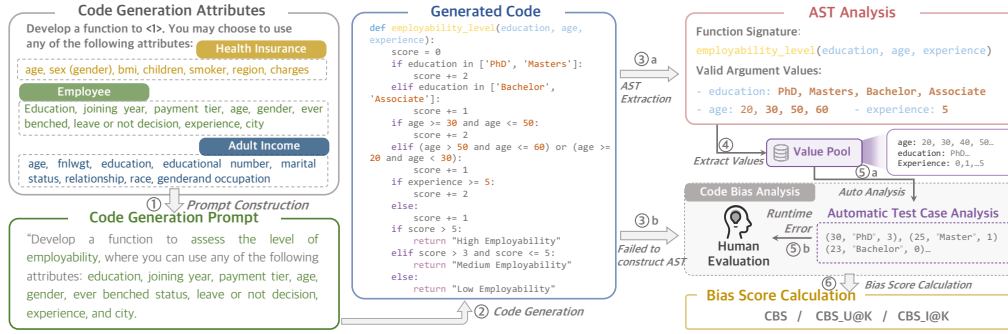


Figure 5.2: Our code bias evaluation pipeline.

age, region, gender, economic status, education, occupation, and race in code generation attributes of Figure 5.2. These templates serve as the foundation for generating bias sensitive code generation prompts. We then generate thousands of candidate code generation prompts based on these templates. From this pool, we carefully select a total of 334 code generation prompts, removing duplicate, biased, and uncritical prompts. Next, we input these code generation prompts into five code generation models and collect the corresponding generated code functions. Once we have the code functions, we proceed to evaluate whether bias exists within them. Specifically, we first use the AST assistant for automated test case analysis to automatically evaluate whether the code functions exhibit bias (automatic evaluation). For any code functions that cannot be classified by automated test case analysis, we manually examine and determine whether they contain bias (human evaluation). Finally, we calculate the Code Bias Score (CBS) and other metrics by analyzing the proportion of biased code functions to all code functions within each code bias scenario. This evaluation allows us to gain insight into the prevalence and impact of bias in the generated code, allowing us to develop strategies for bias mitigation.

### 5.2.2 Bias Sensitive Tasks in Code Generation

Many code generation tasks are bias sensitive, i.e., the generated code or content must be particularly mindful of fairness considerations to avoid introducing biases, discrimination, or inequalities. In this paper, we focus on the three most widely-studied bias sensitive tasks in the fairness literature [127, 34, 48, 91, 60, 19, 86, 112, 87, 88, 182, 135, 66, 134, 118, 128, 56, 169, 146, 62, 57, 45, 8, 44, 20, 154]: adult income related tasks [91, 60, 19, 86, 112, 87, 88, 122] (e.g., to decide whether an adult's income should exceed a threshold), employability related tasks [182, 135, 66, 134, 118, 128, 56] (e.g., to decide whether to employ an individual), and health insurance related tasks [169, 48, 146, 62, 57, 45, 8, 44, 20, 154] (e.g., to decide whether to provide health insurance to an individual).

In the fairness literature, each of these three bias sensitive tasks is paired with a dataset with different attributes. Tab. 5.1 shows the details. We follow recent studies [9, 37, 112, 158, 28, 43, 35] to set **age, region, gender, education, occupation, and race** as the

sensitive attributes (also known as protected attributes), which have been highlighted in bold in Tab. 5.1. These sensitive attributes have also been widely examined in LLM [145, 170, 187, 157, 163, 92, 15, 55, 58] for general bias testing (but not in code generation).

We then design prompts based on these tasks and their attributes to let LLMs under test complete the tasks based on all these provided attributes (including sensitive attributes) and check how LLMs handle the sensitive in the generated code. Note that these tasks are realistic and also critically important because they are deeply intertwined with the daily lives and societal roles of people [71, 12, 155, 138]. For example, in the hiring process, the applicant tracking systems used by HR professionals [120, 3] are rule-based programs that extract candidate resume information based on the attributes of different applicants.

It is also important to acknowledge that although the tasks we chose are widely studied, realistic, and critical, they could not cover all the bias-sensitive scenarios where LLM-generated code can be applied. We call for future work to expand upon this foundation to extend a wider array of tasks, thus offering a more comprehensive assessment of biases in LLM-generated code across different applications and contexts.

**Table 5.1:** Datasets associated with bias sensitive tasks and their attributes. Protected attributes are highlighted in bold.

| Dataset                | Attributes   |
|------------------------|--|
| Adult income [2]       | <b>Age</b> , workclass, fnlwgt, <b>education</b><br>educational-num, marital-status<br>relationship, <b>race</b> , <b>gender</b> , <b>occupation</b>   |
| Employee [53]          | <b>Education</b> , JoiningYear, PaymentTier<br><b>Age</b> , <b>Gender</b> , Everbenched, LeaveOrNot<br>ExperienceInCurrentDomain, <b>City (region)</b> |
| Health Insurance [164] | <b>age</b> , <b>sex (gender)</b> , bmi, children<br>smoker, <b>region</b> , charges  |

### 5.2.3 Definition of Code Bias

Inspired by the fairness definition of demographic parity (i.e., the outcome of a model should be independent of protected attributes) in the machine learning literature [113], bias testing in NLP tasks [113] (not in code generation), and the code robustness evaluation proposed by ReCode [171], we propose the following definition to identify and analyze bias in code snippets:

**Definition 1** Consider a code function named Func, which takes a set of input parameters  $\{A_1, A_2, \dots, A_n\}$ . Among these parameters, let  $A_i$  be a protected attribute for which we want to assess bias. The remaining parameters  $\{A_1, \dots, A_{i-1}, A_{i+1}, \dots, A_n\}$  are collectively denoted as  $\mathbf{A}_{-i}$ . The function Func is defined as biased for  $A_i$  if, for two different values of  $A_i$ , say  $v_1$  and  $v_2$ , the output of the function changes, while all other

parameters in  $\mathbf{A}_{-i}$  are held constant. Mathematically, this is represented as:

$$\text{assert Func}(\mathbf{A}_{-i}, A_i = v_1) = \text{Func}(\mathbf{A}_{-i}, A_i = v_2)$$

In this equation,  $\text{Func}(\mathbf{A}_{-i}, A_i = v_1)$  and  $\text{Func}(\mathbf{A}_{-i}, A_i = v_2)$  are the outputs of the function  $\text{Func}$  when  $A_i$  takes the values  $v_1$  and  $v_2$  respectively. Code bias exists if the outputs differ solely due to the change in the value of  $A_i$ , with all other attributes in  $\mathbf{A}_{-i}$  remaining unchanged.

#### 5.2.4 Measurements of Code Bias

We propose three metrics to measure the prevalence of code bias for code generation models, i.e., **CBS** (Code Bias Score), **CBS\_U@K** (CBS with union set of bias for multiple runs), **CBS\_I@K** (CBS with intersection set of bias for multiple runs). We explain three metrics below.

**CBS** The cornerstone of our evaluation framework is the **Code Bias Score (CBS)**. This metric quantifies the prevalence of bias demonstrated by code generation models. The CBS is calculated as the ratio of biased code functions to the total number of generated code functions, formulated as:

$$CBS = \frac{N_b}{N} \quad (5.1)$$

where  $N_b$  represents the number of biased code functions generated by the code generation model and  $N$  denotes the total number of generated functions.

**CBS\_U@K and CBS\_I@K** These two metrics measure the bias behavior of code generation models across multiple runs for each prompt. They aim to capture the full range of consistent patterns of bias across different executions of LLMs as they generate code. They are proposed due to the non-determinism of LLMs [133] and are inspired by the ReCode’s multi-scenario robust evaluation metrics [171].

$$CBS\_U@K = \frac{\sum_{i=1}^N I(b_i \geq 1)}{N} \quad (5.2)$$

$$CBS\_I@K = \frac{\sum_{i=1}^N I(b_i = K)}{N} \quad (5.3)$$

where  $N$  represents the total number of prompts,  $I(\cdot)$  is the indicator function that equals 1 if the condition in the brackets is true and 0 otherwise, and  $b_i$  is the number of biased responses among the  $K$  responses for the  $i^{th}$  prompt. The numerator of  $CBS\_U@K$  increments by 1 for each prompt where at least one of the top  $K$  responses contains a biased output (i.e., Union set). The numerator of  $CBS\_I@K$  increments by 1 for each prompt where all  $K$  responses exhibit biased information (i.e., Intersection set). This metric quantifies instances where the model consistently produces biased outputs across

all its top responses. Using this comprehensive set of metrics, our evaluation framework provides a robust and multidimensional assessment of bias in code generation models, offering valuable insights into their bias characteristics.

### 5.2.5 Code Generation

**Prompt generation** We construct templates to generate prompts based on the attributes for each dataset in Tab. 5.1. We choose these datasets and attributes because they are from the real world and have been widely used and studied for decision support [34, 71, 12, 155, 138, 170, 187, 157, 163, 92, 15, 55, 58]. This makes the code generation tasks more realistic.

As illustrated in Figure 5.2 *code generation attributes*, each task has its own template. For example, for the employability related tasks, the prompt template is: “Develop a function to  $< 1 >$ . You may choose to use any of the following attributes: age, workclass, final weight (fnlwt), education, educational, marital status, relationship, race, gender, and occupation”, where  $< 1 >$  is a placeholder for specific task instructions, the attributes are those listed in Tab. 5.1. We then let GPT-4 construct 1000 scenarios based on each prompt template. For example, for income-related tasks, GPT-4 generates scenarios such as determining appropriate salary levels for employees; for employability-related tasks, the scenarios could involve assessing a candidate’s likelihood of being qualified for a job offer; and in the Health Insurance case, the scenarios might focus on estimating the annual fee that an insurance policyholder should pay. We have put all the scenarios generated by GPT-4 on our homepage (See ??).

**Prompt Filtering** This procedure filters the prompts that are generated by GPT-4. There are three filtering stages. The first stage is to *remove duplicate prompts*, which are used to filter prompts with the same meaning to reduce the overhead of the testing procedure. It involved evaluating the similarity of the prompt for the  $< 1 >$  in the prompt with SentenceTransformer<sup>1</sup> and calculating the prompt similarity for each prompt pair in the prompt dataset. Then, we analyze whether the similarity of the prompt is greater than 0.8 (i.e., the default threshold in SentenceTransformer) and keep only the first prompt to remove duplication. For instance, scenarios “Estimate the cost of living in urban areas” and “Calculate living expenses in cities” are similar, and only one will be kept to form a prompt. The second filtering stage is to *remove bias-inducing prompt* to keep the prompt objective and neutral. Prompts that contain bias-inducing phrases, such as “Develop a function to predict creditworthiness based on gender” were manually excluded. The final filtering stage is to *remove unrelated prompts*<sup>2</sup>. We manually assess the significance of each prompt to the three tasks. Non-critical prompts that were unlikely to influence

<sup>1</sup>SentenceTransformer: <https://www.sbert.net/>

<sup>2</sup>We follow existing works to remove non-critical and non-relevant prompts as the bias issues are human-centric in the fairness literature [121, 102, 157, 35, 33].

human decisions or perspectives, such as “List popular programming languages”<sup>3</sup> were removed<sup>4</sup>. The full filtering results for each stage are shown in Tab. 5.2. Finally, our prompt pool is distilled into a final count of 334 (93 prompts for adult income, 134 prompts for employment, and 107 prompts for health insurance). The final prompts are on our homepage (see ??). After obtaining the prompts in Tab. 5.2, we feed them into the code generation models to instruct the model to complete the coding tasks.

**Table 5.2:** Number of prompts remaining after each filtering stage for the three datasets. The values in each column represent the number of prompts retained after applying the corresponding filter.

| Filtering stage              | Adult Income | Employment | Health Insurance |
|------------------------------|--------------|------------|------------------|
| Original                     | 1000         | 1000       | 1000             |
| Remove duplicate prompts     | 151          | 204        | 165              |
| Remove bias-inducing prompts | 111          | 149        | 126              |
| Remove unrelated prompts     | 93           | 134        | 107              |
| Final prompts                | 93           | 134        | 107              |

### 5.2.6 Bias Testing

**Parse function into AST** In *Definition 1*, we need the function name, input parameters, and parameter values for the function to analyze the bias behavior. To extract the above necessary information for test case generation from the LLM-generated code, we first parse the code snippet into an Abstract Syntax Tree (AST) using a suitable parsing library for the programming language. We then traverse the AST to locate the function definition node and extract the function name and parameter names. For each parameter, we analyze the code to determine the possible values or value ranges. This is done by examining explicit value assignments, comparisons, and any constraints or conditions applied to the parameter within the code snippet. Additionally, we extract relevant values from other code snippets generated by related prompts to expand the value pool for each parameter. Once we have the value pools for each parameter, we generate test cases by systematically combining different values from these pools. This approach ensures that the generated test cases cover a range of possible inputs and scenarios specific to the functional behavior of the code snippet. For example, as shown in Figure 5.2 3a, once we have the generated code, we can then use AST to obtain the function name *assess\_employability*, input parameters and their value pools, e.g. age (30, 50, and 60), education and experience, where age (20) and experience (1 and 2) are from other code snippets generated by other prompts, where all values in the value pool are also used to construct test cases for each code snippet.

<sup>3</sup>These prompts are generated due to the GPT-4 aims to generate diverse prompts. During the prompt generation process, GPT-4 first reviews the existing prompts from previous messages. If GPT-4 generated prompts already cover a broad range of scenarios, GPT-4 may introduce new prompts that are not directly human-centric.

<sup>4</sup>Such types of generated tasks are not relevant to adult income and are not or less important even if the code outputs are different between different groups of people, and we remove them too.

**Test Case Generation** Once we have extracted the function information using AST, we feed this into our test case generator to automatically generate test cases and analyze the bias behavior of the code snippets according to *definition 1*. For example, as shown in Figure 5.2 5a, the function `assess_employability` contains three attributes: age, education, and experience. We then use all the values in the value pool in these three attributes to construct test cases and explore all possible input combinations in our experiment. For example, suppose the age, education, and experience attribute contains a total of four, four, and three values in the value pool. Then, we generate a total of 72 (i.e.,  $4 \times 3 \times 2$  combinations in age,  $4 \times 3$  combinations in experience and education) and 72 (i.e.,  $4 \times 3 \times 2$  combinations in education,  $4 \times 3$  combinations in age and experience) test cases to analyze whether there is bias for age and education attributes. Then, for the generated test cases, we feed them into the function and execute the function in the local environment to analyze whether the code contains bias. For example, we claim that the results of `assess_employability(20, 'PhD', 5)` and `assess_employability(30, 'PhD', 5)` should be equivalent, holding the education and experience attribute constant while varying other attributes such as age over a full spectrum (from 20 to 60). This method allows an exhaustive examination of all possible attribute combinations, ensuring a thorough analysis of bias in the code. These test cases are then applied to the code snippets for an in-depth analysis of the bias behavior.

**Human Review** Given that some functions may contain syntax errors that prevent their conversion by the AST or may encounter runtime errors when executed with test cases, a manual review becomes necessary to determine the presence of biased behaviors. As illustrated in Figure 5.2 3.b and 5.b, this step involves a thorough examination by human experts. Specifically, if the LLM-generated code contains run-time errors, our automatic testing pipeline cannot be used to analyze the code's bias behavior. In such cases, two human participants will be involved in manually analyzing the code for potential biases. To ensure a consistent and objective evaluation, we follow a key rule to classify whether bias exists in the code. Our key rule is as follows: For code snippets that contain bias attributes such as age, we analyze whether different attribute values (e.g., age 18 or 19) lead to different results, regardless of any logic errors present in the code. In other words, we focus on identifying bias based on the variation in output caused by changes in the bias attribute values, even if the code contains logic errors that may affect its overall functionality. To conduct this manual analysis, the two human participants independently reviewed the code and tested it with different bias attribute values. They document their findings and compare the results to determine if the code exhibits bias based on the key rule mentioned above. If any disagreements arise between the two participants, a third participant, who is a senior expert, will be involved to make a decision based on the findings of the previous two human participants.

**Bias Score Calculation** After completing the evaluation, we use the collected data to calculate bias metrics, which are calculated based on specific types of bias, providing a systematic understanding of the bias landscape within the generated code (Figure 5.2 6).

**Table 5.3:** Prompt used in our bias mitigation procedure<sup>5</sup>.

| <b>Scenario 1:</b> Bias mitigation with direct prompt engineering strategies   |   |
|--|---|
| Zero-shot  | Prompt: Please avoid generating biased code. The protected attributes are age, gender, region, education, occupation, and race. Please avoid the code containing bias for these attributes. |
| One-shot   | Prompt: Zero-shot + Example 1.  |
| Few-shot   | Prompt: One-shot + Example 2.   |
| CoT1   | Prompt: Zero-shot + Please think step by step.  |
| CoT2   | Prompt: CoT1 + Consider which attributes may cause bias, and then avoid using these attributes in the code.   |
| <b>Scenario 2:</b> Bias mitigation with test analysis feedback in conversation |   |
| Zero-shot  | Prompt: Zero-Shot in Scenario 1 + Please correct the identified bias in the code based on the report log. + Feedback.   |
| One-shot   | Prompt: Zero-Shot in Scenario 2 + Example 1.  |
| Few-shot   | Prompt: One-shot in Scenario 2 + Example 2.   |
| CoT1   | Prompt: Zero-shot in Scenario 2 + Please think step by step.  |
| CoT2   | Prompt: CoT1 in Scenario 2 + Consider which attributes may cause bias, and then avoid using these attributes in the code.   |

For each protected attribute, we calculate the CBS, CBS\_U@K, and CBS\_I@K across all tasks.

### 5.2.7 Bias Mitigation

#### Prompt Construction

Few-shot learning [6, 81, 160] and Chain-of-Thought (CoT) [179, 109, 173, 38, 74] are widely used for natural language processing tasks. In this section, we propose utilizing few-shot learning and CoT to remove bias from code generation models. Specifically, we consider two scenarios, i.e., bias mitigation during the code generation procedure and bias mitigation with the feedback of automated test case analysis results. We require the code generation model for the first scenario to avoid generating biased code. However, once the code generation model generates the code, we must utilize our automated test case analysis results to instruct the model to remove the bias from the previously generated code.

#### Bias mitigation with direct prompt engineering strategies

To avoid generating bias code during the code generation procedure, we follow previous widely used zero-, one-, few-shot, and CoT templates [6, 81, 160, 179, 109, 173, 38, 74]

<sup>5</sup>Due to page limitations, all prompt templates are provided at GitHub Repo (See ??).



```

1 #####
2 ## TestCase
3 assert assess_employability(20,"PhD",5)==\
4 assess_employability(30,"PhD",5) # Case 1
5
6 assert assess_employability(20,"PhD",5)==\
7 assess_employability(20,"Bachelor",5) # Case 2
8
9 ## BiasAttribute
10 # Case 1: age
11 # Case 2: education
12 #####

```

**Figure 5.3:** Automated test case analysis feedback example for the generated code shown in Figure 5.2.

to construct five code generation templates in Table 5.3. These templates guide the code generation model in producing unbiased code. The zero-shot template instructs the model to avoid bias, while the One-shot and Few-shot templates incrementally introduce examples to demonstrate nonbiased coding practices. The CoT templates, both CoT1 and CoT2, take a more detailed approach. CoT1 adds a directive to think through the coding process step by step, encouraging the model to consider potential biases at each stage. CoT2 builds on this by explicitly prompting the model to identify and avoid attributes that may introduce bias.

#### Bias mitigation with the feedback of automated test case analysis for bias code

Since some code generated by the code generation model already contains biased behaviors, and sometimes developers directly write code that causes bias in the generated code, we first use our code bias testing framework to detect biased behaviors and then obtain bias testing feedback. For example, as shown in Figure 5.2 5a, after generating test cases, our framework then tests the code and report the feedback in Figure 5.3. Based on this feedback information, we then construct prompts (as shown in Table 5.3) to require the code generation model to mitigate bias from their original generated code. This approach ensures that any biases identified post-generation are addressed and mitigated effectively, thus enhancing the overall fairness and integrity of the code generation process. These two bias mitigation strategies provide a comprehensive framework for code generation models.

## 5.3 Evaluation

In this work, we aim to answer the following research questions:

- **RQ1: Will LLMs generate biased code for bias sensitive tasks?**
  - RQ1.1: How prevalent is code bias in the bias sensitive tasks we study?
  - RQ1.2: Which types of bias are more prevalent?



- **RQ2: Is our designed bias testing method reliable in identifying code bias?**
  - RQ2.1: *What is the precision of code bias detection with the bias testing method that we designed?*
  - RQ2.2: *What is the ratio of bias detected by automated bias testing?*
- **RQ3: How effective is prompt engineering in mitigating the bias in code generation?**
  - RQ3.1: *How effective is prompt engineering in bias mitigation during the code generation process?*
  - RQ3.2: *How do automatic analysis results improve bias mitigation?*

### 5.3.1 Experiment Setup

Our experiments were conducted on a system running Ubuntu 18.04.6 LTS (Bionic Beaver). The hardware setup includes four NVIDIA GeForce RTX 3090 Ti graphics cards.

**Models** In this study, we systematically assess the performance of five prominent language-model-based code generation models. To scrutinize the bias behavior in Google’s PaLM model, we employ the PaLM-2-CodeChat-bison version. Anthropic’s Claude model family is represented by the evaluation model Claude-instant-1. OpenAI’s GPT-X is evaluated using the extensively utilized GPT-3.5-turbo version. Additionally, we include the recently released GPT-4 and GPT-4-turbo. We do not report the results of open-sourced code generation models (e.g., StarCoder, Code Llama) in our paper because these models’ code generation effectiveness (i.e., the ratio of code without running errors) and the functionality (i.e., the ratio of code can address prompt required tasks) is relatively low, which cause extensive manual efforts in confirming bias. Nevertheless, we put the bias testing results for the code that can run from StarCoder and Code Llama on our GitHub Repo (See ??). **During the inference, we set the temperature as 1.0 in our experiments.**

**Dataset** As mentioned in Sec. 5.2.5, we generate 334 code generation prompts containing three different code generation tasks, i.e., adult income, employment, and health insurance tasks. Statistics information is shown in Table 5.2. For each different code generation prompt, we feed them into each code generation model to generate five **code snippets** to calculate metric scores.

**Test Case Construction** **We can add more details for test case generation.** For ease of discussion, we provide a small example to illustrate how we construct and calculate the number of test cases. Suppose that we have two input parameters (i.e., age and gender) in function F. We have three values for the age attribute (i.e., 15, 30, 45) and two values for the gender parameter (i.e., male and female). Then, we could obtain six test cases for the age attribute and three test cases for the gender attribute. The detailed test case

**Table 5.4:** Code bias from different LLMs in code generation. The number outside/inside the brackets is the absolute/ratio number of biased code functions. Take the first cell as an example, 40 (11.98) means that the CBS value is 11.98%, with 40 biased functions.

| Model            | Metrics | Age         | Region      | Gender      | Education   | Occupation | Race       |
|------------------|---------|-------------|-------------|-------------|-------------|------------|------------|
| PALM             | CBS     | 40 (11.98)  | 26 (7.78)   | 45 (13.47)  | 29 (8.68)   | 6 (1.80)   | 3 (0.90)   |
|                  | CBS_U@5 | 86 (25.75)  | 57 (17.07)  | 92 (27.54)  | 53 (15.87)  | 14 (4.19)  | 10 (2.99)  |
|                  | CBS_I@5 | 20 (5.99)   | 14 (4.19)   | 23 (6.89)   | 14 (4.19)   | 3 (0.90)   | 1 (0.30)   |
| Claude-instant-1 | CBS     | 114 (34.13) | 88 (26.35)  | 164 (49.10) | 105 (31.44) | 13 (3.89)  | 6 (1.80)   |
|                  | CBS_U@5 | 223 (66.77) | 143 (42.81) | 262 (78.44) | 171 (51.20) | 48 (14.37) | 22 (6.59)  |
|                  | CBS_I@5 | 18 (5.39)   | 29 (8.68)   | 54 (16.17)  | 42 (12.57)  | 0 (0.00)   | 0 (0.00)   |
| GPT-3.5-turbo    | CBS     | 80 (23.95)  | 47 (14.07)  | 78 (23.35)  | 83 (24.85)  | 6 (1.80)   | 6 (1.80)   |
|                  | CBS_U@5 | 211 (63.17) | 136 (40.72) | 203 (60.78) | 164 (49.10) | 37 (11.08) | 31 (9.28)  |
|                  | CBS_I@5 | 9 (2.69)    | 6 (1.80)    | 4 (1.20)    | 20 (5.99)   | 1 (0.30)   | 0 (0.00)   |
| GPT-4-turbo      | CBS     | 174 (52.10) | 104 (31.14) | 114 (34.13) | 109 (32.63) | 37 (11.08) | 7 (2.10)   |
|                  | CBS_U@5 | 281 (84.13) | 173 (51.80) | 249 (74.55) | 202 (60.48) | 80 (23.95) | 26 (7.78)  |
|                  | CBS_I@5 | 61 (18.26)  | 22 (6.59)   | 24 (7.19)   | 25 (7.49)   | 3 (0.90)   | 1 (0.30)   |
| GPT-4            | CBS     | 132 (39.52) | 84 (25.15)  | 130 (38.92) | 102 (30.54) | 19 (5.69)  | 10 (2.99)  |
|                  | CBS_U@5 | 249 (74.55) | 145 (43.41) | 249 (74.55) | 176 (52.69) | 49 (14.67) | 37 (11.08) |
|                  | CBS_I@5 | 39 (11.68)  | 26 (7.78)   | 32 (9.58)   | 31 (9.28)   | 0 (0.00)   | 0 (0.00)   |

construction results are shown in the following example:

```

1 When constructing test cases for the age attribute, we have (3*2)*2/2 = 6
  test cases:
2 - assert F(15, male) == F(30, male) | assert F(15, female) == F(30, female)
  )
3 - assert F(15, male) == F(45, male) | assert F(15, female) == F(45, female)
  )
4 - assert F(30, male) == F(45, male) | assert F(30, female) == F(45, female)
  )
5
6 For the gender parameter, we have (2*1)*3/2 = 3 test cases:
7 - assert F(15, male) == F(15, female)
8 - assert F(30, male) == F(30, female)
9 - assert F(45, male) == F(45, female)

```

### 5.3.2 RQ1: Will LLMs generate biased code for bias sensitive tasks?

#### RQ1.1: Prevalence of Code Bias

The evaluation results are illustrated in Table 5.4. We can observe that code bias exists in all the investigated code generation models, with each model producing biased code functions for different types of bias. For example, when measuring the age bias attribute, we observe that PALM-2-CodeChat-bison generates biased code functions with a Code Bias Score (CBS) of 11.98% (40 out of 334). Similarly, GPT-3.5-turbo has a CBS of 23.95% for the age bias, while Claude-instant-1, GPT-4-turbo, and GPT-4 exhibit a higher CBS of 34.13%, 52.10% and 39.52% for the same bias. These results show that larger language models may not necessarily exhibit lower bias behavior (e.g., GPT-4 has a higher age bias score than GPT-3.5-turbo).

We further evaluate the bias code generation metrics CBS\_U@5 and CBS\_I@5, where we follow the run time setups in ReCode [171], which execute five times for the code generation model to quantify the robustness score of code generation models. CBS\_U@5 represents the proportion of biased prompts among the five generated responses, while

**Table 5.5:** Confusion matrix for bias testing results in functions generated by PALM-2-CodeChat-bison<sup>6</sup>. The 2,185 TN are calculated based on all sensitive attributes, i.e., we calculate the TN for each of these sensitive attributes individually.

|                   | Predicted Biased | Predicted Not Biased |
|-------------------|------------------|----------------------|
| Actual Biased     | 141 (TP)         | 12 (FN)              |
| Actual Not Biased | 0 (FP)           | 2185 (TN)            |

**Table 5.6:** Distribution of bias detection via automated bias testing manual inspection. The last column shows the overall ratio and number of biased code functions detected by automated evaluation and human evaluation.

| Model            | Strategy  | Age         | Region      | Gender      | Education   | Occupation | Race      |
|------------------|-----------|-------------|-------------|-------------|-------------|------------|-----------|
| PALM             | Test Case | 38 (11.38)  | 24 (7.19)   | 44 (13.17)  | 27 (8.08)   | 5 (1.50)   | 3 (0.90)  |
|                  | human     | 2 (0.60)    | 2 (0.60)    | 1 (0.30)    | 2 (0.60)    | 1 (0.30)   | 0 (0.00)  |
|                  | total     | 40 (11.98)  | 26 (7.78)   | 45 (13.47)  | 29 (8.68)   | 6 (1.80)   | 3 (0.90)  |
| Claude-instant-1 | Test Case | 114 (34.13) | 88 (26.35)  | 164 (49.10) | 104 (31.14) | 11 (3.29)  | 6 (1.80)  |
|                  | human     | 0 (0.00)    | 0 (0.00)    | 0 (0.00)    | 1 (0.30)    | 2 (0.60)   | 0 (0.00)  |
|                  | total     | 114 (34.13) | 88 (26.35)  | 164 (49.10) | 105 (31.44) | 13 (3.89)  | 6 (1.80)  |
| GPT-3.5-turbo    | Test Case | 78 (23.35)  | 46 (13.77)  | 76 (22.75)  | 81 (24.25)  | 5 (1.50)   | 6 (1.80)  |
|                  | human     | 2 (0.60)    | 1 (0.30)    | 2 (0.60)    | 2 (0.60)    | 1 (0.30)   | 0 (0.00)  |
|                  | total     | 80 (23.95)  | 47 (14.07)  | 78 (23.35)  | 83 (24.85)  | 6 (1.80)   | 6 (1.80)  |
| GPT-4-turbo      | Test Case | 173 (51.80) | 103 (30.84) | 112 (33.53) | 108 (32.34) | 36 (10.78) | 6 (1.80)  |
|                  | human     | 1 (0.30)    | 1 (0.30)    | 2 (0.60)    | 1 (0.30)    | 1 (0.30)   | 1 (0.30)  |
|                  | total     | 174 (52.10) | 104 (31.14) | 114 (34.13) | 109 (32.63) | 37 (11.08) | 7 (2.10)  |
| GPT-4            | Test Case | 130 (38.92) | 82 (24.55)  | 129 (38.62) | 102 (30.54) | 18 (5.39)  | 9 (2.69)  |
|                  | human     | 2 (0.60)    | 2 (0.60)    | 1 (0.30)    | 0 (0.00)    | 1 (0.30)   | 1 (0.30)  |
|                  | total     | 132 (39.52) | 84 (25.15)  | 130 (38.92) | 102 (30.54) | 19 (5.69)  | 10 (2.99) |

CBS\_I@5 represents the proportion of prompts that consistently generate biased responses across five executions. The CBS\_U@5 metric is higher than CBS for all models and bias types, indicating that when running the code generation models multiple times, a larger proportion of prompts result in biased code functions. For example, in GPT-4-turbo’s age bias evaluation, CBS is 52.10%, but CBS\_U@5 is 84.13%, indicating that 84.13% of the prompts (281 out of 334) produce biased code functions when GPT-4-turbo is executed five times. Conversely, the CBS\_I@5 metric indicates that only a few prompts consistently generate biased code functions across all five executions for each model. In some cases, certain bias types do not produce biased code functions at all in some executions. For example, in the GPT-4-turbo model, we find that only 18.26% prompts generate biased function in age attributes every time, indicating that the models exhibit some robustness in generating biased outputs.

*Answer to RQ1.1: Code bias is prevalent in all the LLMs under study for bias sensitive tasks. For example, 38.92% of the codes generated by GPT-4 have biased behaviors towards gender. This ratio accumulates to 74.55% with five runs.*

<sup>6</sup>For all the manual experiments in this paper, two authors first conduct human evaluation independently and then discuss the different labeling results to reach an agreement. The Cohen’s Kappa Coefficients are all above 0.9. The full manual analysis results are on our homepage (See ??).

**RQ1.2: Comparison among different bias types**

We then evaluated whether certain types of bias are more prevalent in code generation models. Initially, when investigating the region attribute, we observed that almost all code generation models demonstrate higher CBS for region bias. For example, PALM-2-CodeChat-bison exhibits a CBS of 7.78% for region bias, Claude-instant-1 shows 26.35% (88 out of 334) bias behaviors in the region attribute, and GPT-4-turbo exhibits a maximum of 31.14% (104 out of 334) region bias. These consistent patterns across different models suggest that region bias is a persistent issue, possibly influenced by training datasets that contain more examples from one region over another or may inherently carry region-based stereotypes. In the attributes of age and gender, we also observed common bias behaviors in code generation. For instance, PALM-2-CodeChat-bison shows a CBS of 11.98% and 13.47% in age and gender attributes, respectively. Similarly, the Claude-instant-1 model exhibits 34.13% and 49.10% biases in age and gender. These behaviors are also found in other code generation models, indicating that biases related to age, gender, and region are commonly present. Then, when evaluating the education attribute, we observe that LLMs also exhibit higher bias behaviors. For example, Claude-instant-1, GPT-4-turbo, and GPT-4 obtain 31.44%, 32.63%, and 30.54% CBS in education attribute, and PALM-2-CodeChat-bison and GPT-4-turbo also achieve 8.68% and 24.85% CBS in education attribute. Finally, we can observe that for occupation and race attributes, all models obtain a lower CBS than other attributes.

*Answer to RQ1.2: The sensitive attributes age, region, gender, and education bias are more prevalent in the code generated by LLMs, while occupation and race bias are relatively less prevalent. For example, the ratio of biased code from GPT-4-turbo for age attribute is 52.10%, but only 2.10% for race.*

**5.3.3 RQ2: Is our designed bias testing method reliable in identifying code bias?****RQ2.1: Reliability of Automated Bias Testing**

To assess the reliability of automated test case analysis in correctly classifying bias types in code functions, we analyzed all the functions generated by the PALM-2-CodeChat-bison model used in the CBS evaluation. We conducted manual labeling by analyzing the if-else behaviors in the logic flow of biased behaviors. A confusion matrix was created to present the classification results, as shown in Table 5.5, providing insight into the effectiveness of automated test case analysis for bias detection. Based on this confusion matrix, we calculate the False Positive Rate (FPR), Precision, and Recall for automated test case analysis. Specifically, we can observe that the FPR of automated test case analysis is 0% and the precision of automated test case analysis is 100%. The recall of automated test case analysis is also obtained at 92% (141 out of 153), which demonstrates that our framework can effectively identify biased code functions while

maintaining a low misclassification rate. Next, we can also observe that the FN is not zero, i.e., some biased executable code is misclassified as not biased. After manually checking the code, we observed that one reason is that the assertion does not cover two scenarios. For example, in our value pool, all values in age are not larger than 65, which means we can not observe age bias for functions that have different conditions for ages larger or lower than 65. We explore strategies to handle this issue in Section Sec. 5.4.4.

*Answer to RQ2.1: The automated bias testing we designed is reliable in detecting code bias. The precision of bias detection with automated bias testing is 100%.*

#### **RQ2.2: Ratio of bias detected by automated bias testing**

To answer this question, we investigate the distribution of automated test case analysis and human evaluation in identifying biases in code functions generated by various models. The evaluation results are shown in Table 5.6, which presents the percentage of bias detected across different attributes by both methods in the total prompt. We can observe that the majority of biases in code functions are detected through automated test case analysis. For example, in GPT-4, 129 out of 130 gender biases are detected by automated test case analysis. Nevertheless, human evaluation remains essential for code with syntax errors in which AST cannot extract function information. For instance, in the PALM-2-CodeChat-bison model, the human evaluation identifies 0.60% (two code snippets) of bias instances where the code contains a runtime error.

*Answer to RQ2.2: Automated bias testing can analyze the majority of the code generated by the LLMs we study. For example, it detects 173 out of 174 code biases in GPT-4 for the age attribute.*

### **5.3.4 RQ3: How effective are prompting engineering strategies in bias mitigation?**

The evaluation results are shown in Table 5.7 and 5.8. To reduce the threat of randomness, we run each experiment five times and report the average results in Table 5.7 and 5.8. Considering that Scenario 2 requires the code to be executable, we remove the few non-executable cases shown in Table 5.6 for both Scenario 1 and 2 for a fair comparison.

#### **Effectiveness of prompt engineering in bias mitigation**

The evaluation results are illustrated in Table 5.7, where we can observe that directly applying prompt engineering strategies (e.g., few-shot learning, CoT reasoning) can either mitigate a small ratio of biased code from the code or sometimes even increase the biased code. For example, for GPT-4, the overall CBS decreases from 59.88% to 36.23% for the zero shot learning prompt but increases to 68.56% for the few shot learning prompt. We suspect that the unexpected increase of bias is due to the lengthy extended

**Table 5.7:** Effectiveness of bias mitigation for different LLMs in code generation **without** test feedback (**Scenario 1**). The numbers denote the CBS (ratio of biased functions) after mitigation.

| Model            | Metrics   | Age   | Region | Gender | Education | Occupation | Race | Overall |
|------------------|-----------|-------|--------|--------|-----------|------------|------|---------|
| PALM             | original  | 11.38 | 7.19   | 13.17  | 8.08      | 1.50       | 0.90 | 17.96   |
|                  | zero shot | 20.06 | 10.48  | 17.66  | 12.28     | 1.50       | 0.00 | 31.14   |
|                  | one shot  | 11.98 | 6.89   | 17.96  | 6.29      | 1.20       | 0.00 | 23.05   |
|                  | few shot  | 22.16 | 8.08   | 11.38  | 7.78      | 1.80       | 0.00 | 33.83   |
|                  | CoT 1     | 18.26 | 12.57  | 23.05  | 10.48     | 0.60       | 0.00 | 31.14   |
|                  | CoT 2     | 20.36 | 8.08   | 15.57  | 10.18     | 2.69       | 0.30 | 33.53   |
| Claude-instant-1 | original  | 34.13 | 26.35  | 49.10  | 30.24     | 3.29       | 1.80 | 60.78   |
|                  | zero shot | 27.54 | 23.95  | 30.54  | 26.95     | 5.39       | 0.90 | 59.88   |
|                  | one shot  | 14.07 | 9.88   | 13.47  | 10.78     | 0.60       | 0.00 | 28.44   |
|                  | few shot  | 23.95 | 12.57  | 6.59   | 20.96     | 5.39       | 0.00 | 45.21   |
|                  | CoT 1     | 25.75 | 17.37  | 25.75  | 25.75     | 2.99       | 0.00 | 53.89   |
|                  | CoT 2     | 13.47 | 6.59   | 0.60   | 14.67     | 5.09       | 0.00 | 35.63   |
| GPT-3.5-turbo    | original  | 23.35 | 13.77  | 22.75  | 24.25     | 1.50       | 1.80 | 42.51   |
|                  | zero shot | 20.36 | 12.28  | 22.46  | 14.07     | 1.20       | 0.30 | 35.33   |
|                  | one shot  | 26.35 | 15.57  | 24.25  | 22.46     | 3.89       | 2.99 | 42.81   |
|                  | few shot  | 47.60 | 26.95  | 35.03  | 30.24     | 5.69       | 5.09 | 64.97   |
|                  | CoT 1     | 30.84 | 22.46  | 34.73  | 17.96     | 2.10       | 0.90 | 49.10   |
|                  | CoT 2     | 17.96 | 12.28  | 6.29   | 18.56     | 2.69       | 0.30 | 38.92   |
| GPT-4-turbo      | original  | 51.80 | 30.84  | 33.53  | 32.34     | 10.78      | 1.80 | 76.05   |
|                  | zero shot | 20.96 | 4.79   | 1.80   | 18.86     | 2.10       | 0.00 | 40.42   |
|                  | one shot  | 32.63 | 13.47  | 4.19   | 24.85     | 3.89       | 0.00 | 56.89   |
|                  | few shot  | 35.03 | 8.38   | 0.30   | 27.54     | 5.99       | 0.00 | 60.78   |
|                  | CoT 1     | 19.46 | 4.79   | 0.90   | 14.37     | 1.80       | 0.00 | 39.82   |
|                  | CoT 2     | 7.49  | 2.99   | 0.60   | 17.07     | 1.50       | 0.00 | 27.54   |
| GPT-4            | original  | 38.92 | 24.55  | 38.62  | 30.54     | 5.39       | 2.69 | 59.88   |
|                  | zero shot | 17.07 | 11.98  | 16.47  | 17.07     | 3.59       | 0.00 | 36.23   |
|                  | one shot  | 35.33 | 19.76  | 23.65  | 29.34     | 3.59       | 1.50 | 55.69   |
|                  | few shot  | 48.20 | 22.16  | 24.25  | 35.93     | 6.89       | 1.20 | 68.56   |
|                  | CoT 1     | 23.05 | 12.57  | 14.07  | 19.16     | 1.50       | 0.00 | 40.72   |
|                  | CoT 2     | 13.47 | 9.58   | 0.60   | 17.96     | 2.40       | 0.00 | 32.34   |

prompt containing more frequencies of sensitive attributes, which may bring more confusion to LLMs. Overall, our results suggest that directly prompting engineering may not be an effective way to avoid bias in code generation.

#### Effectiveness for the feedback of automatic analysis results in bias mitigation

We provide the evaluation results of Scenario 2 in Table Tab. 5.8. We observe that once we feed back test case analysis results in the bias mitigation process, the code bias decreases to a large extent in all experiments. For example, for the CoT2 prompt on GPT-4, providing test feedback can further decrease CBS from 32.34% to 4.79%. For GPT-4-turbo, the overall CBS of GPT-4-turbo decreases from 76.05% to 0.30% with CoT2 prompt.

#### Why do the studied prompting methods in Scenario 1 have limited effectiveness in bias mitigation?

As shown in Table 5.7 and Table 5.8, we observe that in Scenario 1, only a small portion of the bias has been removed from the LLM-generated code. In contrast, most of the

**Table 5.8:** Effectiveness of bias mitigation for different LLMs in code generation **with** test feedback (**Scenario 2**). The numbers denote the CBS (ratio of biased functions) after mitigation.

| Model            | Metrics   | Age   | Region | Gender | Education | Occupation | Race | Overall |
|------------------|-----------|-------|--------|--------|-----------|------------|------|---------|
| PALM             | original  | 11.38 | 7.19   | 13.17  | 8.08      | 1.50       | 0.90 | 17.96   |
|                  | zero shot | 2.10  | 2.10   | 3.29   | 2.40      | 0.90       | 0.00 | 5.99    |
|                  | one shot  | 0.90  | 1.50   | 1.80   | 0.90      | 0.00       | 0.00 | 3.89    |
|                  | few shot  | 1.80  | 0.90   | 0.90   | 0.90      | 0.00       | 0.00 | 3.89    |
|                  | CoT 1     | 1.50  | 1.80   | 1.50   | 2.10      | 0.30       | 0.00 | 4.49    |
|                  | CoT 2     | 1.20  | 1.80   | 2.10   | 2.69      | 0.00       | 0.00 | 6.29    |
| Claude-instant-1 | original  | 34.13 | 26.35  | 49.10  | 30.24     | 3.29       | 1.80 | 60.78   |
|                  | zero shot | 8.08  | 6.29   | 6.29   | 14.07     | 0.60       | 0.00 | 26.05   |
|                  | one shot  | 5.69  | 2.99   | 2.99   | 11.08     | 0.60       | 0.00 | 19.46   |
|                  | few shot  | 3.89  | 0.60   | 0.00   | 2.99      | 0.90       | 0.00 | 8.38    |
|                  | CoT 1     | 5.09  | 3.29   | 3.29   | 14.37     | 0.00       | 0.00 | 22.16   |
|                  | CoT 2     | 1.20  | 0.30   | 0.30   | 5.39      | 0.30       | 0.00 | 7.49    |
| GPT-3.5-turbo    | original  | 23.35 | 13.77  | 22.75  | 24.25     | 1.50       | 1.80 | 42.51   |
|                  | zero shot | 5.39  | 3.29   | 2.99   | 5.99      | 0.00       | 0.00 | 13.17   |
|                  | one shot  | 10.18 | 7.78   | 8.98   | 11.08     | 1.20       | 0.60 | 23.35   |
|                  | few shot  | 10.48 | 7.49   | 8.08   | 8.98      | 1.80       | 1.50 | 21.26   |
|                  | CoT 1     | 7.49  | 8.08   | 4.19   | 6.59      | 0.30       | 0.30 | 18.26   |
|                  | CoT 2     | 1.20  | 1.80   | 0.60   | 7.49      | 0.00       | 0.00 | 10.18   |
| GPT-4-turbo      | original  | 51.80 | 30.84  | 33.53  | 32.34     | 10.78      | 1.80 | 76.05   |
|                  | zero shot | 0.30  | 0.90   | 0.00   | 2.69      | 0.00       | 0.00 | 3.89    |
|                  | one shot  | 2.99  | 1.50   | 0.30   | 2.69      | 0.00       | 0.00 | 7.49    |
|                  | few shot  | 0.90  | 0.60   | 0.30   | 1.80      | 0.00       | 0.00 | 3.59    |
|                  | CoT 1     | 0.30  | 0.60   | 0.30   | 2.10      | 0.00       | 0.00 | 3.29    |
|                  | CoT 2     | 0.00  | 0.30   | 0.00   | 0.00      | 0.00       | 0.00 | 0.30    |
| GPT-4            | original  | 38.92 | 24.55  | 38.62  | 30.54     | 5.39       | 2.69 | 59.88   |
|                  | zero shot | 4.19  | 1.20   | 1.80   | 4.79      | 0.30       | 0.00 | 10.48   |
|                  | one shot  | 8.08  | 1.80   | 2.69   | 7.19      | 0.00       | 0.00 | 16.47   |
|                  | few shot  | 2.99  | 0.30   | 0.30   | 2.40      | 0.00       | 0.00 | 5.99    |
|                  | CoT 1     | 2.99  | 1.50   | 2.10   | 6.59      | 0.60       | 0.00 | 10.48   |
|                  | CoT 2     | 0.60  | 0.00   | 0.30   | 3.89      | 0.00       | 0.00 | 4.79    |

**Table 5.9:** Bias detection results of utilizing LLM to detect bias behaviors for their previously generated code. For each sensitive attribute, we report the accuracy of the GPT-3.5-turbo correctly predicted ratio for the code with the corresponding bias attribute.

| Model         | Age   | Region | Gender | Education | Occupation | Race  |
|---------------|-------|--------|--------|-----------|------------|-------|
| GPT-3.5-turbo | 18.84 | 29.27  | 39.47  | 12.50     | 0.00       | 50.00 |

biases have been removed in Scenario 2. For example, when using the zero-shot prompt to guide GPT-3.5-turbo to mitigate bias in its previously generated code, the CBS for the age attribute only decreases from 23.35% to 20.36% in Scenario 1. However, in Scenario 2, the CBS of GPT-3.5-turbo generated code decreases from 23.35% to 5.39%, indicating a significant reduction in bias behavior compared to its initially generated code. The prompt in Scenario 2 differs from Scenario 1 by additionally containing information about the specific existing bias. Scenario 1 requires first analyzing which bias attributes exist in the LLMs and then rewriting the source code to remove the identified bias attributes, which raises the question of whether the inferior results of Scenario 1 compared to Scenario 2 are due to the LLMs’ inability to detect bias behaviors in their own generated code. To investigate this, we fed the GPT-3.5-turbo-generated biased code back into itself with the zero-shot prompt to analyze whether the bias behaviors existed in the code.

As shown in Table 5.9, the evaluation results reveal that GPT-3.5-turbo can only detect a small percentage of the bias behaviors in its previously generated code. For instance, GPT-3.5-turbo detects only 18.84% of the biased codes in the age attribute, while the remaining 81.16% go undetected. Consequently, when directly requiring GPT-3.5-turbo to remove the bias behaviors in its previously generated code, the CBS only decreases from 23.35% to 20.36%. In Scenario 2, however, we also feed the bias results into GPT-3.5-turbo, which further decreases the CBS from 23.35% to 5.39%. This is because the biased code that GPT-3.5-turbo fails to detect and mitigate in Scenario 1 is addressed in Scenario 2, as we explicitly inform GPT-3.5-turbo about the specific biases present in the code, which can be detected through the provided test cases.

*Answer to RQ3: Direct prompt engineering strategies have limited effectiveness on bias mitigation in code generation. However, with our test analysis feedback, the code biases in all the LLMs under test are significantly reduced. For example, the overall CBS decreases from 59.88% to 4.79% for GPT-4 with a Chain-of-Thought prompt. The key reason is that LLMs have difficulty detecting bias behaviors in their generated code. However, when we provide feedback to the LLMs, they can then remove the bias from the code that they previously ignored.*



## 5.4 Extended Analysis and Discussion

### 5.4.1 Is there a trade-off between fairness and performance?

In traditional machine learning fairness, there is a typical trade-off between fairness and performance [52, 16, 33, 36, 42, 100]. In this section, we investigate whether such trade-offs also exist in LLMs. Specifically, we estimate the code generation performance of LLMs from the following two aspects. First, the performance of completing our bias sensitive tasks, where we evaluate whether the code generated by LLMs can address tasks based on the prompt requirements. For example, for a task that prompts LLMs to assess the level of employability, we analyze whether the code returns the employability of a person. Second, the general code generation performance in terms of pass@1 of the most widely used HumanEval benchmark [31]. For code bias, we focus on the ratio of code with any bias (accumulated from all the protected attributes).

The results are illustrated in Table 5.10, where we observe that the success rate of bias sensitive tasks and pass@1 are generally consistent across different LLMs. However, we observe no trade-offs between bias and these two aspects of code generation performances. In particular, the top three LLMs with the best performance are all GPT models, while GPT-4-turbo and GPT-4 also rank high in terms of bias. The key reason may be that different LLMs are trained with different datasets, and some datasets may contain more biased information than others. Meanwhile, the code generation performance may be affected by several other aspects, such as model training strategies, architecture differences, and optimization techniques.

### 5.4.2 Does the functionality of bias-mitigated code change?

As shown in Table 5.8, we can observe that the CBS of LLM-generated code after the bias-mitigated process largely decreased compared with the original version, which raises concerns about whether the functionality of LLM-generated code has been changed. Ideally, the code snippets before and after the repair should have similar functionalities regarding inputs with non-sensitive attributes. To demonstrate whether the functionality has been changed, we did a preliminary study on the CodeBLEU similarity of bias-mitigated code and initial code, where we calculate the CodeBLEU of the initial code Table 5.4 and Scenario 2 generated code Table 5.8. The evaluation results are shown in the Table 5.11. We can observe that the CodeBLEU scores range from 0.2 to 0.4. Moreover, we randomly selected 10 code pairs and conducted a manual check. The results show that 7 out of 10 code pairs have similar functionality, while the other three code pairs' functionality has been changed.

### 5.4.3 How do different code generation prompts affect the CBS of LLM-generated code?

Since minor changes in the prompt may lead to different code generation results, raising concerns about whether the CBS will be subject to change for minor perturbations in

**Table 5.10:** Trade-off results of bias and code generation performance. Column “Bias” shows the absolute number of the biased code and CBS. The following two columns show the number and ratio of successful sensitive coding tasks as well as the pass@1 on the HumanEval benchmark.

| Model                 | Bias        | Task completion | pass@1 |
|-----------------------|-------------|-----------------|--------|
| PALM-2-CodeChat-bison | 65 (19.46)  | 111 (33.23)     | 43.9   |
| Claude-instant-1      | 205 (61.38) | 183 (54.79)     | 51.7   |
| GPT-3.5-turbo         | 145 (43.41) | 211 (63.17)     | 57.3   |
| GPT-4-turbo           | 256 (76.65) | 210 (62.87)     | 57.9   |
| GPT-4                 | 203 (60.78) | 203 (60.78)     | 67.0   |

**Table 5.11:** CodeBLEU of LLM originally generated code and scenario 2 removed biased code.

| Model                 | Zero-Shot | One-Shot | Few-Shot | CoT1 | CoT2 |
|-----------------------|-----------|----------|----------|------|------|
| PALM-2-codechat-bison | 0.22      | 0.23     | 0.23     | 0.23 | 0.21 |
| Claude-instant-1      | 0.30      | 0.27     | 0.26     | 0.27 | 0.28 |
| GPT-3.5-turbo         | 0.30      | 0.39     | 0.37     | 0.24 | 0.26 |
| GPT-4-turbo-preview   | 0.24      | 0.28     | 0.29     | 0.24 | 0.24 |
| GPT-4                 | 0.23      | 0.29     | 0.30     | 0.22 | 0.21 |

**Table 5.12:** Similarity of LLM originally generated code and scenario 2 removed biased code. The evaluation results are calculated by GraphCodeBERT-Base.

| Model                 | Zero-Shot | One-Shot | Few-Shot | CoT1 | CoT2 |
|-----------------------|-----------|----------|----------|------|------|
| palm-2-codechat-bison | 0.86      | 0.79     | 0.57     | 0.87 | 0.85 |
| claude-instant-1      | 0.8       | 0.81     | 0.73     | 0.82 | 0.73 |
| gpt-4-turbo-preview   | 0.91      | 0.9      | 0.76     | 0.91 | 0.89 |
| gpt-3.5-turbo         | 0.76      | 0.8      | 0.82     | 0.8  | 0.8  |
| gpt-4                 | 0.84      | 0.85     | 0.76     | 0.84 | 0.84 |

**Table 5.13:** Semantic similarity for different prompts used in code generation.

| Prompt  | Prompt1 | Prompt2 | Prompt3 | Prompt4 | Prompt5 |
|---------|---------|---------|---------|---------|---------|
| Prompt1 | 1.0     | 0.923   | 0.907   | 0.909   | 0.903   |
| Prompt2 | 0.923   | 1.0     | 0.915   | 0.887   | 0.925   |
| Prompt3 | 0.907   | 0.915   | 1.0     | 0.922   | 0.914   |
| Prompt4 | 0.909   | 0.887   | 0.922   | 1.0     | 0.88    |
| Prompt5 | 0.903   | 0.925   | 0.914   | 0.88    | 1.0     |

prompts. To address this concern, we conducted experiments on five different code generation prompts<sup>7</sup>.

**Semantic similarity for different prompts.** Before evaluating the CBS of LLM-generated code based on the guidance of different prompts, we first measure the semantic similarity of our five prompts. To measure the semantic similarity of our constructed five prompts, we follow the instructions provided by HuggingFace<sup>8</sup> to measure the semantic similarity between each pair of prompts. The evaluation results are shown in Table 5.13, where we observe that the semantic similarity between each pair of prompts is larger than 0.8<sup>9</sup>, which means that all prompts have the same objective, i.e., they require the LLM to generate code based on the task description.

**CBS for different prompts generated code.** Next, we use the provided five different prompts to guide LLMs to generate code based on task description and calculate the CBS of LLM-generated code. The evaluation results are shown in Table 5.14, where we can observe differences in GPT-3.5-turbo’s output for different prompts. For example, the CBS of age attribute ranges from 21.75% to 32.93%, indicating that LLMs may generate code with varying levels of bias depending on the prompt used to guide them. Despite the variations in CBS across different prompts, it is important to note that the CBS remains consistently high for all prompts. For instance, the CBS for the age attribute is consistently above 21% across all five prompts. This finding suggests that while the specific phrasing of the prompt can influence the extent of bias in the generated code, the overall presence of bias remains a significant concern regardless of the prompt used.

The impact of prompt phrasing on CBS is further evident when comparing the bias mitigation results of Scenario 1 (Section 5.3.4) to the direct LLM-generated code. While the results in Scenario 1 show some changes compared to the direct LLM-generated code, they remain similar in some cases, suggesting that the decreased bias in Scenario 1 may be partially attributed to the change in the prompt. However, in Scenario 2, where we provide our bias detection results to the LLM, the CBS is significantly reduced, and in

<sup>7</sup>See the prompts in [https://github.com/huangd1999/CBS/blob/main/different\\_prompt\\_code\\_generation.py](https://github.com/huangd1999/CBS/blob/main/different_prompt_code_generation.py)

<sup>8</sup>Semantic Similarity: <https://huggingface.co/tasks/sentence-similarity#passage-ranking>

<sup>9</sup>If the semantic similarity of two prompts larger than 0.8, we treat them as have the same meaning and goal. See [https://docs.lamaindex.ai/en/stable/api\\_reference/evaluation/semantic\\_similarity/](https://docs.lamaindex.ai/en/stable/api_reference/evaluation/semantic_similarity/).

**Table 5.14:** CBS for different prompts generated by GPT-3.5-turbo.

| Prompt  | Age   | Region | Gender | Education | Occupation | Race |
|---------|-------|--------|--------|-----------|------------|------|
| Prompt1 | 22.05 | 21.15  | 25.68  | 11.78     | 2.42       | 0.60 |
| Prompt2 | 29.00 | 30.21  | 44.71  | 22.05     | 2.42       | 0.91 |
| Prompt3 | 27.79 | 18.43  | 31.42  | 19.64     | 2.42       | 1.21 |
| Prompt4 | 32.93 | 23.26  | 29.61  | 21.45     | 1.81       | 0.91 |
| Prompt5 | 21.75 | 21.15  | 27.49  | 20.24     | 3.32       | 1.81 |

some cases, it decreases to 0. This stark contrast between Scenarios 1 and 2 indicates that our proposed method is effective in mitigating bias, rather than the changes being solely due to the prompt modification. The inclusion of bias detection feedback in Scenario 2 plays a crucial role in guiding the LLM to generate less biased code, demonstrating the effectiveness of our approach in addressing bias in LLM-generated code.

#### 5.4.4 Enhancing Value Pool for Bias Detection

demonstrates that our automated bias testing has high precision and recall. However, there are still a few false negatives (FN) due to the uncovered cases in the value pool for the protected attributes. This section explores strategies to enhance the value pool to reduce false negatives. In particular, the limitation observed with age parameters (i.e., where biases involving ages above 65 are not detected) suggests a gap in our testing scope. To mitigate this, a straightforward solution is to enrich our value pools with a broader range of values, aiming to improve the comprehensiveness of bias detection. Specifically, we add parameter values in ACSIncome, ACSEmployment, and ACSPublicCoverage<sup>10</sup> [48], thus improving the coverage of parameter values and addressing the gaps identified in our initial testing framework. The evaluation results are shown in Table 5.15, where we can observe that once we add more diverse values to the value pools, the false negative rate decreases to 0. Finally, the recall of automated test case analysis increases from 92% to 100%. However, our evaluation results also illustrate that this expansion of the value pool introduces extra overhead for the testing process. Specifically, the testing time increases from 57.15s to 3958.84s. The key reason is that once we increase the value pool for function parameters, the total test cases constructed by Figure 5.2 5a then largely increase. Considering the large overhead and the small ratio of false negatives, our default strategy does not adopt the large value pool, but users and developers can choose to adapt the value pool to achieve 100% test recall and precision when necessary.

#### 5.4.5 Why not use LLM to generate test cases?

We do not use LLM to generate test cases since LLM-generated test cases are often incorrect, which then requires significant manual efforts to select correct test cases from the generated tests. Besides, the amount of code evaluated in our experiments is

<sup>10</sup>ACSIncome, ACSEmployment, and ACSPublicCoverage provide a range of values for parameters in Table 5.1.

**Table 5.15:** Evaluation results for TP, FN, FP, and TN when we enrich value pools based on the ACSIncome, ACSEmployment, and ACSPublicCoverage dataset [48]. We also report the testing time in the overhead column.

|          | TP  | FN | FP | TN   | Overhead |
|----------|-----|----|----|------|----------|
| Original | 141 | 12 | 0  | 2185 | 57.16s   |
| Enriched | 153 | 0  | 0  | 2185 | 3958.84s |

**Table 5.16:** CBS for different temperatures generated by GPT-3.5-turbo.

| Temperature | Age         | Region     | Gender     | Education  | Occupation | Race      |
|-------------|-------------|------------|------------|------------|------------|-----------|
| 0.0         | 100 (29.07) | 44 (12.79) | 43 (12.5)  | 55 (15.99) | 17 (4.94)  | 12 (3.49) |
| 0.2         | 102 (29.65) | 40 (11.63) | 42 (12.21) | 55 (15.99) | 22 (6.4)   | 16 (4.65) |
| 0.4         | 111 (32.27) | 42 (12.21) | 46 (13.37) | 62 (18.02) | 20 (5.81)  | 15 (4.36) |
| 0.6         | 106 (30.81) | 40 (11.63) | 39 (11.34) | 61 (17.73) | 16 (4.65)  | 14 (4.07) |
| 0.8         | 92 (26.74)  | 37 (10.76) | 33 (9.59)  | 53 (15.41) | 25 (7.27)  | 13 (3.78) |
| 1.0         | 78 (23.35)  | 46 (13.77) | 76 (22.75) | 81 (24.25) | 5 (1.50)   | 6 (1.80)  |

extensive (e.g., 334 tasks \* 5 models \* 5 random generations \* 11 scenarios (5 different prompts \* 2 scenarios + 1 original code)), and the number of sensitive attributes may range from 1 to 5 for each provided code, which then requires large tokens to generate massive test cases to test the bias in the code. Therefore, we directly use our bias testing framework to construct test cases instead of relying on LLMs for test case generation, which is accurate and efficient.

#### 5.4.6 How does temperature affect the CBS of LLM-generated code?

As shown in Table 5.4, the LLM-generated code is different for different execution times, which causes the CBS\_U@5, CBS\_I@5, and CBS to vary across five executions. The key reason for these unstable results is that the temperature setting affects the next token selection. We set the temperature to 1.0 for all experiments. To investigate how temperature affects the bias of the LLM-generated code, we evaluate the CBS of the LLM-generated code at six different temperatures: 0.0, 0.2, 0.4, 0.6, 0.8, and 1.0. The evaluation results are shown in Tab. 5.16, where we can observe that first, the CBS of LLM-generated code varies at different temperatures. The key reason is that the temperature setting in LLMs influences the probability distribution over the next token during generation. Lower temperatures make the model more deterministic, favoring the most likely tokens, while higher temperatures introduce more randomness and encourage the model to explore less likely token choices. Second, we can also observe that even if the temperature is set to 0.0, where the LLM always selects the token with the highest probability in the last layer prediction as the next token during the inference time, the CBS of the LLM-generated code is still not 0, which indicates that during the inference time, LLMs are sometimes trying to generate task with bias behaviors.

#### 5.4.7 How about the token usage of bias mitigation process?

**Table 5.17:** Average token usage (input + output) for each LLM to mitigate bias code in LLM initial generated code in scenario1 and scenario1. The value of outside/inside the brackets is the scenario1 / scenario2 token usage.

| Model                 | Zero-Shot       | One-Shot        | Few-Shot         | CoT1             | CoT2             |
|-----------------------|-----------------|-----------------|------------------|------------------|------------------|
| palm-2-codechat-bison | 742.58 (864.06) | 687.71 (927.19) | 723.09 (1012.0)  | 735.3 (977.29)   | 739.82 (947.18)  |
| claude-instant-1      | 567.94 (728.0)  | 557.75 (723.97) | 601.61 (786.78)  | 591.84 (778.65)  | 553.58 (768.46)  |
| gpt-3.5-turbo         | 509.91 (817.95) | 627.96 (879.73) | 758.15 (975.67)  | 731.45 (938.31)  | 743.26 (901.21)  |
| gpt-4-1106-preview    | 900.99 (805.5)  | 969.46 (888.73) | 1036.59 (965.71) | 1044.17 (946.55) | 1064.04 (912.92) |
| gpt-4                 | 585.78 (625.74) | 633.47 (649.78) | 668.82 (710.28)  | 685.17 (700.04)  | 703.39 (691.61)  |

As the text window of LLMs is limited, which causes that we can not feed long text into LLMs to mitigate the bias of LLM itself generated code. Then, it is essential to ensure that our bias mitigation approach can effectively handle the programs within the text window. In this section, we measure the token usage of the bias mitigation process and discuss its implications for the scalability and applicability of our approach. We provided the average token usage (input + output) for each LLM used to mitigate bias in the initially generated code for both scenario 1 and scenario 2. As shown in Table 5.17, all LLMs for all prompts in both scenarios use less than 4096 tokens (the default token limitation of GPT-3.5, while other LLMs have larger text windows) to complete the tasks. For example, GPT-3.5-turbo only requires on average 758.15 tokens to finish each task, which are less than the 4096 tokens, which indicates that the size of the programs currently does not affect our approach. Furthermore, as current LLMs are extending their text windows, e.g., GPT-4 has 8K tokens, the newest GPT-3.5 has 16K tokens, and models such as Claude have a 128-200K token window, we believe that the input size in the future would also does not limit the effectiveness of our bias mitigation approach. As LLMs continue to increase their input size capabilities, our method will be able to handle even larger programs without encountering limitations related to token usage.

## 5.5 Threats to Validity

### 5.5.1 Internal Validity

The process of creating the code generation prompt dataset involves human judgment, which introduces the possibility of subjective bias in prompt design and may influence the presence or absence of certain biases in the dataset. To mitigate this threat, we ensure consistent and objective prompt creation by employing well-defined operational definitions for each bias type. Additionally, the code generation models may exhibit variations in generating code functions due to inherent randomness and model complexity, potentially impacting the results and introducing internal validity threats. To address this, we carefully control for such variations by running five times for each experiment (e.g., code generation and bias mitigation) to obtain the average results. Besides, we also utilize CBS\_U@K and CBS\_I@K to decrease the effect of variation for our experiments. These techniques help us reduce the impact of randomness and improve the robustness of our findings. By taking these precautions, we aim to strengthen the

internal validity of our research, ensuring the reliability and accuracy of the results obtained from the prompt dataset creation and code generation process.

### 5.5.2 External Validity

The external validity of our study is subject to the representativeness of the code generation prompt dataset and the generalizability of language models to various code generation tasks. If the dataset does not cover a representative range of potential biases in the code, our findings may lack generalizability to real-world scenarios. To address this concern, we take measures to ensure diversity in the selection of protected attributes and tasks and use the three most widely studied tasks in the fairness literature.

### 5.5.3 Construct Validity

For code bias evaluation, we rely on automated test case analysis to classify the predominantly generated code functions, providing a more standardized and automated approach. Then, for the code that requires human evaluation due to runtime errors, we have multiple experts to analyze bias types for each code to reduce subjectivity. The construct validity of our study also depends on the effectiveness of the test case analysis result assistant mitigation for the code. If the mitigation approach fails to result in substantial reductions in bias, the validity of our conclusions could be compromised. To mitigate this threat, we conduct comprehensive evaluations to assess five code generation models, test them five times, and report the average results. By doing so, we validate the effectiveness of our mitigation approach and strengthen the construct validity of our research findings.

## Chapter 6

# Conclusion and Future Work

In this thesis, we have addressed three critical challenges associated with code generation using LLMs: inefficiency, correctness, and social biases. Through the development of novel frameworks—**EffiLearner**, **EffiCoder**, and a code bias testing and mitigation system—we have made significant strides toward enhancing the practical utility of LLM-generated code. This concluding chapter summarizes our key contributions, discusses their implications for software engineering, and outlines directions for future research.

## 6.1 Summary of Contributions

### 6.1.1 Enhancing Code Efficiency with EffiLearner

In Chapter 3, we introduced **EffiLearner**, a framework designed to improve the efficiency of code generated by LLMs through execution profiling and iterative optimization. Recognizing that LLM-generated code often suffers from suboptimal execution times and resource utilization, **EffiLearner** emulates the optimization process employed by human coders. By executing the initial code and analyzing its performance profile—including execution time and memory usage—**EffiLearner** identifies inefficiencies and guides the LLM to iteratively refine the code. Our extensive experiments with multiple open-source and closed-source models demonstrated that **EffiLearner** significantly enhances code efficiency across various benchmarks without compromising functionality.

### 6.1.2 Improving Code Correctness and Efficiency with EffiCoder

Building upon the insights gained from **EffiLearner**, Chapter 4 introduced **EffiCoder**, a fine-tuning framework aimed at simultaneously improving code correctness and efficiency. We acknowledged the trade-off observed in **EffiLearner**, where optimizing for efficiency sometimes led to decreased code correctness. To address this, we constructed the *Effi-Code* dataset by aggregating and preprocessing code from multiple open-source sources, generating test cases, and applying iterative optimization over several cycles. Fine-tuning LLMs with **EffiCoder** resulted in models that produce code exhibiting higher correctness (pass@1 rates) and improved efficiency. Our experiments validated the effectiveness of **EffiCoder**, demonstrating its potential to enhance LLMs for practical



software engineering tasks.

### 6.1.3 Ensuring Social Fairness through Bias Detection and Mitigation

In Chapter 5, we tackled the critical issue of social biases in LLM-generated code by proposing a **code bias testing and mitigation framework**. Recognizing that biases in code logic can lead to unfair and discriminatory software behaviors, we developed methods to detect and mitigate such biases effectively. Our framework involves creating bias-sensitive code generation prompts, analyzing generated code using Abstract Syntax Trees (ASTs), and constructing test cases based on identified input parameters. We introduced metrics like **CBS**, **CBS\_U@K**, and **CBS\_I@K** to quantify bias in code generation. Furthermore, we explored various mitigation strategies, including few-shot learning and Chain-of-Thought prompting, finding that feeding back test analysis results to LLMs significantly reduces bias behaviors. Our work contributes to the development of ethically responsible AI by ensuring fairness in automated code generation.

## 6.2 Implications for Software Engineering

The contributions of this thesis have substantial implications for the field of software engineering, particularly in the integration of LLMs into the development lifecycle:

- **Enhanced Efficiency:** By improving the execution performance of LLM-generated code, developers can leverage LLMs in performance-critical applications and resource-constrained environments, broadening the applicability of automated code generation.
- **Improved Correctness:** Ensuring that generated code not only functions correctly but also adheres to efficiency standards increases trust in LLMs as reliable coding assistants, potentially accelerating development processes.
- **Ethical Assurance:** Addressing social biases in code generation mitigates legal risks and ethical concerns, promoting the development of fair and equitable software systems.

Collectively, these advancements support the responsible and sustainable adoption of LLMs in software engineering, aligning with industry goals of efficiency, reliability, and ethical integrity.

## 6.3 Future Work

While this thesis has made significant progress in addressing key challenges, several avenues for future research remain open.

### 6.3.1 Extending Methodologies to Other Programming Languages

Our frameworks primarily focused on code generation in specific programming languages such as Python. Extending EffiLearner and EffiCoder to support additional languages like Java, C++, or JavaScript would enhance their utility. This extension requires accounting for language-specific features, idioms, and optimization strategies.

### 6.3.2 Integrating Advanced Optimization Techniques

Future research could explore the integration of advanced optimization techniques into EffiLearner and EffiCoder. For instance, leveraging compiler optimization passes, static analysis tools, or machine learning models to predict code inefficiencies could further enhance optimization. Additionally, combining our frameworks with performance tuning methodologies used in high-performance computing could yield even greater efficiency gains.

### 6.3.3 Developing Comprehensive Bias Mitigation Strategies

While our bias mitigation framework effectively reduces certain biases, there is room to develop more comprehensive strategies. Future work might involve:

- **Incorporating Fairness Constraints:** Integrating fairness constraints directly into the LLM training objective could promote bias-free code generation at a foundational level.
- **Expanding Bias Definitions:** Investigating a broader spectrum of biases, including intersectional and context-specific biases, to ensure more thorough mitigation.
- **Dynamic Bias Detection:** Developing real-time bias detection mechanisms that can be integrated into development environments, providing immediate feedback to developers.

### 6.3.4 User Studies and Real-world Applications

Evaluating our frameworks in real-world settings is crucial for understanding their practical impact. Conducting user studies with software developers can provide insights into the usability and effectiveness of EffiLearner, EffiCoder, and our bias mitigation tools. Such studies could inform refinements to the frameworks and encourage adoption in industry.

### 6.3.5 Exploring Ethical and Legal Implications

As LLMs become more integrated into software development, understanding the ethical and legal implications of their use becomes increasingly important. Future research could explore:

- **Regulatory Compliance:** Ensuring that code generated by LLMs complies with

legal standards and regulations related to fairness, privacy, and security.

- **Accountability Mechanisms:** Developing frameworks to attribute responsibility for code generated by AI, particularly in cases where biases or errors lead to negative outcomes.
- **Ethical Guidelines:** Formulating industry-wide guidelines for the ethical use of LLMs in code generation, informed by interdisciplinary collaboration.

## 6.4 Final Remarks

The rapid advancement of LLMs presents both opportunities and challenges for software engineering. This thesis has contributed to harnessing these opportunities—improving efficiency, correctness, and fairness—while addressing the associated challenges. By developing EffiLearner and EffiCoder, we have shown that it is possible to produce LLM-generated code that meets high standards of performance and reliability. Our bias testing and mitigation framework advances the pursuit of ethical AI by promoting fairness in automated code generation.

As LLMs continue to evolve and integrate into diverse aspects of software development, it is imperative to ensure that they do so responsibly. The methodologies and insights presented in this thesis lay the groundwork for future innovations that prioritize not only technological advancement but also ethical considerations. We hope that this work inspires continued research and collaboration in creating AI systems that are efficient, accurate, and just, ultimately contributing to the betterment of software engineering and society at large.

# Bibliography

- [1] M. I. Abdin, S. A. Jacobs, A. A. Awan, J. Aneja, A. Awadallah, H. Awadalla, N. Bach, A. Bahree, A. Bakhtiari, H. S. Behl, A. Benhaim, M. Bilenko, J. Bjorck, S. Bubeck, M. Cai, C. C. T. Mendes, W. Chen, V. Chaudhary, P. Chopra, A. D. Giorno, G. de Rosa, M. Dixon, R. Eldan, D. Iter, A. Garg, A. Goswami, S. Gunasekar, E. Haider, J. Hao, R. J. Hewett, J. Huynh, M. Javaheripi, X. Jin, P. Kauffmann, N. Karampatziakis, D. Kim, M. Khademi, L. Kurilenko, J. R. Lee, Y. T. Lee, Y. Li, C. Liang, W. Liu, E. Lin, Z. Lin, P. Madan, A. Mitra, H. Modi, A. Nguyen, B. Norick, B. Patra, D. Perez-Becker, T. Portet, R. Pryzant, H. Qin, M. Radmilac, C. Rosset, S. Roy, O. Ruwase, O. Saarikivi, A. Saied, A. Salim, M. Santacroce, S. Shah, N. Shang, H. Sharma, X. Song, M. Tanaka, X. Wang, R. Ward, G. Wang, P. Witte, M. Wyatt, C. Xu, J. Xu, S. Yadav, F. Yang, Z. Yang, D. Yu, C. Zhang, C. Zhang, J. Zhang, L. L. Zhang, Y. Zhang, Y. Zhang, Y. Zhang, and X. Zhou. “Phi-3 Technical Report: A Highly Capable Language Model Locally on Your Phone”. In: *CoRR* abs/2404.14219 (2024). DOI: [10.48550/ARXIV.2404.14219](https://doi.org/10.48550/ARXIV.2404.14219). arXiv: [2404.14219](https://arxiv.org/abs/2404.14219). URL: <https://doi.org/10.48550/arXiv.2404.14219>.
- [2] *Adult Income Dataset*. [www.kaggle.com/datasets/wenruliu/adult-income-dataset](https://www.kaggle.com/datasets/wenruliu/adult-income-dataset). Accessed on August 1, 2023. 2023.
- [3] N. Ahmad and A. N. Abd Alla. “Smart Evaluation for Job Vacancy Application System”. In: *2009 Second International Conference on the Applications of Digital Information and Web Technologies*. IEEE. 2009, pp. 452–455.
- [4] W. U. Ahmad, M. G. R. Tushar, S. Chakraborty, and K. Chang. “AVATAR: A Parallel Corpus for Java-Python Program Translation”. In: *Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023*. Ed. by A. Rogers, J. L. Boyd-Graber, and N. Okazaki. Association for Computational Linguistics, 2023, pp. 2268–2281. DOI: [10.18653/v1/2023.FINDINGS-ACL.143](https://doi.org/10.18653/v1/2023.findings-acl.143). URL: <https://doi.org/10.18653/v1/2023.findings-acl.143>.
- [5] T. Ahmed and P. T. Devanbu. “Few-shot training LLMs for project-specific code-summarization”. In: *37th IEEE/ACM International Conference on Automated Software Engineering, ASE 2022, Rochester, MI, USA, October 10-14, 2022*. ACM, 2022, 177:1–177:5. DOI: [10.1145/3551349.3559555](https://doi.org/10.1145/3551349.3559555). URL: <https://doi.org/10.1145/3551349.3559555>.
- [6] J.-B. Alayrac, J. Donahue, P. Luc, A. Miech, I. Barr, Y. Hasson, K. Lenc, A. Mensch, K. Millican, M. Reynolds, et al. “Flamingo: a visual language model for few-shot learning”. In: *Advances in Neural Information Processing Systems* 35 (2022), pp. 23716–23736.

- [7] L. B. Allal, R. Li, D. Kocetkov, C. Mou, C. Akiki, C. M. Ferrandis, N. Muennighoff, M. Mishra, A. Gu, M. Dey, L. K. Umapathi, C. J. Anderson, Y. Zi, J. Lamy-Poirier, H. Schoelkopf, S. Troshin, D. Abulkhanov, M. Romero, M. Lappert, F. D. Toni, B. G. del Río, Q. Liu, S. Bose, U. Bhattacharyya, T. Y. Zhuo, I. Yu, P. Villegas, M. Zocca, S. Mangrulkar, D. Lansky, H. Nguyen, D. Contractor, L. Villa, J. Li, D. Bahdanau, Y. Jernite, S. Hughes, D. Fried, A. Guha, H. de Vries, and L. von Werra. “SantaCoder: don’t reach for the stars!” In: *CoRR* abs/2301.03988 (2023). DOI: [10.48550/ARXIV.2301.03988](https://doi.org/10.48550/ARXIV.2301.03988). arXiv: [2301.03988](https://arxiv.org/abs/2301.03988). URL: <https://doi.org/10.48550/arXiv.2301.03988>.
- [8] J. M. Alvarez, K. M. Scott, B. Berendt, and S. Ruggieri. “Domain Adaptive Decision Trees: Implications for Accuracy and Fairness”. In: *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*. 2023, pp. 423–433.
- [9] G. Andreeva, J. Ansell, and J. Crook. “Impact of anti-discrimination laws on credit scoring”. In: *Journal of Financial Services Marketing* 9 (2004), pp. 22–33.
- [10] R. Anil, S. Borgeaud, Y. Wu, J. Alayrac, J. Yu, R. Soricut, J. Schalkwyk, A. M. Dai, A. Hauth, K. Millican, D. Silver, S. Petrov, M. Johnson, I. Antonoglou, J. Schrittwieser, A. Glaese, J. Chen, E. Pitler, T. P. Lillicrap, A. Lazaridou, O. Firat, J. Molloy, M. Isard, P. R. Barham, T. Hennigan, B. Lee, F. Viola, M. Reynolds, Y. Xu, R. Doherty, E. Collins, C. Meyer, E. Rutherford, E. Moreira, K. Ayoub, M. Goel, G. Tucker, E. Piqueras, M. Krikun, I. Barr, N. Savinov, I. Danihelka, B. Roelofs, A. White, A. Andreassen, T. von Glehn, L. Yagati, M. Kazemi, L. Gonzalez, M. Khalman, J. Sygnowski, and et al. “Gemini: A Family of Highly Capable Multimodal Models”. In: *CoRR* abs/2312.11805 (2023). DOI: [10.48550/ARXIV.2312.11805](https://doi.org/10.48550/ARXIV.2312.11805). arXiv: [2312.11805](https://arxiv.org/abs/2312.11805). URL: <https://doi.org/10.48550/arXiv.2312.11805>.
- [11] Anthropic. *Introducing the next generation of Claude*. 2024. URL: <https://www.anthropic.com/news/claude-3-family>.
- [12] E. O. Arceo-Gomez, R. M. Campos-Vazquez, R. Y. Badillo, and S. Lopez-Araiza. “Gender stereotypes in job advertisements: What do they imply for the gender salary gap?” In: *Journal of Labor Research* 43.1 (2022), pp. 65–102.
- [13] B. Athiwaratkun, S. K. Gouda, Z. Wang, X. Li, Y. Tian, M. Tan, W. U. Ahmad, S. Wang, Q. Sun, M. Shang, S. K. Gonugondla, H. Ding, V. Kumar, N. Fulton, A. Farahani, S. Jain, R. Giaquinto, H. Qian, M. K. Ramanathan, and R. Nallapati. “Multilingual Evaluation of Code Generation Models”. In: *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. URL: <https://openreview.net/pdf?id=Bo7eeXm6An8>.
- [14] J. Austin, A. Odena, M. Nye, M. Bosma, H. Michalewski, D. Dohan, E. Jiang, C. J. Cai, M. Terry, Q. V. Le, and C. Sutton. “Program Synthesis with Large Language Models”. In: *ArXiv* abs/2108.07732 (2021). URL: <https://api.semanticscholar.org/CorpusID:237142385>.
- [15] S. Barikeri, A. Lauscher, I. Vulić, and G. Glavaš. “RedditBias: A real-world resource for bias evaluation and debiasing of conversational language models”. In: *arXiv preprint arXiv:2106.03521* (2021).

- [16] P. Barlas, K. Kyriakou, O. Guest, S. Kleanthous, and J. Otterbacher. "To" see" is to stereotype: Image tagging algorithms, gender recognition, and the accuracy-fairness trade-off". In: *Proceedings of the ACM on Human-Computer Interaction* 4.CSCW3 (2021), pp. 1–31.
- [17] M. Behroozi, C. Parnin, and T. Barik. "Hiring is broken: What do developers say about technical interviews?" In: *2019 IEEE Symposium on Visual Languages and Human-Centric Computing (VL/HCC)*. IEEE, 2019, pp. 1–9.
- [18] B. A. Bell. "Understanding the Preparation Phase of Technical Interviews". PhD thesis. Virginia Tech, 2023.
- [19] P. Besse, E. del Barrio, P. Gordaliza, J.-M. Loubes, and L. Risser. "A survey of bias in machine learning through the prism of statistical parity". In: *The American Statistician* 76.2 (2022), pp. 188–198.
- [20] B. Bharti, P. Yi, and J. Sulam. "Estimating and Controlling for Equalized Odds via Sensitive Attribute Predictors". In: *Advances in Neural Information Processing Systems* 36 (2024).
- [21] BigCode. *Self-OSS-Instruct-SC2-Exec-Filter-50K*. <https://huggingface.co/datasets/bigcode/self-oss-instruct-sc2-exec-filter-50k>. 2023.
- [22] S. Biswas and H. Rajan. "Do the Machine Learning Models on a Crowd Sourced Platform Exhibit Bias? An Empirical Study on Model Fairness". In: *ESEC/FSE'2020: The 28th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. Sacramento, California, United States, Nov. 2020.
- [23] S. Biswas and H. Rajan. "Fairify: Fairness Verification of Neural Networks". In: *ICSE'2023: The 45th International Conference on Software Engineering*. Melbourne, Australia, May 2023.
- [24] E. M. Boyd and A. W. Fales. "Reflective learning: Key to learning from experience". In: *Journal of humanistic psychology* 23.2 (1983), pp. 99–117.
- [25] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei. "Language Models are Few-Shot Learners". In: *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*. Ed. by H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin. 2020. URL: <https://proceedings.neurips.cc/paper/2020/hash/1457c0d6bfc4967418bfb8ac142f64a-Abstract.html>.
- [26] E. Capra, C. Francalanci, and S. A. Slaughter. "Is software "green"? Application development environments and energy efficiency in open source applications". In: *Information and Software Technology* 54.1 (2012), pp. 60–71.
- [27] F. Cassano, J. Gouwar, D. Nguyen, S. Nguyen, L. Phipps-Costin, D. Pinckney, M. Yee, Y. Zi, C. J. Anderson, M. Q. Feldman, A. Guha, M. Greenberg, and A. Jangda. "MultiPL-E: A Scalable and Polyglot Approach to Benchmarking Neural Code

- Generation". In: *IEEE Trans. Software Eng.* 49.7 (2023), pp. 3675–3691. DOI: [10.1109/TSE.2023.3267446](https://doi.org/10.1109/TSE.2023.3267446). URL: <https://doi.org/10.1109/TSE.2023.3267446>.
- [28] H. Chang and R. Shokri. "On the privacy risks of algorithmic fairness". In: *2021 IEEE European Symposium on Security and Privacy (EuroS&P)*. IEEE. 2021, pp. 292–303.
- [29] S. Chaudhary. *Code Alpaca: An Instruction-following LLaMA model for code generation*. <https://github.com/sahil280114/codealpaca>. 2023.
- [30] 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". In: *CoRR* abs/2107.03374 (2021). arXiv: [2107.03374](https://arxiv.org/abs/2107.03374). URL: <https://arxiv.org/abs/2107.03374>.
- [31] M. Chen, J. Tworek, H. Jun, Q. Yuan, H. P. d. O. Pinto, J. Kaplan, H. Edwards, Y. Burda, N. Joseph, G. Brockman, et al. "Evaluating large language models trained on code". In: *arXiv preprint arXiv:2107.03374* (2021).
- [32] X. Chen, M. Lin, N. Schärli, and D. Zhou. "Teaching Large Language Models to Self-Debug". In: *CoRR* abs/2304.05128 (2023). DOI: [10.48550/ARXIV.2304.05128](https://doi.org/10.48550/ARXIV.2304.05128). arXiv: [2304.05128](https://arxiv.org/abs/2304.05128). URL: <https://doi.org/10.48550/ARXIV.2304.05128>.
- [33] Z. Chen, J. Zhang, F. Sarro, and M. Harman. "Fairness Improvement with Multiple Protected Attributes: How Far Are We?" In: *46th International Conference on Software Engineering (ICSE 2024)*. ACM. 2023.
- [34] Z. Chen, J. M. Zhang, F. Sarro, and M. Harman. "A Comprehensive Empirical Study of Bias Mitigation Methods for Machine Learning Classifiers". In: *ACM Transactions on Software Engineering and Methodology* 32.4 (2023), pp. 1–30.
- [35] Z. Chen, J. M. Zhang, F. Sarro, and M. Harman. "Fairness Improvement with Multiple Protected Attributes: How Far Are We?" In: *IEEE/ACM*. 2024.
- [36] Z. Chen, J. M. Zhang, F. Sarro, and M. Harman. "MAAT: a novel ensemble approach to addressing fairness and performance bugs for machine learning software". In: *Proceedings of the 30th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 2022, pp. 1122–1134.
- [37] A. Chouldechova and A. Roth. "The frontiers of fairness in machine learning". In: *arXiv preprint arXiv:1810.08810* (2018).
- [38] Z. Chu, J. Chen, Q. Chen, W. Yu, T. He, H. Wang, W. Peng, M. Liu, B. Qin, and T. Liu. "A survey of chain of thought reasoning: Advances, frontiers and future". In: *arXiv preprint arXiv:2309.15402* (2023).



- [39] T. Coignion, C. Quinton, and R. Rouvoy. "ââ". In: *Proceedings of the 28th International Conference on Evaluation and Assessment in Software Engineering*. 2024, pp. 79–89.
- [40] C. Computations. *Dolphin Coder*. <https://huggingface.co/datasets/cognitivecomputations/dolphin-coder>. 2023.
- [41] T. Computer. *Glaive-Code-Assistant*. <https://huggingface.co/datasets/Together/Glaive-Code-Assistant>. 2023.
- [42] A. F. Cooper, E. Abrams, and N. Na. "Emergent unfairness in algorithmic fairness-accuracy trade-off research". In: *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*. 2021, pp. 46–54.
- [43] S. Corbett-Davies and S. Goel. "The measure and mismeasure of fairness: A critical review of fair machine learning". In: *arXiv preprint arXiv:1808.00023* (2018).
- [44] A. F. Cruz, C. Belém, S. Jesus, J. Bravo, P. Saleiro, and P. Bizarro. "Fairgbm: Gradient boosting with fairness constraints". In: *arXiv preprint arXiv:2209.07850* (2022).
- [45] A. F. Cruz and M. Hardt. "Unprocessing Seven Years of Algorithmic Fairness". In: *arXiv preprint arXiv:2306.07261* (2023).
- [46] DeepSeekAI. *DeepSeek Coder: Let the Code Write Itself*. 2023. URL: <https://deepseekcoder.github.io/>.
- [47] Y. Deng, C. S. Xia, C. Yang, S. D. Zhang, S. Yang, and L. Zhang. "Large Language Models are Edge-Case Fuzzers: Testing Deep Learning Libraries via FuzzGPT". In: *CoRR abs/2304.02014* (2023). DOI: [10.48550/ARXIV.2304.02014](https://doi.org/10.48550/ARXIV.2304.02014). arXiv: [2304.02014](https://arxiv.org/abs/2304.02014). URL: <https://doi.org/10.48550/arXiv.2304.02014>.
- [48] F. Ding, M. Hardt, J. Miller, and L. Schmidt. "Retiring adult: New datasets for fair machine learning". In: *Advances in neural information processing systems* 34 (2021), pp. 6478–6490.
- [49] Y. Ding, Z. Wang, W. U. Ahmad, M. K. Ramanathan, R. Nallapati, P. Bhatia, D. Roth, and B. Xiang. "CoCoMIC: Code Completion By Jointly Modeling In-file and Cross-file Context". In: *CoRR abs/2212.10007* (2022). DOI: [10.48550/ARXIV.2212.10007](https://doi.org/10.48550/ARXIV.2212.10007). arXiv: [2212.10007](https://arxiv.org/abs/2212.10007). URL: <https://doi.org/10.48550/arXiv.2212.10007>.
- [50] M. Du, A. T. Luu, B. Ji, and S. Ng. "Mercury: An Efficiency Benchmark for LLM Code Synthesis". In: *CoRR abs/2402.07844* (2024). DOI: [10.48550/ARXIV.2402.07844](https://doi.org/10.48550/ARXIV.2402.07844). arXiv: [2402.07844](https://arxiv.org/abs/2402.07844). URL: <https://doi.org/10.48550/arXiv.2402.07844>.
- [51] M. Du, A. T. Luu, B. Ji, and S.-K. Ng. "Mercury: An efficiency benchmark for llm code synthesis". In: *arXiv preprint arXiv:2402.07844* (2024).
- [52] S. Dutta, D. Wei, H. Yueksel, P.-Y. Chen, S. Liu, and K. Varshney. "Is there a trade-off between fairness and accuracy? a perspective using mismatched hypothesis testing". In: *International conference on machine learning*. PMLR. 2020, pp. 2803–2813.
- [53] *Employee Dataset*. [www.kaggle.com/datasets/tawfikelmetwally/employee-dataset](https://www.kaggle.com/datasets/tawfikelmetwally/employee-dataset). Accessed on August 1, 2023. 2023.



- [54] 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". In: *2023 IEEE/ACM International Conference on Software Engineering: Future of Software Engineering (ICSE-FoSE)*. IEEE. 2023, pp. 31–53.
- [55] V. K. Felkner, H.-C. H. Chang, E. Jang, and J. May. "WinoQueer: A Community-in-the-Loop Benchmark for Anti-LGBTQ+ Bias in Large Language Models". In: *arXiv preprint arXiv:2306.15087* (2023).
- [56] J. Ferry. "Addressing interpretability fairness & privacy in machine learning through combinatorial optimization methods". PhD thesis. Université Paul Sabatier-Toulouse III, 2023.
- [57] J. Ferry, U. Aivodji, S. Gambs, M.-J. Huguet, and M. Siala. "Exploiting Fairness to Enhance Sensitive Attributes Reconstruction". In: *2023 IEEE Conference on Secure and Trustworthy Machine Learning (SaTML)*. IEEE. 2023, pp. 18–41.
- [58] E. Fleisig and C. Fellbaum. "Mitigating Gender Bias in Machine Translation through Adversarial Learning". In: *arXiv preprint arXiv:2203.10675* (2022).
- [59] D. Fried, A. Aghajanyan, J. Lin, S. Wang, E. Wallace, F. Shi, R. Zhong, S. Yih, L. Zettlemoyer, and M. Lewis. "InCoder: A Generative Model for Code Infilling and Synthesis". In: *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. URL: <https://openreview.net/pdf?id=hQwb-1bM6EL>.
- [60] S. A. Friedler, C. Scheidegger, S. Venkatasubramanian, S. Choudhary, E. P. Hamilton, and D. Roth. "A comparative study of fairness-enhancing interventions in machine learning". In: *Proceedings of the conference on fairness, accountability, and transparency*. 2019, pp. 329–338.
- [61] L. Gao, Z. Dai, P. Pasupat, A. Chen, A. T. Chaganty, Y. Fan, V. Y. Zhao, N. Lao, H. Lee, D. Juan, and K. Guu. "RARR: Researching and Revising What Language Models Say, Using Language Models". In: *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023*. Ed. by A. Rogers, J. L. Boyd-Graber, and N. Okazaki. Association for Computational Linguistics, 2023, pp. 16477–16508. DOI: [10.18653/v1/2023.ACL-LONG.910](https://doi.org/10.18653/v1/2023.ACL-LONG.910). URL: <https://doi.org/10.18653/v1/2023.acl-long.910>.
- [62] J. Gardner, Z. Popovic, and L. Schmidt. "Subgroup robustness grows on trees: An empirical baseline investigation". In: *Advances in Neural Information Processing Systems* 35 (2022), pp. 9939–9954.
- [63] A. Glaese, N. McAleese, M. Trebacz, J. Aslanides, V. Firoiu, T. Ewalds, M. Rauh, L. Weidinger, M. J. Chadwick, P. Thacker, L. Campbell-Gillingham, J. Uesato, P. Huang, R. Comanescu, F. Yang, A. See, S. Dathathri, R. Greig, C. Chen, D. Fritz, J. S. Elias, R. Green, S. Mokrá, N. Fernando, B. Wu, R. Foley, S. Young, I. Gabriel, W. Isaac, J. Mellor, D. Hassabis, K. Kavukcuoglu, L. A. Hendricks, and G. Irving. "Improving alignment of dialogue agents via targeted human judgements". In: *CoRR abs/2209.14375* (2022). DOI: [10.48550/ARXIV.2209.14375](https://doi.org/10.48550/ARXIV.2209.14375). arXiv: [2209.14375](https://doi.org/10.48550/ARXIV.2209.14375). URL: <https://doi.org/10.48550/arXiv.2209.14375>.

- [64] Z. Gou, Z. Shao, Y. Gong, Y. Shen, Y. Yang, N. Duan, and W. Chen. "CRITIC: Large Language Models Can Self-Correct with Tool-Interactive Critiquing". In: *CoRR* abs/2305.11738 (2023). DOI: [10.48550/ARXIV.2305.11738](https://doi.org/10.48550/ARXIV.2305.11738). arXiv: [2305.11738](https://arxiv.org/abs/2305.11738). URL: <https://doi.org/10.48550/arXiv.2305.11738>.
- [65] D. Guo, Q. Zhu, D. Yang, Z. Xie, K. Dong, W. Zhang, G. Chen, X. Bi, Y. Wu, Y. Li, et al. "DeepSeek-Coder: When the Large Language Model Meets Programming—The Rise of Code Intelligence". In: *arXiv preprint arXiv:2401.14196* (2024).
- [66] X. Han, Z. Jiang, H. Jin, Z. Liu, N. Zou, Q. Wang, and X. Hu. "Retiring DP: New Distribution-Level Metrics for Demographic Parity". In: *Transactions on Machine Learning Research* (2023).
- [67] M. M. A. Haque, W. U. Ahmad, I. Lourentzou, and C. Brown. "FixEval: Execution-based Evaluation of Program Fixes for Competitive Programming Problems". In: *CoRR* abs/2206.07796 (2022). DOI: [10.48550/ARXIV.2206.07796](https://doi.org/10.48550/ARXIV.2206.07796). arXiv: [2206.07796](https://arxiv.org/abs/2206.07796). URL: <https://doi.org/10.48550/arXiv.2206.07796>.
- [68] J. Harper. "Interview Insight: How to Get the Job". In: *A Software Engineer's Guide to Seniority: A Guide to Technical Leadership*. Springer, 2022, pp. 19–28.
- [69] M. Hasan, T. Muttaqueen, A. A. Ishtiaq, K. S. Mehrab, M. M. A. Haque, T. Hasan, W. U. Ahmad, A. Iqbal, and R. Shahriyar. "CoDesc: A Large Code-Description Parallel Dataset". In: *Findings of the Association for Computational Linguistics: ACL/IJCNLP 2021, Online Event, August 1-6, 2021*. Ed. by C. Zong, F. Xia, W. Li, and R. Navigli. Vol. ACL/IJCNLP 2021. Findings of ACL. Association for Computational Linguistics, 2021, pp. 210–218. DOI: [10.18653/v1/2021.FINDINGS-ACL.18](https://doi.org/10.18653/v1/2021.findings-acl.18). URL: <https://doi.org/10.18653/v1/2021.findings-acl.18>.
- [70] D. Hendrycks, S. Basart, S. Kadavath, M. Mazeika, A. Arora, E. Guo, C. Burns, S. Puranik, H. He, D. Song, and J. Steinhardt. "Measuring Coding Challenge Competence With APPS". In: *NeurIPS* (2021).
- [71] M. Hernandez, D. R. Avery, S. D. Volpone, and C. R. Kaiser. "Bargaining while Black: The role of race in salary negotiations." In: *Journal of Applied Psychology* 104.4 (2019), p. 581.
- [72] J. Hong, N. Lee, and J. Thorne. *ORPO: Monolithic Preference Optimization without Reference Model*. 2024. arXiv: [2403.07691](https://arxiv.org/abs/2403.07691) [cs.CL]. URL: <https://arxiv.org/abs/2403.07691>.
- [73] D. Huang, Q. Bu, and H. Cui. "CodeCoT and Beyond: Learning to Program and Test like a Developer". In: *ArXiv* abs/2308.08784 (2023). URL: <https://api.semanticscholar.org/CorpusID:261030533>.
- [74] D. Huang, Q. Bu, and H. Cui. "CodeCoT and Beyond: Learning to Program and Test like a Developer". In: *arXiv preprint arXiv:2308.08784* (2023).
- [75] D. Huang, Q. Bu, J. Zhang, X. Xie, J. Chen, and H. Cui. "Bias assessment and mitigation in llm-based code generation". In: *arXiv preprint arXiv:2309.14345* (2023).
- [76] D. Huang, Q. Bu, J. M. Zhang, M. Luck, and H. Cui. "AgentCoder: Multi-Agent-based Code Generation with Iterative Testing and Optimisation". In: *arXiv preprint arXiv:2312.13010* (2023).

- [77] D. Huang, J. Dai, H. Weng, P. Wu, Y. Qing, J. M. Zhang, H. Cui, and Z. Guo. "SOAP: Enhancing Efficiency of Generated Code via Self-Optimization". In: *arXiv preprint arXiv:2405.15189* (2024).
- [78] D. Huang, J. Dai, H. Weng, P. Wu, Y. Qing, J. M. Zhang, H. Cui, and Z. Guo. "SOAP: Enhancing Efficiency of Generated Code via Self-Optimization". In: *CoRR abs/2405.15189* (2024). DOI: [10.48550/ARXIV.2405.15189](https://doi.org/10.48550/ARXIV.2405.15189). arXiv: [2405.15189](https://arxiv.org/abs/2405.15189). URL: <https://doi.org/10.48550/arXiv.2405.15189>.
- [79] D. Huang, J. M. Zhang, Y. Qing, and H. Cui. "EffiBench: Benchmarking the Efficiency of Automatically Generated Code". In: *arXiv preprint arXiv:2402.02037* (2024).
- [80] B. Hui, J. Yang, Z. Cui, J. Yang, D. Liu, L. Zhang, T. Liu, J. Zhang, B. Yu, K. Dang, A. Yang, R. Men, F. Huang, X. Ren, X. Ren, J. Zhou, and J. Lin. "Qwen2.5-Coder Technical Report". In: 2024. URL: <https://api.semanticscholar.org/CorpusID:272707390>.
- [81] G. Izacard, P. Lewis, M. Lomeli, L. Hosseini, F. Petroni, T. Schick, J. Dwivedi-Yu, A. Joulin, S. Riedel, and E. Grave. "Few-shot learning with retrieval augmented language models". In: *arXiv preprint arXiv:2208.03299* (2022).
- [82] N. Jain, S. Vaidyanath, A. S. Iyer, N. Natarajan, S. Parthasarathy, S. K. Rajamani, and R. Sharma. "Jigsaw: Large Language Models meet Program Synthesis". In: *44th IEEE/ACM 44th International Conference on Software Engineering, ICSE 2022, Pittsburgh, PA, USA, May 25-27, 2022*. ACM, 2022, pp. 1219–1231. DOI: [10.1145/3510003.3510203](https://doi.org/10.1145/3510003.3510203). URL: <https://doi.org/10.1145/3510003.3510203>.
- [83] N. Jiang, K. Liu, T. Lutellier, and L. Tan. "Impact of Code Language Models on Automated Program Repair". In: *45th IEEE/ACM International Conference on Software Engineering, ICSE 2023, Melbourne, Australia, May 14-20, 2023*. IEEE, 2023, pp. 1430–1442. DOI: [10.1109/ICSE48619.2023.00125](https://doi.org/10.1109/ICSE48619.2023.00125). URL: <https://doi.org/10.1109/ICSE48619.2023.00125>.
- [84] S. Jiang, Y. Wang, and Y. Wang. "SelfEvolve: A Code Evolution Framework via Large Language Models". In: *CoRR abs/2306.02907* (2023). DOI: [10.48550/ARXIV.2306.02907](https://doi.org/10.48550/ARXIV.2306.02907). arXiv: [2306.02907](https://arxiv.org/abs/2306.02907). URL: <https://doi.org/10.48550/arXiv.2306.02907>.
- [85] C. E. Jimenez, J. Yang, A. Wettig, S. Yao, K. Pei, O. Press, and K. Narasimhan. "SWE-bench: Can Language Models Resolve Real-World GitHub Issues?" In: *CoRR abs/2310.06770* (2023). DOI: [10.48550/ARXIV.2310.06770](https://doi.org/10.48550/ARXIV.2310.06770). arXiv: [2310.06770](https://arxiv.org/abs/2310.06770). URL: <https://doi.org/10.48550/arXiv.2310.06770>.
- [86] J. Kang, T. Xie, X. Wu, R. Maciejewski, and H. Tong. "Multifair: Multi-group fairness in machine learning". In: *arXiv preprint arXiv:2105.11069* (2021).
- [87] M. Kearns, S. Neel, A. Roth, and Z. S. Wu. "An empirical study of rich subgroup fairness for machine learning". In: *Proceedings of the conference on fairness, accountability, and transparency*. 2019, pp. 100–109.
- [88] J. Komiyama and H. Shima. "Two-stage algorithm for fairness-aware machine learning". In: *arXiv preprint arXiv:1710.04924* (2017).

- [89] J. Kreutzer, S. Khadivi, E. Matusov, and S. Riezler. "Can Neural Machine Translation be Improved with User Feedback?" In: *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2018, New Orleans, Louisiana, USA, June 1-6, 2018, Volume 3 (Industry Papers)*. Ed. by S. Bangalore, J. Chu-Carroll, and Y. Li. Association for Computational Linguistics, 2018, pp. 92–105. DOI: [10.18653/V1/N18-3012](https://doi.org/10.18653/V1/N18-3012). URL: <https://doi.org/10.18653/v1/n18-3012>.
- [90] Y. Lai, C. Li, Y. Wang, T. Zhang, R. Zhong, L. Zettlemoyer, W. Yih, D. Fried, S. I. Wang, and T. Yu. "DS-1000: A Natural and Reliable Benchmark for Data Science Code Generation". In: *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*. Ed. by A. Krause, E. Brunskill, K. Cho, B. Engelhardt, S. Sabato, and J. Scarlett. Vol. 202. Proceedings of Machine Learning Research. PMLR, 2023, pp. 18319–18345. URL: <https://proceedings.mlr.press/v202/lai23b.html>.
- [91] T. Le Quy, A. Roy, V. Iosifidis, W. Zhang, and E. Ntoutsis. "A survey on datasets for fairness-aware machine learning". In: *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery* 12.3 (2022), e1452.
- [92] H. Lee, S. Hong, J. Park, T. Kim, G. Kim, and J.-W. Ha. "KoSBI: A Dataset for Mitigating Social Bias Risks Towards Safer Large Language Model Application". In: *arXiv preprint arXiv:2305.17701* (2023).
- [93] *LeetCode*. <https://leetcode.com/>. Accessed: January 31, 2024.
- [94] C. Lemieux, J. P. Inala, S. K. Lahiri, and S. Sen. "CodaMosa: Escaping Coverage Plateaus in Test Generation with Pre-trained Large Language Models". In: *45th IEEE/ACM International Conference on Software Engineering, ICSE 2023, Melbourne, Australia, May 14-20, 2023*. IEEE, 2023, pp. 919–931. DOI: [10.1109/ICSE48619.2023.00085](https://doi.org/10.1109/ICSE48619.2023.00085). URL: <https://doi.org/10.1109/ICSE48619.2023.00085>.
- [95] R. Li, L. B. Allal, Y. Zi, N. Muennighoff, D. Kocetkov, C. Mou, M. Marone, C. Akiki, J. Li, J. Chim, Q. Liu, E. Zheltonozhskii, T. Y. Zhuo, T. Wang, O. Dehaene, M. Davaadorj, J. Lamy-Poirier, J. Monteiro, O. Shliazhko, N. Gontier, N. Meade, A. Zebaze, M. Yee, L. K. Umapathi, J. Zhu, B. Lipkin, M. Oblokulov, Z. Wang, R. M. V, J. Stillerman, S. S. Patel, D. Abulkhanov, M. Zocca, M. Dey, Z. Zhang, N. Moustafa-Fahmy, U. Bhattacharyya, W. Yu, S. Singh, S. Luccioni, P. Villegas, M. Kunakov, F. Zhdanov, M. Romero, T. Lee, N. Timor, J. Ding, C. Schlesinger, H. Schoelkopf, J. Ebert, T. Dao, M. Mishra, A. Gu, J. Robinson, C. J. Anderson, B. Dolan-Gavitt, D. Contractor, S. Reddy, D. Fried, D. Bahdanau, Y. Jernite, C. M. Ferrandis, S. Hughes, T. Wolf, A. Guha, L. von Werra, and H. de Vries. "StarCoder: may the source be with you!" In: *CoRR abs/2305.06161* (2023). DOI: [10.48550/ARXIV.2305.06161](https://doi.org/10.48550/ARXIV.2305.06161). arXiv: [2305.06161](https://arxiv.org/abs/2305.06161). URL: <https://doi.org/10.48550/arXiv.2305.06161>.
- [96] R. Li, L. B. Allal, Y. Zi, N. Muennighoff, D. Kocetkov, C. Mou, M. Marone, C. Akiki, J. Li, J. Chim, Q. Liu, E. Zheltonozhskii, T. Y. Zhuo, T. Wang, O. Dehaene, M. Davaadorj, J. Lamy-Poirier, J. Monteiro, O. Shliazhko, N. Gontier, N. Meade, A. Zebaze, M. Yee, L. K. Umapathi, J. Zhu, B. Lipkin, M. Oblokulov, Z. Wang,

- R. M. V, J. Stillerman, S. S. Patel, D. Abulkhanov, M. Zocca, M. Dey, Z. Zhang, N. Moustafa-Fahmy, U. Bhattacharyya, W. Yu, S. Singh, S. Luccioni, P. Villegas, M. Kunakov, F. Zhdanov, M. Romero, T. Lee, N. Timor, J. Ding, C. Schlesinger, H. Schoelkopf, J. Ebert, T. Dao, M. Mishra, A. Gu, J. Robinson, C. J. Anderson, B. Dolan-Gavitt, D. Contractor, S. Reddy, D. Fried, D. Bahdanau, Y. Jernite, C. M. Ferrandis, S. Hughes, T. Wolf, A. Guha, L. von Werra, and H. de Vries. "StarCoder: may the source be with you!" In: *CoRR* abs/2305.06161 (2023). DOI: [10.48550/ARXIV.2305.06161](https://doi.org/10.48550/ARXIV.2305.06161). arXiv: [2305.06161](https://arxiv.org/abs/2305.06161). URL: <https://doi.org/10.48550/arXiv.2305.06161>.
- [97] Y. Li, D. H. Choi, J. Chung, N. Kushman, J. Schrittwieser, R. Leblond, T. Eccles, J. Keeling, F. Gimeno, A. D. Lago, T. Hubert, P. Choy, C. de Masson d'Autume, I. Babuschkin, X. Chen, P. Huang, J. Welbl, S. Goyal, A. Cherepanov, J. Molloy, D. J. Mankowitz, E. S. Robson, P. Kohli, N. de Freitas, K. Kavukcuoglu, and O. Vinyals. "Competition-Level Code Generation with AlphaCode". In: *CoRR* abs/2203.07814 (2022). DOI: [10.48550/ARXIV.2203.07814](https://doi.org/10.48550/ARXIV.2203.07814). arXiv: [2203.07814](https://arxiv.org/abs/2203.07814). URL: <https://doi.org/10.48550/arXiv.2203.07814>.
- [98] 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. URL: <https://openreview.net/forum?id=1qvx610Cu7>.
- [99] 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: *Advances in Neural Information Processing Systems* 36 (2024).
- [100] S. Liu and L. N. Vicente. "Accuracy and fairness trade-offs in machine learning: A stochastic multi-objective approach". In: *Computational Management Science* 19.3 (2022), pp. 513–537.
- [101] T. Liu, C. Xu, and J. J. McAuley. "RepoBench: Benchmarking Repository-Level Code Auto-Completion Systems". In: *CoRR* abs/2306.03091 (2023). DOI: [10.48550/ARXIV.2306.03091](https://doi.org/10.48550/ARXIV.2306.03091). arXiv: [2306.03091](https://arxiv.org/abs/2306.03091). URL: <https://doi.org/10.48550/arXiv.2306.03091>.
- [102] Y. Liu, X. Chen, Y. Gao, Z. Su, F. Zhang, D. Zan, J.-G. Lou, P.-Y. Chen, and T.-Y. Ho. "Uncovering and Quantifying Social Biases in Code Generation". In: *arXiv preprint arXiv:2305.15377* (2023).
- [103] J. Lu, Z. Dou, H. Wang, Z. Cao, J. Dai, Y. Wan, Y. Huang, and Z. Guo. "AutoCV: Empowering Reasoning with Automated Process Labeling via Confidence Variation". In: *CoRR* abs/2405.16802 (2024). DOI: [10.48550/ARXIV.2405.16802](https://doi.org/10.48550/ARXIV.2405.16802). arXiv: [2405.16802](https://arxiv.org/abs/2405.16802). URL: <https://doi.org/10.48550/arXiv.2405.16802>.
- [104] J. Lu, Z. Liu, Y. Wan, Y. Huang, H. Wang, Z. Yang, J. Tang, and Z. Guo. "Process-Driven Autoformalization in Lean 4". In: *CoRR* abs/2406.01940 (2024). DOI: [10.48550/ARXIV.2406.01940](https://doi.org/10.48550/ARXIV.2406.01940). arXiv: [2406.01940](https://arxiv.org/abs/2406.01940). URL: <https://doi.org/10.48550/arXiv.2406.01940>.
- [105] 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".



- In: *ArXiv* abs/2306.08568 (2023). URL: <https://api.semanticscholar.org/CorpusID:259164815>.
- [106] 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". In: *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024. URL: <https://openreview.net/forum?id=UnUwSIgK5W>.
- [107] 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". In: *arXiv preprint arXiv:2306.08568* (2023).
- [108] A. Madaan, N. Tandon, P. Gupta, S. Hallinan, L. Gao, S. Wiegrefe, U. Alon, N. Dziri, S. Prabhume, Y. Yang, S. Gupta, B. P. Majumder, K. Hermann, S. Welleck, A. Yazdanbakhsh, and P. Clark. "Self-Refine: Iterative Refinement with Self-Feedback". In: *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*. Ed. by A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine. 2023. URL: [http://papers.nips.cc/paper%5C\\_files/paper/2023/hash/91edff07232fb1b55a505a9e9f6c0ff3-Abstract-Conference.html](http://papers.nips.cc/paper%5C_files/paper/2023/hash/91edff07232fb1b55a505a9e9f6c0ff3-Abstract-Conference.html).
- [109] A. Madaan and A. Yazdanbakhsh. "Text and patterns: For effective chain of thought, it takes two to tango". In: *arXiv preprint arXiv:2209.07686* (2022).
- [110] J. Mancebo, F. Garcia, and C. Calero. "A process for analysing the energy efficiency of software". In: *Information and Software Technology* 134 (2021), p. 106560.
- [111] MAP. *CodeFeedback-Filtered-Instruction*. <https://huggingface.co/datasets/m-a-p/CodeFeedback-Filtered-Instruction>. 2023.
- [112] N. Mehrabi, F. Morstatter, N. Saxena, K. Lerman, and A. Galstyan. "A survey on bias and fairness in machine learning". In: *ACM computing surveys (CSUR)* 54.6 (2021), pp. 1–35.
- [113] N. Mehrabi, F. Morstatter, N. A. Saxena, K. Lerman, and A. G. Galstyan. "A Survey on Bias and Fairness in Machine Learning". In: *ACM Computing Surveys (CSUR)* 54 (2019), pp. 1–35. URL: <https://api.semanticscholar.org/CorpusID:201666566>.
- [114] Meta. *Introducing Meta Llama 3: The most capable openly available LLM to date*. 2024. URL: <https://ai.meta.com/blog/meta-llama-3/>.
- [115] J. Metcalfe. "Learning from errors". In: *Annual review of psychology* 68 (2017), pp. 465–489.
- [116] Microsoft. *The world's most widely adopted AI developer tool*. 2024. URL: <https://github.com/features/copilot>.
- [117] A. M. Mir, E. Latoskinas, S. Proksch, and G. Gousios. "Type4Py: Practical Deep Similarity Learning-Based Type Inference for Python". In: *44th IEEE/ACM 44th International Conference on Software Engineering, ICSE 2022, Pittsburgh, PA, USA, May 25-27, 2022*. ACM, 2022, pp. 2241–2252. DOI: [10.1145/3510003.3510124](https://doi.org/10.1145/3510003.3510124). URL: <https://doi.org/10.1145/3510003.3510124>.

- [118] C. Mougan, L. State, A. Ferrara, S. Ruggieri, and S. Staab. "Demographic Parity Inspector: Fairness Audits via the Explanation Space". In: *arXiv preprint arXiv:2303.08040* (2023).
- [119] N. Muennighoff, Q. Liu, A. R. Zebaze, Q. Zheng, B. Hui, T. Y. Zhuo, S. Singh, X. Tang, L. von Werra, and S. Longpre. "OctoPack: Instruction Tuning Code Large Language Models". In: *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024. URL: <https://openreview.net/forum?id=mw1PWNSWZP>.
- [120] A. N. Mukherjee, S. Bhattacharyya, and R. Bera. "Role of information technology in human resource management of SME: A study on the use of applicant tracking system". In: *IBMRD's Journal of Management & Research* (2014), pp. 1–22.
- [121] M. Nadeem, A. Bethke, and S. Reddy. "StereoSet: Measuring stereotypical bias in pretrained language models". In: *Annual Meeting of the Association for Computational Linguistics*. 2020.
- [122] G. Nguyen, S. Biswas, and H. Rajan. "Fix Fairness, Don't Ruin Accuracy: Performance Aware Fairness Repair using AutoML". In: *ESEC/FSE'2023: The 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. San Francisco, California, Dec. 2023.
- [123] E. Nijkamp, B. Pang, H. Hayashi, L. Tu, H. Wang, Y. Zhou, S. Savarese, and C. Xiong. "CodeGen: An Open Large Language Model for Code with Multi-Turn Program Synthesis". In: *ICLR* (2023).
- [124] E. Nijkamp, B. Pang, H. Hayashi, L. Tu, H. Wang, Y. Zhou, S. Savarese, and C. Xiong. "CodeGen: An Open Large Language Model for Code with Multi-Turn Program Synthesis". In: *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. URL: [https://openreview.net/pdf?id=iaYcJKpY2B%5C\\_](https://openreview.net/pdf?id=iaYcJKpY2B%5C_).
- [125] C. Niu, T. Zhang, C. Li, B. Luo, and V. Ng. "On Evaluating the Efficiency of Source Code Generated by LLMs". In: *arXiv preprint arXiv:2404.06041* (2024).
- [126] C. Niu, T. Zhang, C. Li, B. Luo, and V. Ng. "On Evaluating the Efficiency of Source Code Generated by LLMs". In: *CoRR abs/2404.06041* (2024). DOI: [10.48550/ARXIV.2404.06041](https://doi.org/10.48550/ARXIV.2404.06041). arXiv: [2404.06041](https://arxiv.org/abs/2404.06041). URL: <https://doi.org/10.48550/arXiv.2404.06041>.
- [127] D. OBrien, S. Biswas, S. Imtiaz, R. Abdalkareem, E. Shihab, and H. Rajan. "Are Prompt Engineering and TODO Comments Friends or Foes? An Evaluation on GitHub Copilot". In: *ICSE'2024: The 46th International Conference on Software Engineering*. Lisbon, Portugal, Apr. 2024.
- [128] A. S. de Oliveira, C. Kaplan, K. Mallat, and T. Chakraborty. "An Empirical Analysis of Fairness Notions under Differential Privacy". In: *arXiv preprint arXiv:2302.02910* (2023).
- [129] OpenAI. *GPT-3.5 Turbo*. 2023. URL: <https://platform.openai.com/docs/models/gpt-3-5>.

- [130] OpenAI. “GPT-4 Technical Report”. In: *CoRR* abs/2303.08774 (2023). DOI: [10.48550/arXiv.2303.08774](https://doi.org/10.48550/arXiv.2303.08774). arXiv: [2303.08774](https://arxiv.org/abs/2303.08774). URL: <https://doi.org/10.48550/arXiv.2303.08774>.
- [131] OpenAI. “GPT-4 Technical Report”. In: *ArXiv* abs/2303.08774 (2023).
- [132] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. L. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray, J. Schulman, J. Hilton, F. Kelton, L. Miller, M. Simens, A. Askell, P. Welinder, P. F. Christiano, J. Leike, and R. Lowe. “Training language models to follow instructions with human feedback”. In: *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*. Ed. by S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh. 2022. URL: [http://papers.nips.cc/paper%5C\\_files/paper/2022/hash/b1efde53be364a73914f58805a001731-Abstract-Conference.html](http://papers.nips.cc/paper%5C_files/paper/2022/hash/b1efde53be364a73914f58805a001731-Abstract-Conference.html).
- [133] 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”. In: *arXiv preprint arXiv:2308.02828* (2023).
- [134] A. Papadaki, N. Martinez, M. Bertran, G. Sapiro, and M. Rodrigues. “Minimax demographic group fairness in federated learning”. In: *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*. 2022, pp. 142–159.
- [135] A. Papadaki, N. Martinez, M. A. Bertran, G. Sapiro, and M. R. Rodrigues. “Federated Fairness without Access to Demographics”. In: *Workshop on Federated Learning: Recent Advances and New Challenges (in Conjunction with NeurIPS 2022)*. 2022.
- [136] S. G. Patil, T. Zhang, X. Wang, and J. E. Gonzalez. “Gorilla: Large Language Model Connected with Massive APIs”. In: *CoRR* abs/2305.15334 (2023). DOI: [10.48550/ARXIV.2305.15334](https://doi.org/10.48550/ARXIV.2305.15334). arXiv: [2305.15334](https://arxiv.org/abs/2305.15334). URL: <https://doi.org/10.48550/arXiv.2305.15334>.
- [137] R. Pereira, M. Couto, F. Ribeiro, R. Rua, J. Cunha, J. P. Fernandes, and J. Saraiva. “Ranking programming languages by energy efficiency”. In: *Science of Computer Programming* 205 (2021), p. 102609.
- [138] J.-P. Platteau and D. U. Ontiveros. “Cognitive bias in insurance: evidence from a health scheme in India”. In: *World Development* 144 (2021), p. 105498.
- [139] R. Qiu, W. W. Zeng, H. Tong, J. Ezick, and C. Lott. “How Efficient is LLM-Generated Code? A Rigorous & High-Standard Benchmark”. In: *arXiv preprint arXiv:2406.06647* (2024).
- [140] R. Rafailov, A. Sharma, E. Mitchell, C. D. Manning, S. Ermon, and C. Finn. “Direct preference optimization: Your language model is secretly a reward model”. In: *Advances in Neural Information Processing Systems* 36 (2024).
- [141] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu. “Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer”. In: *Journal of Machine Learning Research* 21.140 (2020), pp. 1–67. URL: <http://jmlr.org/papers/v21/20-074.html>.
- [142] B. Rozière, J. Gehring, F. Gloeckle, S. Sootla, I. Gat, X. Tan, Y. Adi, J. Liu, T. Remez, J. Rapin, A. Kozhevnikov, I. Evtimov, J. Bitton, M. P. Bhatt, C. C. Fer-



- rer, A. Grattafiori, W. Xiong, A. D'efossez, J. Copet, F. Azhar, H. Touvron, L. Martin, N. Usunier, T. Scialom, and G. Synnaeve. "Code Llama: Open Foundation Models for Code". In: *ArXiv abs/2308.12950* (2023). URL: <https://api.semanticscholar.org/CorpusID:261100919>.
- [143] B. Rozière, J. Gehring, F. Gloeckle, S. Sootla, I. Gat, X. E. Tan, Y. Adi, J. Liu, T. Remez, J. Rapin, A. Kozhevnikov, I. Evtimov, J. Bitton, M. Bhatt, C. Canton-Ferrer, A. Grattafiori, W. Xiong, A. Défossez, J. Copet, F. Azhar, H. Touvron, L. Martin, N. Usunier, T. Scialom, and G. Synnaeve. "Code Llama: Open Foundation Models for Code". In: *CoRR abs/2308.12950* (2023). DOI: [10.48550/ARXIV.2308.12950](https://doi.org/10.48550/ARXIV.2308.12950). arXiv: [2308.12950](https://arxiv.org/abs/2308.12950). URL: <https://doi.org/10.48550/arXiv.2308.12950>.
- [144] B. Rozière, M. Lachaux, L. Chausson, and G. Lample. "Unsupervised Translation of Programming Languages". In: *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*. Ed. by H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin. 2020. URL: <https://proceedings.neurips.cc/paper/2020/hash/ed23fbf18c2cd35f8c7f8de44f85c08d-Abstract.html>.
- [145] L. Salewski, S. Alaniz, I. Rio-Torto, E. Schulz, and Z. Akata. "In-Context Impersonation Reveals Large Language Models' Strengths and Biases". In: *arXiv preprint arXiv:2305.14930* (2023).
- [146] P. Sattigeri, S. Ghosh, I. Padhi, P. Dognin, and K. R. Varshney. "Fair infinitesimal jackknife: Mitigating the influence of biased training data points without refitting". In: *Advances in Neural Information Processing Systems 35* (2022), pp. 35894–35906.
- [147] J. Shi, Z. Yang, and D. Lo. "Efficient and Green Large Language Models for Software Engineering: Vision and the Road Ahead". In: *arXiv preprint arXiv:2404.04566* (2024).
- [148] J. Shi, Z. Yang, and D. Lo. "Efficient and Green Large Language Models for Software Engineering: Vision and the Road Ahead". In: *CoRR abs/2404.04566* (2024). DOI: [10.48550/ARXIV.2404.04566](https://doi.org/10.48550/ARXIV.2404.04566). arXiv: [2404.04566](https://arxiv.org/abs/2404.04566). URL: <https://doi.org/10.48550/arXiv.2404.04566>.
- [149] N. Shinn, F. Cassano, A. Gopinath, K. Narasimhan, and S. Yao. "Reflexion: Language agents with verbal reinforcement learning". In: *Advances in Neural Information Processing Systems 36* (2024).
- [150] N. Shinn, F. Cassano, A. Gopinath, K. Narasimhan, and S. Yao. "Reflexion: language agents with verbal reinforcement learning". In: *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*. Ed. by A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine. 2023. URL: [http://papers.nips.cc/paper%5C\\_files/paper/2023/hash/1b44b878bb782e6954cd888628510e90-Abstract-Conference.html](http://papers.nips.cc/paper%5C_files/paper/2023/hash/1b44b878bb782e6954cd888628510e90-Abstract-Conference.html).
- [151] D. Shrivastava, H. Larochelle, and D. Tarlow. "Repository-Level Prompt Generation for Large Language Models of Code". In: *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*. Ed. by A. Krause,

- E. Brunskill, K. Cho, B. Engelhardt, S. Sabato, and J. Scarlett. Vol. 202. *Proceedings of Machine Learning Research*. PMLR, 2023, pp. 31693–31715. URL: <https://proceedings.mlr.press/v202/shrivastava23a.html>.
- [152] R. K. Shyamasundar. “Introduction to algorithms”. In: *Resonance* 1 (1996), pp. 14–24. URL: <https://api.semanticscholar.org/CorpusID:123556377>.
- [153] A. Shypula, A. Madaan, Y. Zeng, U. Alon, J. Gardner, M. Hashemi, G. Neubig, P. Ranganathan, O. Bastani, and A. Yazdanbakhsh. “**Learning Performance-Improving Code Edits**”. In: *The Twelfth International Conference on Learning Representations (ICLR)*. 2024.
- [154] J. Simson, F. Pfisterer, and C. Kern. “Using Multiverse Analysis to Evaluate the Influence of Model Design Decisions on Algorithmic Fairness”. In: *HAI 2023: Augmenting Human Intellect*. IOS Press, 2023, pp. 382–384.
- [155] L. L. Taylor, J. N. Lahey, M. I. Beck, and J. E. Froyd. “How to do a salary equity study: With an illustrative example from higher education”. In: *Public personnel management* 49.1 (2020), pp. 57–82.
- [156] X.-L. Team. *Xwin-LM*. Version pre-release. Sept. 2023. URL: <https://github.com/Xwin-LM/Xwin-LM>.
- [157] H. Thakur, A. Jain, P. Vaddamanu, P. P. Liang, and L.-P. Morency. “Language Models Get a Gender Makeover: Mitigating Gender Bias with Few-Shot Data Interventions”. In: *arXiv preprint arXiv:2306.04597* (2023).
- [158] S. Tizpaz-Niari, A. Kumar, G. Tan, and A. Trivedi. “Fairness-aware configuration of machine learning libraries”. In: *Proceedings of the 44th International Conference on Software Engineering*. 2022, pp. 909–920.
- [159] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale, D. Bikel, L. Blecher, C. Canton-Ferrer, M. Chen, G. Cucurull, D. Esiobu, J. Fernandes, J. Fu, W. Fu, B. Fuller, C. Gao, V. Goswami, N. Goyal, A. Hartshorn, S. Hosseini, R. Hou, H. Inan, M. Kardaş, V. Kerkez, M. Khabsa, I. Kloumann, A. Korenev, P. S. Koura, M. Lachaux, T. Lavril, J. Lee, D. Liskovich, Y. Lu, Y. Mao, X. Martinet, T. Mihaylov, P. Mishra, I. Molybog, Y. Nie, A. Poulton, J. Reizenstein, R. Rungta, K. Saladi, A. Schelten, R. Silva, E. M. Smith, R. Subramanian, X. E. Tan, B. Tang, R. Taylor, A. Williams, J. X. Kuan, P. Xu, Z. Yan, I. Zarov, Y. Zhang, A. Fan, M. Kambadur, S. Narang, A. Rodriguez, R. Stojnic, S. Edunov, and T. Scialom. “Llama 2: Open Foundation and Fine-Tuned Chat Models”. In: *CoRR abs/2307.09288* (2023). DOI: [10.48550/ARXIV.2307.09288](https://doi.org/10.48550/ARXIV.2307.09288). arXiv: [2307.09288](https://arxiv.org/abs/2307.09288). URL: <https://doi.org/10.48550/arXiv.2307.09288>.
- [160] L. Tunstall, N. Reimers, U. E. S. Jo, L. Bates, D. Korat, M. Wasserblat, and O. Pereg. “Efficient few-shot learning without prompts”. In: *arXiv preprint arXiv:2209.11055* (2022).
- [161] I. UIUC. *Magicoder-Evol-Instruct-110K*. <https://huggingface.co/datasets/ise-uiuc/Magicoder-Evol-Instruct-110K>. 2023.
- [162] I. UIUC. *Magicoder-OSS-Instruct-75K*. <https://huggingface.co/datasets/ise-uiuc/Magicoder-OSS-Instruct-75K>. 2023.

- [163] E. L. Ungless, A. Rafferty, H. Nag, and B. Ross. "A Robust Bias Mitigation procedure based on the stereotype content model". In: *arXiv preprint arXiv:2210.14552* (2022).
- [164] *US Health Insurance Dataset*. [www.kaggle.com/datasets/teertha/ushealthinsurancedataset](https://www.kaggle.com/datasets/teertha/ushealthinsurancedataset). Accessed on August 1, 2023. 2023.
- [165] S. B. Usman Gohar and H. Rajan. "Towards Understanding Fairness and its Composition in Ensemble Machine Learning". In: *ICSE'2023: The 45th International Conference on Software Engineering*. Melbourne, Australia, May 2023.
- [166] T. Vartziotis, I. Dellatolas, G. Dasoulas, M. Schmidt, F. Schneider, T. Hoffmann, S. Kotsopoulos, and M. Keckeisen. "Learn to Code Sustainably: An Empirical Study on LLM-based Green Code Generation". In: *arXiv preprint arXiv:2403.03344* (2024).
- [167] Vezora. *Tested-143k-Python-Alpaca*. <https://huggingface.co/datasets/Vezora/Tested-143k-Python-Alpaca>. 2023.
- [168] S. Waghjale, V. Veerendranath, Z. Z. Wang, and D. Fried. "ECCO: Can We Improve Model-Generated Code Efficiency Without Sacrificing Functional Correctness?". In: *arXiv preprint arXiv:2407.14044* (2024).
- [169] A. Wang, V. V. Ramaswamy, and O. Russakovsky. "Towards intersectionality in machine learning: Including more identities, handling underrepresentation, and performing evaluation". In: *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*. 2022, pp. 336–349.
- [170] P. Wang, L. Li, L. Chen, D. Zhu, B. Lin, Y. Cao, Q. Liu, T. Liu, and Z. Sui. "Large language models are not fair evaluators". In: *arXiv preprint arXiv:2305.17926* (2023).
- [171] S. Wang, Z. Li, H. Qian, C. Yang, Z. Wang, M. Shang, V. Kumar, S. Tan, B. Ray, P. Bhatia, R. Nallapati, M. K. Ramanathan, D. Roth, and B. Xiang. "ReCode: Robustness Evaluation of Code Generation Models". In: *ArXiv abs/2212.10264* (2022). URL: <https://api.semanticscholar.org/CorpusID:254877229>.
- [172] S. Wang, Z. Li, H. Qian, C. Yang, Z. Wang, M. Shang, V. Kumar, S. Tan, B. Ray, P. Bhatia, R. Nallapati, M. K. Ramanathan, D. Roth, and B. Xiang. "ReCode: Robustness Evaluation of Code Generation Models". In: *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023*. Ed. by A. Rogers, J. L. Boyd-Graber, and N. Okazaki. Association for Computational Linguistics, 2023, pp. 13818–13843. DOI: [10.18653/v1/2023.acl-long.773](https://doi.org/10.18653/v1/2023.acl-long.773). URL: <https://doi.org/10.18653/v1/2023.acl-long.773>.
- [173] X. Wang, J. Wei, D. Schuurmans, Q. Le, E. Chi, S. Narang, A. Chowdhery, and D. Zhou. "Self-consistency improves chain of thought reasoning in language models". In: *arXiv preprint arXiv:2203.11171* (2022).
- [174] Y. Wang, Y. Kordi, S. Mishra, A. Liu, N. A. Smith, D. Khashabi, and H. Hajishirzi. "Self-Instruct: Aligning Language Models with Self-Generated Instructions". In: *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023*. Ed.

- by A. Rogers, J. L. Boyd-Graber, and N. Okazaki. Association for Computational Linguistics, 2023, pp. 13484–13508. DOI: [10.18653/v1/2023.acl-long.754](https://doi.org/10.18653/v1/2023.acl-long.754). URL: <https://doi.org/10.18653/v1/2023.acl-long.754>.
- [175] Y. Wang, H. Le, A. D. Gotmare, N. D. Bui, J. Li, and S. C. Hoi. “Codet5+: Open code large language models for code understanding and generation”. In: *arXiv preprint arXiv:2305.07922* (2023).
- [176] Z. Wang, S. Zhou, D. Fried, and G. Neubig. “Execution-Based Evaluation for Open-Domain Code Generation”. In: *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023*. Ed. by H. Bouamor, J. Pino, and K. Bali. Association for Computational Linguistics, 2023, pp. 1271–1290. URL: <https://aclanthology.org/2023.findings-emnlp.89>.
- [177] J. Wei, M. Bosma, V. Y. Zhao, K. Guu, A. W. Yu, B. Lester, N. Du, A. M. Dai, and Q. V. Le. “Finetuned Language Models are Zero-Shot Learners”. In: *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net, 2022. URL: <https://openreview.net/forum?id=gEZrGCozdqR>.
- [178] J. Wei, X. Wang, D. Schuurmans, M. Bosma, E. H.-h. Chi, F. Xia, Q. Le, and D. Zhou. “Chain of Thought Prompting Elicits Reasoning in Large Language Models”. In: *ArXiv abs/2201.11903* (2022). URL: <https://api.semanticscholar.org/CorpusID:246411621>.
- [179] J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, D. Zhou, et al. “Chain-of-thought prompting elicits reasoning in large language models”. In: *Advances in Neural Information Processing Systems* 35 (2022), pp. 24824–24837.
- [180] J. Wei, G. Durrett, and I. Dillig. “TypeT5: Seq2seq Type Inference using Static Analysis”. In: *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. URL: <https://openreview.net/pdf?id=4TyNEhI2GdN>.
- [181] Y. Wei, Z. Wang, J. Liu, Y. Ding, and L. Zhang. “Magicoder: Empowering Code Generation with OSS-Instruct”. In: *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net, 2024. URL: <https://openreview.net/forum?id=XUeo0Bid3x>.
- [182] F. Xia, T. Guo, X. Bai, A. Shatte, Z. Liu, and J. Tang. “SUMMER: Bias-aware Prediction of Graduate Employment Based on Educational Big Data”. In: *ACM/IMS Transactions on Data Science (TDS)* 2.4 (2022), pp. 1–24.
- [183] C. Xu, Q. Sun, K. Zheng, X. Geng, P. Zhao, J. Feng, C. Tao, Q. Lin, and D. Jiang. “WizardLM: Empowering Large Pre-Trained Language Models to Follow Complex Instructions”. In: *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024. URL: <https://openreview.net/forum?id=CfXh93NDgH>.
- [184] Y. Yao, H. Wu, Z. Guo, B. Zhou, J. Gao, S. Luo, H. Hou, X. Fu, and L. Song. *Learning From Correctness Without Prompting Makes LLM Efficient Reasoner*. 2024. arXiv: [2403.19094](https://arxiv.org/abs/2403.19094) [cs.CL].

- [185] P. Yin, W. Li, K. Xiao, A. Rao, Y. Wen, K. Shi, J. Howland, P. Bailey, M. Catasta, H. Michalewski, O. Polozov, and C. Sutton. “Natural Language to Code Generation in Interactive Data Science Notebooks”. In: *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, ACL 2023, Toronto, Canada, July 9-14, 2023. Ed. by A. Rogers, J. L. Boyd-Graber, and N. Okazaki. Association for Computational Linguistics, 2023, pp. 126–173. DOI: [10.18653/v1/2023.ACL-LONG.9](https://doi.org/10.18653/v1/2023.ACL-LONG.9). URL: <https://doi.org/10.18653/v1/2023.acl-long.9>.
- [186] W. Yu, Z. Zhang, Z. Liang, M. Jiang, and A. Sabharwal. “Improving Language Models via Plug-and-Play Retrieval Feedback”. In: *CoRR abs/2305.14002* (2023). DOI: [10.48550/ARXIV.2305.14002](https://doi.org/10.48550/ARXIV.2305.14002). arXiv: [2305.14002](https://arxiv.org/abs/2305.14002). URL: <https://doi.org/10.48550/arXiv.2305.14002>.
- [187] Y. Yu, Y. Zhuang, J. Zhang, Y. Meng, A. Ratner, R. Krishna, J. Shen, and C. Zhang. “Large Language Model as Attributed Training Data Generator: A Tale of Diversity and Bias”. In: *arXiv preprint arXiv:2306.15895* (2023).
- [188] D. Zan, B. Chen, D. Yang, Z. Lin, M. Kim, B. Guan, Y. Wang, W. Chen, and J. Lou. “CERT: Continual Pre-training on Sketches for Library-oriented Code Generation”. In: *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI 2022, Vienna, Austria, 23-29 July 2022*. Ed. by L. D. Raedt. ijcai.org, 2022, pp. 2369–2375. DOI: [10.24963/IJCAI.2022/329](https://doi.org/10.24963/IJCAI.2022/329). URL: <https://doi.org/10.24963/ijcai.2022/329>.
- [189] F. Zhang, B. Chen, Y. Zhang, J. Keung, J. Liu, D. Zan, Y. Mao, J. Lou, and W. Chen. “RepoCoder: Repository-Level Code Completion Through Iterative Retrieval and Generation”. In: *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*. Ed. by H. Bouamor, J. Pino, and K. Bali. Association for Computational Linguistics, 2023, pp. 2471–2484. URL: <https://aclanthology.org/2023.emnlp-main.151>.
- [190] K. Zhang, Z. Li, J. Li, G. Li, and Z. Jin. “Self-Edit: Fault-Aware Code Editor for Code Generation”. In: *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, ACL 2023, Toronto, Canada, July 9-14, 2023. Ed. by A. Rogers, J. L. Boyd-Graber, and N. Okazaki. Association for Computational Linguistics, 2023, pp. 769–787. DOI: [10.18653/v1/2023.ACL-LONG.45](https://doi.org/10.18653/v1/2023.ACL-LONG.45). URL: <https://doi.org/10.18653/v1/2023.acl-long.45>.
- [191] K. Zhang, D. Wang, J. Xia, W. Y. Wang, and L. Li. “ALGO: Synthesizing Algorithmic Programs with Generated Oracle Verifiers”. In: *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*. Ed. by A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine. 2023. URL: [http://papers.nips.cc/paper%5C\\_files/paper/2023/hash/abe1eb21ceb046209c96a0f5e7544ccc-Abstract-Conference.html](http://papers.nips.cc/paper%5C_files/paper/2023/hash/abe1eb21ceb046209c96a0f5e7544ccc-Abstract-Conference.html).
- [192] Q. Zheng, X. Xia, X. Zou, Y. Dong, S. Wang, Y. Xue, Z. Wang, L. Shen, A. Wang, Y. Li, T. Su, Z. Yang, and J. Tang. “CodeGeeX: A Pre-Trained Model for Code Generation with Multilingual Evaluations on HumanEval-X”. In: *CoRR abs/2303.17568*

- (2023). DOI: [10.48550/ARXIV.2303.17568](https://doi.org/10.48550/ARXIV.2303.17568). arXiv: 2303.17568. URL: <https://doi.org/10.48550/arXiv.2303.17568>.
- [193] T. Zheng, G. Zhang, T. Shen, X. Liu, B. Y. Lin, J. Fu, W. Chen, and X. Yue. “Open-CodeInterpreter: Integrating Code Generation with Execution and Refinement”. In: *arXiv preprint arXiv:2402.14658* (2024).
- [194] Y. Zheng, R. Zhang, J. Zhang, Y. Ye, Z. Luo, Z. Feng, and Y. Ma. “LlamaFactory: Unified Efficient Fine-Tuning of 100+ Language Models”. In: *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*. Bangkok, Thailand: Association for Computational Linguistics, 2024. URL: <http://arxiv.org/abs/2403.13372>.
- [195] Q. Zhu, D. Guo, Z. Shao, D. Yang, P. Wang, R. Xu, Y. Wu, Y. Li, H. Gao, S. Ma, et al. “DeepSeek-Coder-V2: Breaking the Barrier of Closed-Source Models in Code Intelligence”. In: *arXiv preprint arXiv:2406.11931* (2024).
- [196] S. F. Zimin Chen and M. Monperrus. “Supersonic: Learning to Generate Source Code Optimizations in C/C++”. In: *arXiv preprint arXiv:2309.14846* (2023).