

Code Comparison Tuning for Code Large Language Models

Yufan Jiang¹, Qiaozhi He¹, Xiaomin Zhuang¹, Zhihua Wu¹

¹National Supercomputing Center in Wuxi

jiangyufan2018@outlook.com

Abstract

We present Code Comparison Tuning (CCT), a simple and effective tuning method for code large language models (Code LLMs) to better handle subtle code errors. Specifically, we integrate the concept of comparison into instruction tuning, both at the token and sequence levels, enabling the model to discern even the slightest deviations in code. To compare the original code with an erroneous version containing manually added code errors, we use token-level preference loss for detailed token-level comparisons. Additionally, we combine code segments to create a new instruction tuning sample for sequence-level comparisons, enhancing the model’s bug-fixing capability. Experimental results on the HumanEvalFix benchmark show that CCT surpasses instruction tuning in pass@1 scores by up to 4 points across diverse code LLMs, and extensive analysis demonstrates the effectiveness of our method.

1 Introduction

Fixing bugs with neural models has become popular among programmers for its powerful capabilities. The earliest of these approaches typically consist of multiple individual stages, such as the detection stage and generation stage (Lutellier et al., 2020; Allamanis et al., 2021; Yasunaga and Liang, 2021; Mashhadi and Hemmati, 2021; Bui et al., 2022), whereas Code LLMs successfully address the problem with a simple instruction “Fix the bugs in the code” and achieve competitive performance.

Closed-source LLMs like GPT-4 (OpenAI, 2023) have already shown promising results in these code-related tasks. However, due to high API fees and security problems, exploring how to achieve similar performance using open-source Code LLMs has become a highly meaningful research direction that we focus on in this work. To ensure responsiveness to human requests, open-source Code LLMs usually undergo a two-step process. First, they are

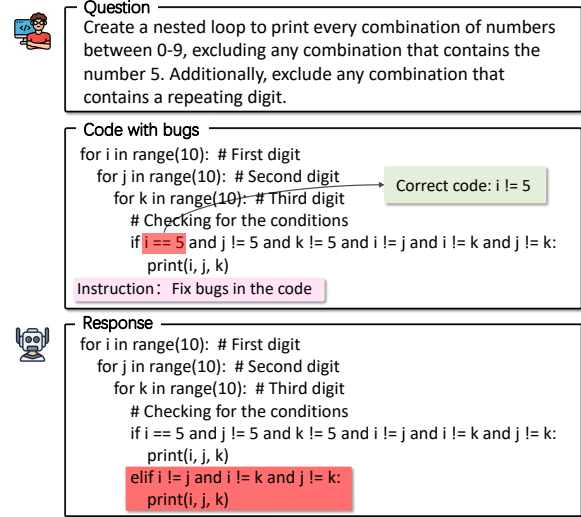


Figure 1: An erroneous bug fix example. Given the code-related issues, users or code language models generate code with bugs. The fine-tuned models tend to introduce additional errors when attempting to fix bugs (red).

pre-trained on extensive raw code data, enabling them to acquire a foundational understanding of code patterns and structures (Nijkamp et al., 2022; Fried et al., 2022; Li et al., 2023; Roziere et al., 2023; Di et al., 2023). Following pre-training, instruction tuning (Wei et al., 2021; Ouyang et al., 2022) is employed to align Code LLMs with specific code task instructions provided by humans, such as code completion, bug fixing, or code interpretation (Luo et al., 2023; Shen et al., 2023; Wang et al., 2023b).

To further enhance the bug-fixing capabilities of open-source Code LLMs, some approaches construct specific code-fixing datasets, aiming to bridge the gap between instruction tuning and actual bug fixing (Zhang et al., 2023; Muennighoff et al., 2023). Other approaches attempt to integrate code interpreters into the Code LLMs in the form of APIs, enabling real-time code inspection (Wang et al., 2023a; Bai et al., 2023; Gou et al.,

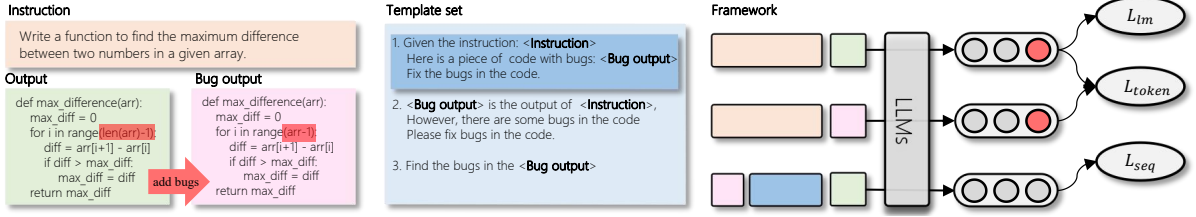


Figure 2: The overall framework of our proposed CCT.

2023; Chen et al., 2023). While these solutions have demonstrated effectiveness in practice, teaching Code LLMs to fix bugs remains a challenge. Constructing datasets necessitates careful design and collection, making it impractical to cover all error types. Furthermore, the fine-tuned code models have been proven ineffective in dealing with small changes in the codes (Muennighoff et al., 2023). When instructed to fix bugs in codes, the models often regenerate the erroneous code or introduce new bugs. Take the code in Figure 1 as an example. Additionally, while code interpreters can assist in identifying syntactic errors, they are unable to detect logical errors within the code.

Here, we present a simple and effective tuning method, namely Code Comparison Tuning (CCT). This method is specifically designed to heighten the sensitivity of Code LLMs to nuanced variations in code structures. Central to CCT is the integration of a *comparison* mechanism into instruction tuning, realized by creating erroneous versions of each instructive code example. These versions undergo token-level comparative analysis, significantly improving the model’s ability to discern and differentiate erroneous code. Additionally, the training dataset is augmented by pairing these generated erroneous codes with their correct forms, as demonstrated in constructs like “Fixing the error in A results in B”, which further enhances the model’s capability of fixing bugs. Experiments and analysis conducted on the HumanEvalFix benchmark well validate the effectiveness of CCT. Specifically, we observe a substantial improvement over 4 points in pass@1 scores compared to standard instruction tuning on different backbones.

2 Method

To make code generative models more sensitive to the errors in the code, we incorporate two levels of code comparison (token-level and sequence-level) into the instruction tuning of code pre-trained models. We first give a brief introduction to instruction

tuning. Then, we introduce two kinds of code comparisons in detail. While we conducted research on Python in this paper, our approach can be applied to any programming language.

2.1 Background: Instruction Tuning

The goal of instruction tuning is to improve the capability of language models in effectively processing instructions expressed in natural languages. In general, each instance of instruction-following data begins with "instructions" denoted as c , which describes a task, accompanied by a corresponding output y that represents the answer to the given instruction. The “input” x , is the optional context or input for the task. Given the instruction data, the language models are optimized by minimizing the negative log-likelihood of the output y :

$$\mathcal{L}_{lm} = -\frac{1}{|y|} \sum_i \log p(y_i | c, x), \quad (1)$$

2.2 Code Comparison Tuning

We propose two kinds of code comparisons from different perspectives to improve the model’s ability to handle error codes. Specifically, for the code block t in the output y , we obtain its counterpart t' by introducing code errors manually. Then, we perform comparisons between t and t' at both token level and sequence level.

To construct code containing bugs, we initially extract code blocks from the output y . Subsequently, we introduce bugs by either randomly replacing or deleting elements such as variables, functions, and operators within these code segments. Bug examples are in Appendix A. With examples of the correct code and error code, the model is optimized to locate the bugs and fix them.

Token-level Comparison Previous studies usually provide supervision signals to code language models by training samples in the format of bug fixes, in order to guide models on how to repair bugs. However, this type of sequence-based

training sample causes the model to ignore more granular-level differences between code snippets, which results in the degeneration of the model’s ability to repair errors. To tackle this problem, we adopt a token-level comparison loss (Zeng et al., 2023) to teach models to be aware of the changes in each token.

Formally, given code t and its counterpart t' , the token-level comparison loss is defined as:

$$\mathcal{L}_{token} = -\frac{1}{M-I} \sum_{i=I}^N \max(0, -r_{\theta}(h_i^t) + r_{\theta}(h_i^{t'}) + 1.0), \quad (2)$$

where I represents the index starting from the first differing segment between sequences t and t' , and M is the maximum length of two sequences. The hidden state of each token i is denoted as h_i and we add a linear head r_{θ} that converts the hidden state to a scalar.

Sequence-level Comparison Beyond mastering token-level distinctions, our approach integrates both t and t' within a single sentence, facilitating the model’s acquisition of sequence-level repair skills. Specifically, we first create a set of templates designed to transform comparative code pairs into coherent instructional data. Then we convert the code pairs to instruction-tuning style by randomly sampling a template from T . All the templates are illustrated in Appendix B. Finally, the sequence-level comparison example is used to optimize the language model via Eq.(1) with the associated loss denoted as \mathcal{L}_{seq} .

2.3 Overall Training Objective

The overall training objective is defined as:

$$\mathcal{L} = \mathcal{L}_{lm} + \alpha * \mathcal{L}_{token} + \beta * \mathcal{L}_{seq}, \quad (3)$$

where α and β are non-negative hyper-parameters to balance the effect of each loss term. In this paper, we set α and β to 2.0 and 0.5, respectively

3 Experiments

3.1 Datasets

We conducted experiments on Evol-Instruct-Code-80k dataset¹ licensed by Apache-2.0. The dataset is created following the process described in the WizardCoder Paper (Luo et al., 2023). We extracted data from code written in Python to use as

Model	Params	Pass@1
<i>Closed-source LLMs</i>		
ChatGPT	-	39.6
GPT-4	$\geq 175B$	47.0
<i>Open-source LLMs</i>		
InstructCodeT5+*	16B	2.7
BLOOMZ*	176B	16.6
StarCoder*	15.5B	8.7
CodeLlama*	13B	15.2
CodeGeeX2*	6B	15.9
OctoCoder*	15.5B	30.4
WizardCoder	15.5B	31.8
WizardCoder-Python-13B	13B	47.7
<i>StarCoder backbone</i>		
Instruct tuning	15.5B	33.7
CCT-StarCoder (Ours)	15.5B	38.3
<i>CodeLlama-Python-13B backbone</i>		
Instruct tuning	13B	43.5
CCT-CodeLlama (Ours)	13B	47.7

Table 1: **Pass@1 (%) performance on HumanEval-Fix.** Models with * denote that we directly report the scores from the corresponding paper

our instruction data. To verify the effectiveness of our proposed approach, we evaluated CCT on the HumanEvalFix (Muennighoff et al., 2023) which is proposed to task models to fix the bugs in function. It contains 164 HumanEval solutions across all 6 languages (984 total bugs) and the errors are manually inserted into the code.

3.2 Baselines & Settings

We mainly experimented on CodeLlama-13b-Python (Roziere et al., 2023) and StarCoder (Li et al., 2023) in this work. Additionally, we report the results of InstructCodeT5+, BLOOMZ, CodeGeeX2, StarCoder, OctoCoder and WizardCoder-Python-13B (Muennighoff et al., 2022; Wang et al., 2023b; Li et al., 2023; Zheng et al., 2023; Luo et al., 2023). We also report the results of closed-source models such as ChatGPT and GPT-4 which can be accessed via API.

To facilitate a fair and consistent evaluation, we fine-tuned all models for 1 epoch with a batch size of 64. The learning rate was set to $2e-5$ and the weight decay parameter was set to 0.0. For evaluation, we used the pass@1 metric (Chen et al., 2021). Similar to Muennighoff et al. (2023), we used a sampling temperature of 0.2 and $top_p = 0.95$ to estimate pass@1. We generated $n = 20$ samples, which is enough to get reliable pass@1 estimates (Li et al., 2023).

¹<https://github.com/nickrosh/evol-teacher>

Model	Pass@1
Instruct tuning	43.53±0.49
w Sequence-level data	45.76±0.23
CCT-CodeLlama	47.71±0.39
w/o \mathcal{L}_{seq}	44.56±0.61
w/o \mathcal{L}_{token}	45.83±0.32

Table 2: **Ablation study.** We run each experiment 3 times with different random seeds and report mean and standard deviation .

3.3 Results

Table 1 shows the results of several models on HumanEvalFix. We can see that most open-source code LLMs struggle with handling subtle code changes and instructing tuning can substantially enhance their performance. Our Code comparison tuning significantly outperforms instruct tuning on both StarCoder and CodeLlama-Python-13B backbone, leading to an average of 4 Pass@1 scores improvement

At the same time, CCT achieves comparable results to its open-source competitors of the same size. These results demonstrate the effectiveness of the proposed code comparison method. Although CCT has surpassed GPT4 on HumanEvalFix, we still need to conduct further testing for evaluation. We leave this issue for future study.

3.4 Ablation Study

To analyze the impact of different components of CCT, we investigate the following variants: 1) *CCT* w/o \mathcal{L}_{seq} , removing the sequence-level comparison; 2) *CCT* w/o \mathcal{L}_{token} , removing the token-level comparison; Additionally, we utilize the data generated from the sequence-level comparison phase to create instruction fine-tuning data and mix it together with original data to fine-tune the model which denotes as *w Sequence-level data*. We take CodeLlama-Python-13B as the backbone.

The results are listed in Table 2. The degradation of *CCT* w/o \mathcal{L}_{seq} and *CCT* w/o \mathcal{L}_{token} indicate that code LLMs can improve their ability to learn how to fix errors in code by leveraging the code comparison in both token and sequence levels. While *w Sequence-level data* performs better than the standard instruct tuning, there is still room for improvement as our CCT achieved even better results. This suggests that our proposed method goes beyond just data augmentation, as it incorporates comparison during fine-tuning to enhance the effectiveness and efficiency of code LLMs.

Model	Pass@1
<i>Closed-source LLMs</i>	
GPT-4	88.4
<i>Open-source LLMs</i>	
WizardCoder-Python-13B	60.37
Instruct tuning	63.26
CCT-CodeLlama	66.1

Table 3: **Pass@1 (%) performance on HumanEval-FixDocs.**

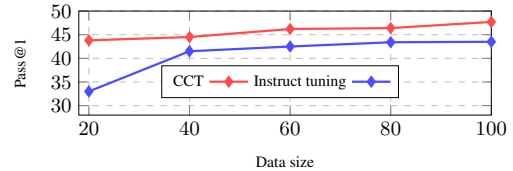


Figure 3: **Effect of Instruction dataset size.** We report pass@1 under different sizes of instructing datasets.

3.5 Results on HumanEvalFixDocs

HumanEvalFixDocs (Muennighoff et al., 2023) provides docstrings as the source of ground truth for the model to fix the buggy function which is generally easier for models than HumanEvalFix. From Table 3, we see that our CCT performs significantly better than instruction fine-tuning and other open-source models. However, it also reveals a notable performance gap compared with GPT4, an aspect we aim to explore in our future research.

3.6 Effect of Corpus Size

In this experiment, we study the impact of data sizes on CCT by sampling different percentages of the instructing dataset. Figure 3 shows the comparison between our CCT and instruct tuning under different data sizes. We see that, when the amount of data used gradually decreases, our CCT still maintains a strong performance. Surprisingly, with only 20% of the data, CCT can achieve a pass@1 score of 43, demonstrating the data efficiency of our proposed method.

4 Conclusions

In this work, we enhance the ability of code LLMs to fix bugs by integrating code comparison during instruct tuning. We consider both token-level and sequence-level comparisons to make code models more sensitive to the small changes in the code. Experiments and analyses validate the effectiveness of our model. We plan to extend our method to more programming languages and conduct tests on a wider range of test sets in our future study.

Limitations

There are still a few drawbacks of our approach that need further investigation. The construction method we use for generating error code snippets is relatively simple. Introducing more complex construction methods is necessary to provide the model with additional comparative information. Second, more bug-fixing tests are needed, including a wider range of programming languages and a greater variety of error types. We leave these investigations for future work. While we have achieved remarkable results in the evaluation metrics of the code repair task, there is still an ongoing need for continuous research and dedicated efforts to enhance how code-pretrained models can better assist programmers in handling code-related tasks.

References

- Miltiadis Allamanis, Henry Jackson-Flux, and Marc Brockschmidt. 2021. Self-supervised bug detection and repair. *Advances in Neural Information Processing Systems*, 34:27865–27876.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. 2023. Qwen technical report. *arXiv preprint arXiv:2309.16609*.
- Nghi DQ Bui, Yue Wang, and Steven Hoi. 2022. Detect-localize-repair: A unified framework for learning to debug with codet5. *arXiv preprint arXiv:2211.14875*.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. Evaluating large language models trained on code.(2021). *arXiv preprint arXiv:2107.03374*.
- Xinyun Chen, Maxwell Lin, Nathanael Schärli, and Denny Zhou. 2023. Teaching large language models to self-debug. *arXiv preprint arXiv:2304.05128*.
- Peng Di, Jianguo Li, Hang Yu, Wei Jiang, Wenting Cai, Yang Cao, Chaoyu Chen, Dajun Chen, Hongwei Chen, Liang Chen, et al. 2023. Codefuse-13b: A pretrained multi-lingual code large language model. *arXiv preprint arXiv:2310.06266*.
- Daniel Fried, Armen Aghajanyan, Jessy Lin, Sida Wang, Eric Wallace, Freda Shi, Ruiqi Zhong, Wen-tau Yih, Luke Zettlemoyer, and Mike Lewis. 2022. InCoder: A generative model for code infilling and synthesis. *arXiv preprint arXiv:2204.05999*.
- Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen, Yujia Yang, Nan Duan, and Weizhu Chen. 2023. Critic: Large language models can self-correct with tool-interactive critiquing. *arXiv preprint arXiv:2305.11738*.
- Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou, Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, et al. 2023. StarCoder: may the source be with you! *arXiv preprint arXiv:2305.06161*.
- Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xiubo Geng, Wenxiang Hu, Chongyang Tao, Jing Ma, Qingwei Lin, and Daxin Jiang. 2023. WizardCoder: Empowering code large language models with evolve-instruct. *arXiv preprint arXiv:2306.08568*.
- Thibaud Lutellier, Hung Viet Pham, Lawrence Pang, Yitong Li, Moshi Wei, and Lin Tan. 2020. Coconut: combining context-aware neural translation models using ensemble for program repair. In *Proceedings of the 29th ACM SIGSOFT international symposium on software testing and analysis*, pages 101–114.
- Ehsan Mashhadi and Hadi Hemmati. 2021. Applying codebert for automated program repair of java simple bugs. In *2021 IEEE/ACM 18th International Conference on Mining Software Repositories (MSR)*, pages 505–509. IEEE.
- Niklas Muennighoff, Qian Liu, Armel Zebaze, Qinkai Zheng, Binyuan Hui, Terry Yue Zhuo, Swayam Singh, Xiangru Tang, Leandro von Werra, and Shayne Longpre. 2023. Octopack: Instruction tuning code large language models. *arXiv preprint arXiv:2308.07124*.
- Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, M Saiful Bari, Sheng Shen, Zheng-Xin Yong, Hailey Schoelkopf, et al. 2022. Crosslingual generalization through multitask finetuning. *arXiv preprint arXiv:2211.01786*.
- Erik Nijkamp, Bo Pang, Hiroaki Hayashi, Lifu Tu, Huan Wang, Yingbo Zhou, Silvio Savarese, and Caiming Xiong. 2022. Codegen: An open large language model for code with multi-turn program synthesis. *arXiv preprint arXiv:2203.13474*.
- OpenAI. 2023. [Gpt-4 technical report](#).
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744.
- Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, et al. 2023. Code llama: Open foundation models for code. *arXiv preprint arXiv:2308.12950*.
- Bo Shen, Jiabin Zhang, Taihong Chen, Daoguang Zan, Bing Geng, An Fu, Muhan Zeng, Ailun Yu, Jichuan

Ji, Jingyang Zhao, et al. 2023. Pangu-coder2: Boosting large language models for code with ranking feedback. *arXiv preprint arXiv:2307.14936*.

Xingyao Wang, Hao Peng, Reyhaneh Jabbarvand, and Heng Ji. 2023a. Leti: Learning to generate from textual interactions. *arXiv preprint arXiv:2305.10314*.

Yue Wang, Hung Le, Akhilesh Deepak Gotmare, Nghi DQ Bui, Junnan Li, and Steven CH Hoi. 2023b. Codet5+: Open code large language models for code understanding and generation. *arXiv preprint arXiv:2305.07922*.

Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. 2021. Finetuned language models are zero-shot learners. *arXiv preprint arXiv:2109.01652*.

Michihiro Yasunaga and Percy Liang. 2021. Break-it-fix-it: Unsupervised learning for program repair. In *International Conference on Machine Learning*, pages 11941–11952. PMLR.

Jiali Zeng, Fandong Meng, Yongjing Yin, and Jie Zhou. 2023. Tim: Teaching large language models to translate with comparison. *arXiv preprint arXiv:2307.04408*.

Kechi Zhang, Zhuo Li, Jia Li, Ge Li, and Zhi Jin. 2023. Self-edit: Fault-aware code editor for code generation. *arXiv preprint arXiv:2305.04087*.

Qinkai Zheng, Xiao Xia, Xu Zou, Yuxiao Dong, Shan Wang, Yufei Xue, Zihan Wang, Lei Shen, Andi Wang, Yang Li, et al. 2023. Codegeex: A pre-trained model for code generation with multilingual evaluations on humaneval-x. *arXiv preprint arXiv:2303.17568*.

A Construction of erroneous codes

In this work, we focus on incorporating token-level bugs into codes. We add the following types of bugs: 1) Misuse variables in the code, as shown in 4. 2) Misuse operators in the code, as shown in 5. 3) Miss functions in the code, as shown in 6. More methods can be tried to create erroneous examples, such as using GPT-4 for generation. We will leave this part of the work for the future.

B Templates for Sequence-level Comparison

The templates we used to construct sequence-level comparison examples are illustrated in 4

Templates

Given the instruction: <Instruction>

Here is a piece of code with bugs: <Bug output>

Fix the bugs in the code.

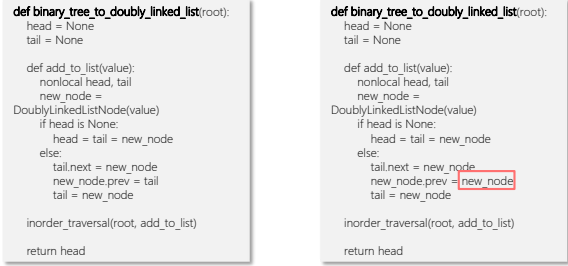
<Bug output> is the code implementation of <Instruction>.

However, there are some bugs in the code

Please fix bugs in the code.

Find the bugs in the <Bug output>

Table 4: Templates for constructing sequence-level comparison examples.



```
def binary_tree_to_doubly_linked_list(root):
    head = None
    tail = None

    def add_to_list(value):
        nonlocal head, tail
        new_node = DoublyLinkedListNode(value)
        if head is None:
            head = tail = new_node
        else:
            tail.next = new_node
            new_node.prev = tail
            tail = new_node

    inorder_traversal(root, add_to_list)

    return head
```

Figure 4: Variable misuse bug example. The buggy code (right) incorrectly uses 'newcode'.



```
def below_zero(operations: List[int]):
    balance = 0

    for op in operations:
        balance += op
        if balance < 0:
            return True
    return False
```

Figure 5: Operator misuse bug example. The buggy code (right) incorrectly uses 'greater than'.



```
def max_difference(arr):
    max_diff = 0
    for i in range(len(arr)-1):
        diff = arr[i+1] - arr[i]
        if diff > max_diff:
            max_diff = diff
    return max_diff
```

Figure 6: Function missing bug example. The buggy code (right) removes 'len()' function.