

CodePatchLLM: Configuring code generation using a static analyzer

Danil Shaikhelislamov
Ivannikov Institute for
System Programming of the
Russian Academy of Sciences
Moscow, Russia
shaykhelislamov.ds@ispras.ru

Mikhail Drobyshvskiy
ISP RAS Research Center for
Trusted Artificial Intelligence
Moscow, Russia

Andrey Belevantsev
Ivannikov Institute for
System Programming of the
Russian Academy of Sciences
Moscow, Russia

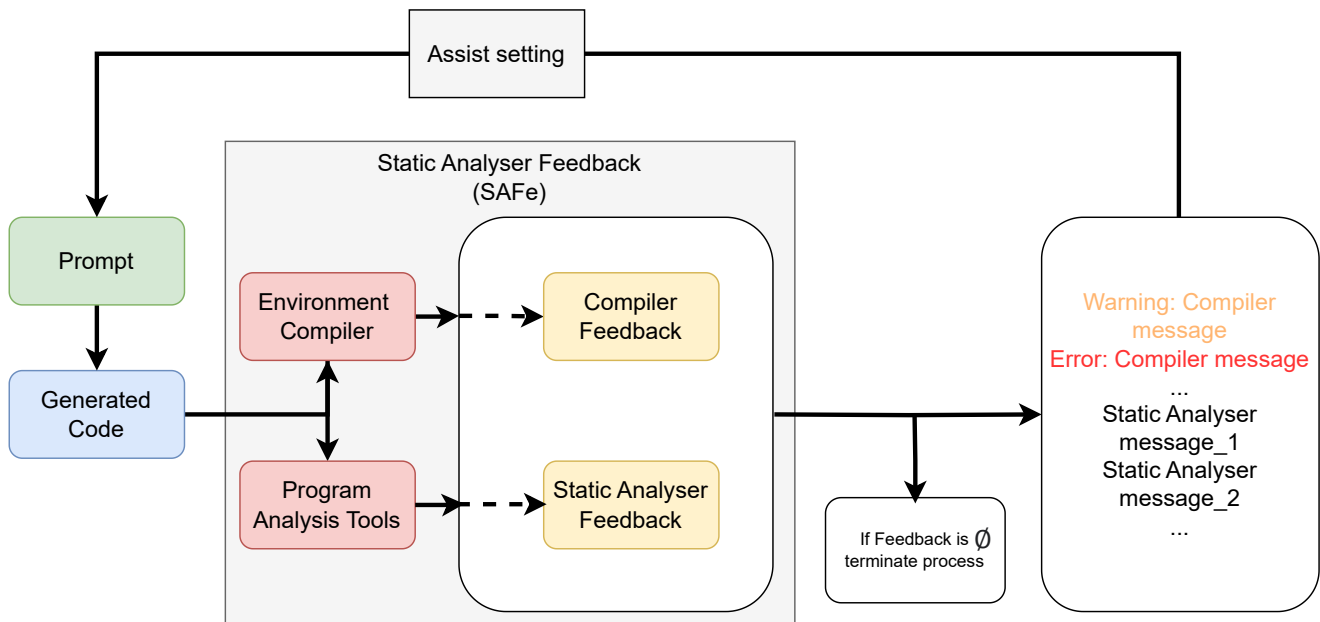


Figure 1: Overview of the model-agnostic CodePatchLLM.

ABSTRACT

The development of large language models (LM) has significantly advanced the field of code generation. A survey of developers by Stack Overflow has found that 70% of respondents are using or plan to use AI coding tools this year [44]. Current approaches mainly rely on supervised fine-tuning objectives borrowed from text generation, neglecting unique sequence-level characteristics of code, including but not limited to compilability as well as syntactic and functional correctness. To address this limitation, we propose a new approach to code generation that synergistically combines pre-trained LLM models with software analysis tools, which are widely used to check

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

GenAI Evaluation KDD2024, August 26, 2024, Barcelona, Spain

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

for vulnerabilities while validating the code. By utilizing expanded messages from code compilation and analysing, proposed approach seamlessly integrates external code-specific knowledge into the prompt chaining process. We develop CodePatchLLM, an extension for LLM that utilizes Svmc feedback for code generation. It is important to note that CodePatchLLM is a model-agnostic framework that can be used across different program languages. Extensive experiments on LeetCode dataset demonstrate the effectiveness of our proposed approach compared to backbone model, CodeLlama, achieving significant improvements in compilation success rates and functional correctness across Java, Python and Kotlin languages. Our CodePatchLLM code is available online¹.

KEYWORDS

Large Language Model, Static Analyzer, Code Quality

ACM Reference Format:

Danil Shaikhelislamov, Mikhail Drobyshvskiy, and Andrey Belevantsev. 2024. CodePatchLLM: Configuring code generation using a static analyzer. In

¹<https://github.com/dsshay/CodePatchLLM>

Proceedings of KDD workshop on Evaluation and Trustworthiness of Generative AI Models (GenAI Evaluation KDD2024). ACM, New York, NY, USA, 7 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

Code generation or program synthesis aims to automatically generate source code that adheres to a specified programming requirement, which is typically described in a natural language [10, 42]. Recently, with the development of large language models (LLM), techniques based on LLMs [1, 35, 36] have demonstrated impressive ability in code generation. However, challenges persist with the use of generated code in complex systems [7, 8, 12, 54], indicating a remaining gap in fully meeting user expectations.

In this context, learning from automatic defect detection tools demonstrates exciting potential to enhance the comprehension of complicated technical specifications and the quality of generated codes [27]. Feedback from compilation and execution results is instrumental in directly ascertaining the functional correctness of programs [8, 45]. Researchers have introduced leveraging compiler feedback from unit tests to guide the exploration of the output space of LLMs [29, 40] using reinforcement learning techniques. In other words, authors fine-tuned the model so that the output program was built successfully and passed tests.

Nevertheless, optimizing LLMs for code generation via compiler feedback presents several challenges. Firstly, the increasing complexity of human requests to LLMs often results in the generation of longer code sequences, and this worsens the final program quality [11, 22]. Secondly, feedback solely from independent unit tests and a compiler is not enough for reliability of such a program. Static analysis tools conduct more thorough source code checks than compilers, which usually only detect syntax errors [4].

Automated code generation (or program synthesis) has attracted much attention over the past few years [30], because of its potential to improve the productivity of developers, as well as to speed up the software development [33]. Companies that hastily implement generative solutions or succumb to the “AI hype” may face a potential increase in cybersecurity risk. Recent work [15, 43] highlights the importance of using robust implementations of generative AI in a business environment to mitigate such risks. The demand for reliable and secure code has never been higher.

To tackle these challenges, several approaches are proposed, such as filtering and repairing the non-executable synthesized programs [21], using energy-based generation models with execution constraints [20], and reinforcement learning (RL) fine-tuning mechanisms [8, 23, 40]. However, existing approaches are often tailored to a specific programming language (PL) or task, e.g., [23] is exclusively designed for program synthesis task in Python.

We introduce a new approach, illustrated in Fig. 1, combining results of program analysis and LLM for code generation. An overview of the proposed approach. Generated Code are first initialized from the pre-trained LLM for the designed task. The generated code is completely transferred to the compiler and static analyzer (SAFE module) suitable for the selected programming language. Detected warnings and errors are collected in a single pool of messages that are transmitted to Assist Setting. Finally, the LLM model prompt is updated based on the obtained values and returns.

Initially, the output of any pre-trained LLM tailored for code generation is transferred to the program analysis module, SAFE (Software Analysis Feedback). The module includes but is not limited to the use of a compiler and a static analyzer. It is assumed that it is possible to include other analysis tools such as DAST [9], IAST [32]. In particular, SAFE analyzes the entire file generated by the program to identify potential vulnerabilities in the architecture of the solution that may occur when using external libraries and classes. The usual classic static analysis stages are performed on the generated program, namely, capturing the program build to generate automatically the required intermediate representation, lightweight analysis of the program’s syntax trees (AST-level analysis), and interprocedural dataflow analysis (which is both context-sensitive and path-sensitive based on symbolic execution). Eventually, the goal is to identify all possible errors in the generated program, before it is used by humans. One distinctive feature of this approach is the utilization of feedback from program analysis tools as additional contextual clues for the LLM.

Specifically, messages generated by the compiler and static analyzer are incorporated into the LLM’s prompt as supplementary comments (in Assist Setting). By enriching the prompt with insights gathered from program analysis, the LLM gains a deeper understanding of the desired code’s requirements, constraints, and potential pitfalls. Once the prompt is augmented with relevant analysis feedback, the LLM is prompted to generate a new code, embedding the provided comments into its output. This iterative process fosters a symbiotic relationship between automated program analysis and advanced language modeling, facilitating the generation of code that not only adheres to syntactic rules but also aligns with best practices, security guidelines, and architectural constraints.

In addition, we propose a novel framework, CodePatchLLM, which integrates the Svac static analyzer [2, 16] feedback into CodeLlama [36] to enhance the reliability and security of generated code. Through a series of experiments and evaluations on Leetcode dataset [13] we seek to demonstrate the effectiveness of our approach in improving code quality, reducing the risk of defects, and ultimately enhancing the trustworthiness of software systems in real-world applications.

To summarize, the major contributions of this paper are as follows:

- We introduce a novel approach that utilizes code-specific feedback as the external source of knowledge in model instructions. The approach is independent of models architecture and generates higher-quality codes.
- We develop the CodePatchLLM extension for CodeLlama that applies the Svac static analyzer to the generated code, and iteratively corrects the model’s prompt using all warnings and errors detected by Svac;
- We demonstrate the effectiveness of CodePatchLLM through experiments across diverse programming tasks (from the Leetcode platform²) and program languages (Java, Python, Kotlin). Using CodeLlama with CodePatchLLM improves the compilation rate in 50% for Java and 10% for Kotlin more

²<https://leetcode.com>

cases and functional correctness over different languages by 5%.

The remainder of this paper is organized as follows. In Section 2 we describe existing code generation models utilizing external knowledge and structure-based approaches. Section 3 delves into the specifics of our proposed approach and new CodePatchLLM framework that includes the Svace static analyzer. The experimental evaluation of CodePatchLLM on programming tasks written in three program languages (Java, Python, Kotlin) and the case study can be found in Section 4. Finally, Section 5 concludes the paper.

2 RELATED WORK

2.1 Fine-tuning large language models for code generation

Recently, LLMs have shown remarkable ability in understanding natural language and code generation by training on large text corpora containing code data. Several pre-trained language models (PLMs) demonstrate significant potential for code generation including CodeGPT [30], PanGu-Coder [38], SantaCoder [3]. In addition, supervised fine-tuning models achieve more competitive performance such as CodeX [6], CodeLlama Instruct [36].

Reinforcement Learning is a method of learning the optimal policy by exploring the environment and obtaining rewards [48]. Recently, some researchers have introduced RL to LLMs and improved the quality of the generated code by utilizing the unit test feedback to explore the output space of the policy model [8, 23, 28, 29, 40]. For instance, Coberly [23] leverages signal from unit tests as rewards and utilizes the actor-critic approach [19] to enhance models on code generation. PPOCoder [40] refines Code by employing the PPO algorithm [37] and RLTF [29] provides fine-grained rewards through the error locations, but the reward space is still sparse. However, the exploration of complex tasks in an environment characterized by a sparse reward is challenging. These methods still fall short of effectively using RL to enhance the model’s performance in code generation [54].

2.2 Chain-of-thought prompting

With the recent advancements in large language models, researchers have discovered that utilizing the chain-of-thought (CoT) [14, 46] techniques can significantly improve reasoning abilities. Authors [46] introduced the concept of few-shot CoT, which involves generating intermediate reasoning steps before arriving at the final answer with in-context demonstrations. This approach deviates from traditional few-shot prompting (also called in-context learning [5]) that directly generates the final answer. Zero-shot CoT [18] is another method leveraging chain-of-thought which adding the prompt “Let’s think step by step.” after the task description to activate LLMs to generate rationales in order for improved task performance. Other researchers also propose various prompting methods to enhance model capabilities, including auto-cot [52], least-to-more [53], decomposing prompting [17] and tree-of-thought [49]. In our work, outputs of program analysis tools can be seen as kind of prompt chaining, since all this data serves as intermediate steps for fixing bugs in the code.

2.3 Prompting with feedback

Despite the remarkable capabilities of large language models, it can still be challenging sometimes to generate the correct answer in a single attempt. Recently, people found that LLMs can receive feedback from external environment or generated by themselves and iteratively refine according to the feedback. Self-refine [31] launches a novel approach that allows LLMs to iteratively refine outputs with the feedback without any labeled data. Reflexion [39] proposes a “verbal reinforcement learning” that LLMs reflect on failures based on feedback and store reflexion in a text style for future trials. REMEMBERER [50] employs a method that allows LLMs to learn experience which is stored in an external memory from the feedback in the training set and transfer that experience to the test set for a better performance. In our work, we focus on code generation task and teach LLMs to improve code based on feedback from program analysis tools such as a static analyzer.

2.4 Prompting for code

Prompting techniques have been extensively utilized in tasks related to code. Some works including [24–26] focus on leveraging prompting to enhance code generation. In [51] prompting is used to facilitate code selection and develop a reviewer model. We focus on using program analysis tools to improve the compilability and security of the code without losing the effectiveness of solving the problem. We use a message from the compiler and static analyzer as a signal for the model.

3 METHOD

In this section, we focus on the methodological details of our approach, which ensures the generation of a compiled program and a program tested by a static analyzer, respectively, as shown in Fig. 1. And we describe CodePatchLLM, a novel framework that uses Svace [2, 16] feedback for correcting prompt for CodeLlama [36], illustrated in Fig. 3.

3.1 Comparing prompting and fine-tuning approaches for code generation tasks

We aim to dissect several key aspects underlying the efficacy of prompting as a methodology over traditional fine-tuning practices. We begin by examining insights gleaned from empirical studies [17, 25, 46].

Better understanding. Consistent improvement of parts of the generated code can facilitate better understanding, error correction, and optimization, ultimately leading to improved overall performance of the LLM [25].

Activating specific internal knowledge. Using prompting activates specific internal knowledge.

In the single step approach, the LLM is constrained to solve the entire problem in a single step. Therefore, the maximum number of tokens it can process is limited by s . In the n intermediary steps approach, the LLM solves the problem in multiple steps, with each step handling a portion of the task. Since the LLM can process tokens in each of the n intermediary steps, the total number of tokens processed across all steps is $s + n \times m$.

Comparing the two approaches, we can see that the n intermediary steps approach allows the LLM to process a total of $s + n \times m$

tokens, which is significantly higher than the token processing capacity of s in the single step approach.

Independence from the model architecture and software tools. Fine-tuning a model with specific static analyzer can lead to overfitting, where the model performs well on errors and warnings that can only be detected by this static analyzer but poorly on unseen data.

Moreover, fine-tuning can sometimes lead to a loss of generalization ability, meaning the model becomes too specialized for the specific task it was fine-tuned for and performs poorly on related tasks or in different environments.

The transferability of a fine-tuned model to different tasks or domains may be limited. While fine-tuning may improve performance on a particular task, it might not transfer well to other tasks without further adjustments or retraining.

Customizing without computational cost. Fine-tuning can be computationally expensive, especially if the pre-trained model is large and the fine-tuning dataset is extensive. This can require significant computational resources and time, making fine-tuning less practical in some scenarios [47].

Previous works [8, 45] by fine-tuning a model based on compiler feedback shows that a reward or binary signal is usually returned as feedback, the context itself does not change. Thus, fine-tuning the model requires a lot of resources, since the model has to explore all the output spaces in order to find the best combination that will lead to a rare reward, program compilability.

External signals do not detract the model from solving the problem. If the previous code snippet is compilable, the generator can fool the compiler easily. The RL is good at making use of this, resulting in the generated code can be compiled, but seriously deviating from the generation likelihood objective.

Previous works [8, 29] to avoid active model being too far away from reference model added a Kullback-Leibler [54] penalty with expectation. To alleviate the imbalance between the reward term and the penalty term and improve the stability of training, authors in [34] used autoregressive fine-tuning. This general setup leads to the fact that the original model corrects only a small number of tokens at the input, which in turn may not be optimal for quickly improving the compilability and security of the code, since the entire solution architecture has to be changed.

Based on this, we propose a simple and convenient framework CodePatchLLM to extend CodeLlama [36] capabilities for code generation through feedback from Svace [16]. The flexible configuration of the framework allows to use any architecture of the LM with CodePatchLLM that solves the task. The authors believe that frameworks like CodePatchLLM can be used as add-ons on an arbitrary code generation model. And can be used when the specification requirement is higher than under normal conditions.

3.2 CodePatchLLM

Our proposed method CodePatchLLM enables large language models to use the code debugging method through a static analyzer. The CodePatchLLM can be used in cases where the requirements for the generated code are much higher than in normal conditions: the code must be vulnerability-tested and executable. It involves testing generated code through the Svace static analyzer. Programmers can

analyze the output generated by the dialogue of LLM and Svace to understand how the code has been changed and what errors were in the first version.

Figure 3 shows an example of the first iteration using CodePatchLLM on the program generation task. We first let LLMs attempt solving the programming problem based solely on the problem description, without any extra information. The generated code, regardless of its size, is checked by the Svace static analyzer (Stage SAFE) and code regeneration based on feedback from program analysis tools (Stage Assist Setting). The above steps will be repeated until Svace finds errors and vulnerabilities in the code or until several rounds of debugging attempts still fail to fix the issues.

We will now provide a detailed discussion of each step:

Generating code. For code generation we use Code Llama 70B Instruct³, the largest and best-performing open-source model. We use a typical usage prompt for many models [14]: *You are a helpful and honest code assistant expert in {LANGUAGE}. Please, provide all answers to programming questions in {LANGUAGE}.* where {LANGUAGE} denotes the selected programming language. When communicating with the model, we use the following chat prompt format:

```
chat = [
  {"role": "system", "content": "System prompt  "},
  {"role": "user", "content": "First user query"},
  {"role": "assistant", "content": "Model response to first query"},
  {"role": "user", "content": "Second user query"},
]
```

Figure 2: Chat Prompt for CodeLlama-70B-Instruct. Source: <https://huggingface.co/codellama/CodeLlama-70b-Instruct-hf>

SAFE stage In this stage, we proceed with code analysis that includes the running the Svace static analyzer. Following the Svace user manual, we initialize the directory in which the program file is located with the command: `svace init`. We capture any error messages provided by the Svace or the original compiler. It is important to note that even if the execution encounters errors, we still collect the log generated prior to the occurrence of the error. For instance, before running the analysis, Svace needs to build intermediate representation generated when monitoring the original compilation process (e.g. via the `javac` compiler for Java). And if an error occurs at the compilation stage, we collect the event and message in the log.

In our research, we use Svace as a base component for code analysis and optimization. Svace identifies security vulnerabilities inherent in software code, including potential buffer overflows, memory leaks, and other unsafe coding practices. By proactively detecting these vulnerabilities at an early stage of development, Svace empowers users to mitigate security risks effectively, safeguarding software systems from potential cyber threats and breaches. Additionally, Svace facilitates the identification of performance bottlenecks within the codebase, such as inefficient algorithms, redundant computations, or memory-intensive operations.

³CodeLlama-70B is free for research: <https://ai.meta.com/blog/code-llama-large-language-model-coding/>

It is assumed that the module SAFE can be flexibly expanded with other code analysis tools, for example, a dynamic analyzer.

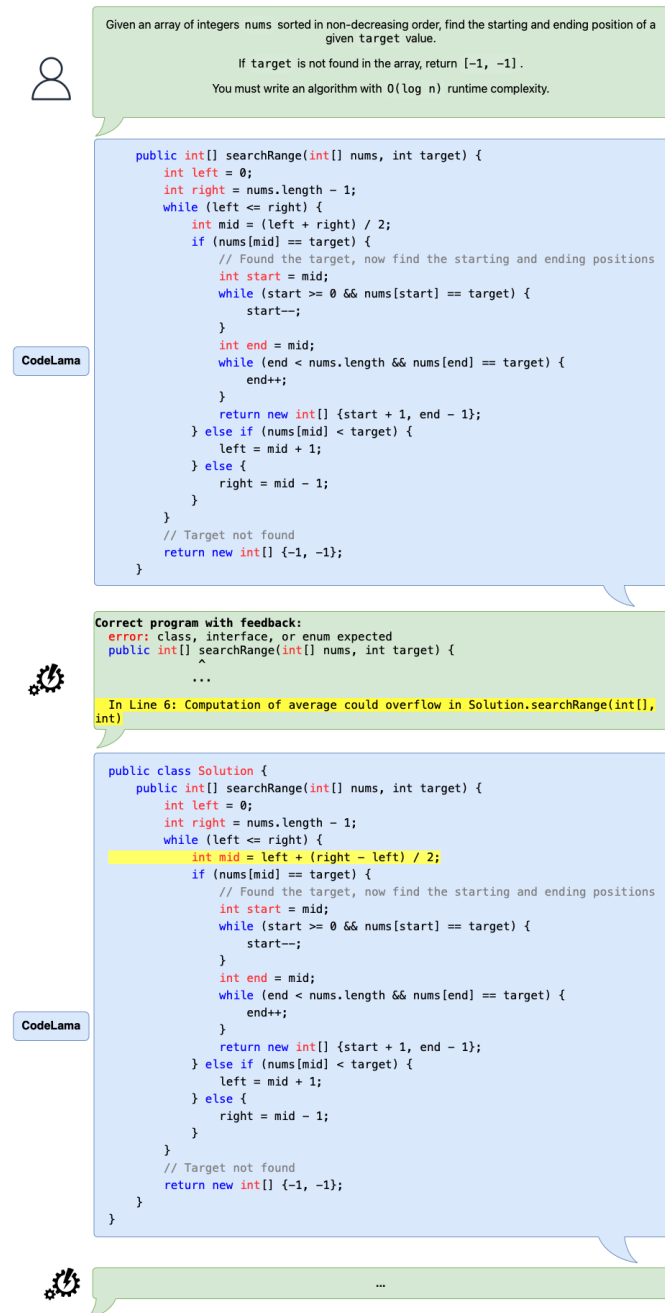


Figure 3: An illustration of Java program generation with CodePatchLLM, utilizing the Svsace feedback.

Assist settings In the final step, we provide error messages to the model turning all feedback into a sequence. Our approach aims to take into account all feedback and instructions provided by program analysis tools. This means that at each timestep, the model

can only utilize the past time steps data and itself. We instruct LLMs to regenerate the code considering the comments found, as illustrated in Figure 3. LLMs are prompted to fix an error in the code in a specific place (the line number and the method used are known). In the Figure 3 two messages generated by the SAFE module are specified: a class declaration is expected for the compilation of the program and the calculation of the average value may be overflowed.

4 EXPERIMENTS

We evaluate our CodePatchLLM framework with CodeLlama [36] as backbone. We investigate (1) how CodePatchLLM framework can boost LLMs’ performance on real-world programming tasks benchmarks; (2) the impact of the framework on code compilability and the reduction of defects; and (3) the effect of the number of iterations.

Benchmark. We consider Leetcode datasets [13] for our evaluations. LeetCode⁴ is one of the most visited platform for practicing programming. We selected LeetCode as the primary source for our evaluation dataset, as it programming tasks can be directly compiled by copying and pasting contents from the LeetCode website. A dataset consists of 2 612 programming tasks. Problems in the dataset are categorized into three levels according to their difficulties.

Metrics. To evaluate the generated codes, we employ the *pass@1* metric following [6], which calculates the percentage of problems for which all unit tests are passed using 1 synthetically generated program sample per problem.

Implementation details. The experiment process was conducted on a device with three NVIDIA A100 80G GPUs.

The weights of the model are loaded from HuggingFace. The maximum output token length is set to 2048.

At each step of the code update, we submit the solution to the platform Leetcode. And also at each stage reinitialize existing Svsace project directory from scratch.

4.1 Experimental results on Leetcode

In our study, we evaluate CodePatchLLM with Java, Python, and Kotlin languages.

The experimental results are illustrated in Table 1. The reported results indicate that the model with CodePatchLLM improve performance the model with basic settings for Java and Kotlin. And experiment shows that in Python the extended model does not improve the quality. Mechanism feedback from Svsace and original compiler designed to executable and feasibility of the program. Therefore, in part, the extension does not change the overall quality, since a Python program does not require a compiler. The lack of feedback from the compiler and the analyzer explains the minimal deviations of the results from the original model. Experimental results show that CodePatchLLM improves executability by 45% for Java and by 10% for Kotlin. Additionally, it enhances the “Accepted” rate by 50% for Java.

In Table 1, the percentage of wrong answers and runtime errors increased after applying CodePatchLLM for the Java and Kotlin languages. Specifically, wrong answers increased from 22% to 37% for Java and from 2% to 10% for Kotlin. This is due to the distribution

⁴<https://leetcode.com>

Table 1: Performance results. Overall, CodePatchLLM boosts the performance and compilability of CodeLlama on the discussed metric.

| | | Accepted | Wrong Answer | Runtime Error | Compiler Error |
|--------|--------------------------------|----------|--------------|---------------|----------------|
| Java | CodeLlama[36] | 12% | 22% | 3% | 62% |
| | CodeLlama + CodePatchLLM (Our) | 25% | 37% | 5% | 33% |
| Python | CodeLlama[36] | 36% | 47% | 17% | 0% |
| | CodeLlama + CodePatchLLM (Our) | | same | result | |
| Kotlin | CodeLlama[36] | 10% | 2% | 24% | 64% |
| | CodeLlama + CodePatchLLM (Our) | 12% | 10% | 24% | 54% |

of compiler error cases among other cases. In particular, some of the tasks became compiled, but had wrong answer.

In 37.3% of cases, CodeLlama incorrectly names the implemented class method at the first request. For example, `swapAdjacentNodes` instead `swapPairs`. It is noteworthy that further dialogue with the model allows you to correct this although in the task formulation there are no clarifications in the naming of methods and classes.

Since CodePatchLLM constrains the model on an example when predicting another one, the model can simply “copy” the example without learning to understand the underlying task. In future, to address this, we can randomly mask between 0% and 5% of past tokens during training, which help regularize the model and prevent it from overfitting to the specific examples seen during training [28, 41].

The Table 2 presents the percentage of generated programs with errors and vulnerabilities, categorized based on the type of feedback received from the original compiler and Svmc static analyzer. the majority of the feedback (98%) has been addressed by CodePatchLLM, indicating its effectiveness in correcting errors. In Java, 12.5% of the generated programs had Svmc feedback indicating vulnerabilities, while for Kotlin, this percentage was lower at 3.1%. No errors were detected in the generated Python programs, likely due to the Python support in Svmc being in the first release, so that not many Svmc checkers are supported for Python. It is worth noting that the analyzer verification stage comes after the compiler verification, and Svmc identifies errors that do not influence whether the program can be successfully built, as all programs being analyzed are already compiled correctly. But the errors found can affect functional correctness.

Table 2: The percentage of generated programs with errors and vulnerabilities. 98% of the reviews have been corrected by CodePatchLLM.

| | Compiler feedback | Svmc feedback |
|--------|-------------------|---------------|
| Java | 62,3% | 12,5% |
| Python | — | 0% |
| Kotlin | 64,5% | 3,1% |

Our approach allows to create a ready-made extension that can be used in conjunction with any LLM to generate code. Experiments with CodePatchLLM demonstrate that the method of using

feedback from the original compiler and static analyzer improves compilability and prevents errors for Kotlin and Java programs.

5 CONCLUSION

In this paper, we introduce CodePatchLLM, a novel framework that leverages feedback from compilers and static analyzers to correct prompts. We identified some limitations of fine-tuning LLM for code generation tasks and designed a new framework that is geared towards program languages as opposed to natural language. We incorporated compiler and static analyzer’s feedback into our framework to encourage the model to generate more syntactically and logically correct codes. Results of our experiments show the effectiveness of our method compared to the basic model without CodePatchLLM in improving the syntactic / functional correctness of the generated codes. It is important to understand that one of the limitations of CodePatchLLM is the additional time spent on data exchange, which can increase computational requirements. However, CodePatchLLM is primarily motivated by its ability to enhance the performance of pre-trained models which is a more cost-effective strategy than fine-tuning ones to specific task.

By addressing the critical need for trustworthy code generation, this research has the potential to contribute to the development of more reliable and secure software systems.

REFERENCES

- [1] Dmitry Abulkhanov, Nikita Sorokin, Sergey Nikolenko, and Valentin Malykh. 2023. Lapca: Language-agnostic pretraining with cross-lingual alignment. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 2098–2102.
- [2] Andrey Belevantsev, Alexey Borodin, Irina Dudina, Valery Ignatiev, Alexey Izbyshchev, Sergey Polyakov, Evgeny Vesevich, and Dmitry Zhurikhin. 2018. Design and Development of Svmc Static Analyzers. In *2018 Ivannikov Memorial Workshop (IVMEM)*. 3–9. <https://doi.org/10.1109/IVMEM.2018.00008>
- [3] Loubna Ben Allal, Raymond Li, Denis Kocetkov, Chenghao Mou, Christopher Akiki, Carlos Munoz Ferrandis, Niklas Muennighoff, Mayank Mishra, Alex Gu, Manan Dey, et al. 2023. SantaCoder: don’t reach for the stars! *arXiv e-prints* (2023), arXiv-2301.
- [4] Bruno Blanchet, Patrick Cousot, Radhia Cousot, Jérôme Feret, Laurent Mauborgne, Antoine Miné, David Monniaux, and Xavier Rival. 2003. A static analyzer for large safety-critical software. In *Proceedings of the ACM SIGPLAN 2003 conference on Programming language design and implementation*. 196–207.
- [5] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems* 33 (2020), 1877–1901.
- [6] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374* (2021).

- [7] Fenia Christopoulou, Gerasimos Lampouras, Milan Gritta, Guchun Zhang, Yinpeng Guo, Zhongqi Li, Qi Zhang, Meng Xiao, Bo Shen, Lin Li, et al. 2022. Pangu-coder: Program synthesis with function-level language modeling. *arXiv preprint arXiv:2207.11280* (2022).
- [8] Shihan Dou, Yan Liu, Haoxiang Jia, Limao Xiong, Enyu Zhou, Junjie Shan, Caishuang Huang, Wei Shen, Xiaoran Fan, Zhiheng Xi, et al. 2024. StepCoder: Improve Code Generation with Reinforcement Learning from Compiler Feedback. *arXiv preprint arXiv:2402.01391* (2024).
- [9] Michael Felderer, Matthias Büchler, Martin Johns, Achim D Brucker, Ruth Brey, and Alexander Pretschner. 2016. Security testing: A survey. In *Advances in Computers*. Vol. 101. Elsevier, 1–51.
- [10] Sumit Gulwani, Oleksandr Polozov, Rishabh Singh, et al. 2017. Program synthesis. *Foundations and Trends® in Programming Languages* 4, 1-2 (2017), 1–119.
- [11] Jianye Hao, Tianpei Yang, Hongyao Tang, Chenjia Bai, Jinyi Liu, Zhaopeng Meng, Peng Liu, and Zhen Wang. 2023. Exploration in deep reinforcement learning: From single-agent to multiagent domain. *IEEE Transactions on Neural Networks and Learning Systems* (2023).
- [12] Dan Hendrycks, Steven Basart, Saurav Kadavath, Mantas Mazeika, Akul Arora, Ethan Guo, Collin Burns, Samir Puranik, Horace He, Dawn Song, et al. 2021. Measuring coding challenge competence with apps. *arXiv preprint arXiv:2105.09938* (2021).
- [13] Wenpin Hou and Zhicheng Ji. 2024. A systematic evaluation of large language models for generating programming code. *arXiv preprint arXiv:2403.00894* (2024).
- [14] Xueyu Hu, Kun Kuang, Jiankai Sun, Hongxia Yang, and Fei Wu. 2024. Leveraging print debugging to improve code generation in large language models. *arXiv preprint arXiv:2401.05319* (2024).
- [15] Declan Humphreys, Abigail Koay, Dennis Desmond, and Erica Mealy. 2024. AI hype as a cyber security risk: the moral responsibility of implementing generative AI in business. *AI and Ethics* (2024), 1–14.
- [16] VP Ivannikov, AA Belevantsev, AE Borodin, VN Ignatiev, DM Zhurikhin, and AI Avetisyan. 2014. Static analyzer Svace for finding defects in a source program code. *Programming and Computer Software* 40 (2014), 265–275.
- [17] Tushar Khot, Harsh Trivedi, Matthew Finlayson, Yao Fu, Kyle Richardson, Peter Clark, and Ashish Sabharwal. 2022. Decomposed prompting: A modular approach for solving complex tasks. *arXiv preprint arXiv:2210.02406* (2022).
- [18] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. *Advances in neural information processing systems* 35 (2022), 22199–22213.
- [19] Vijay Konda and John Tsitsiklis. 1999. Actor-critic algorithms. *Advances in neural information processing systems* 12 (1999).
- [20] Tomasz Korbak, Hady Elsahar, Marc Dymetman, and Germán Kruszewski. 2021. Energy-based models for code generation under compilability constraints. *arXiv preprint arXiv:2106.04985* (2021).
- [21] Sumith Kulal, Panupong Pasupat, Kartik Chandra, Mina Lee, Oded Padon, Alex Aiken, and Percy S Liang. 2019. Spoc: Search-based pseudocode to code. *Advances in Neural Information Processing Systems* 32 (2019).
- [22] Pawel Ladosz, Lilian Weng, Minwoo Kim, and Hyondong Oh. 2022. Exploration in deep reinforcement learning: A survey. *Information Fusion* 85 (2022), 1–22.
- [23] Hung Le, Yue Wang, Akhilesh Deepak Gotmare, Silvio Savarese, and Steven Chu Hong Hoi. 2022. Coderl: Mastering code generation through pretrained models and deep reinforcement learning. *Advances in Neural Information Processing Systems* 35 (2022), 21314–21328.
- [24] Jia Li, Ge Li, Yongmin Li, and Zhi Jin. 2023. Enabling programming thinking in large language models toward code generation. *arXiv preprint arXiv:2305.06599* (2023).
- [25] Jierui Li, Szymon Tworkowski, Yingying Wu, and Raymond Mooney. 2023. Explaining competitive-level programming solutions using llms. *arXiv preprint arXiv:2307.05337* (2023).
- [26] Xin-Ye Li, Jiang-Tian Xue, Zheng Xie, and Ming Li. 2023. Think outside the code: Brainstorming boosts large language models in code generation. *arXiv preprint arXiv:2305.10679* (2023).
- [27] Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, et al. 2022. Competition-level code generation with alphacode. *Science* 378, 6624 (2022), 1092–1097.
- [28] Hao Liu, Xinyang Geng, Lisa Lee, Igor Mordatch, Sergey Levine, Sharan Narang, and Pieter Abbeel. 2022. Towards Better Few-Shot and Finetuning Performance with Forgetful Causal Language Models. *arXiv preprint arXiv:2210.13432* (2022).
- [29] Jiatae Liu, Yiqin Zhu, Kaiwen Xiao, Qiang Fu, Xiao Han, Wei Wang, and Deheng Ye. 2023. Rlhf: Reinforcement learning from unit test feedback. *arXiv preprint arXiv:2307.04349* (2023).
- [30] Shuai Lu, Daya Guo, Shuo Ren, Junjie Huang, Alexey Svyatkovskiy, Ambrosio Blanco, Colin Clement, Dawn Drain, Daxin Jiang, Duyu Tang, et al. 2021. Codeglue: A machine learning benchmark dataset for code understanding and generation. *arXiv preprint arXiv:2102.04664* (2021).
- [31] Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. 2024. Self-refine: Iterative refinement with self-feedback. *Advances in Neural Information Processing Systems* 36 (2024).
- [32] Yuanyuan Pan. 2019. Interactive application security testing. In *2019 International Conference on Smart Grid and Electrical Automation (ICSGEA)*. IEEE, 558–561.
- [33] Md Rizwan Parvez, Wasi Uddin Ahmad, Saikat Chakraborty, Baishakhi Ray, and Kai-Wei Chang. 2021. Retrieval augmented code generation and summarization. *arXiv preprint arXiv:2108.11601* (2021).
- [34] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog* 1, 8 (2019), 9.
- [35] Anton Razzhigaev, Mikhail Salnikov, Valentin Malykh, Pavel Braslavski, and Alexander Panchenko. 2023. A system for answering simple questions in multiple languages. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*. 524–537.
- [36] 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* (2023).
- [37] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347* (2017).
- [38] Bo Shen, Jiaxin 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* (2023).
- [39] Noah Shinn, Beck Labash, and Ashwin Gopinath. 2023. Reflexion: an autonomous agent with dynamic memory and self-reflection. *arXiv preprint arXiv:2303.11366* (2023).
- [40] Parshin Shojae, Aneesh Jain, Sindhu Tipirneni, and Chandan K Reddy. 2023. Execution-based code generation using deep reinforcement learning. *arXiv preprint arXiv:2301.13816* (2023).
- [41] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research* 15, 1 (2014), 1929–1958.
- [42] Alexey Svyatkovskiy, Shao Kun Deng, Shengyu Fu, and Neel Sundaresan. 2020. Intellicode compose: Code generation using transformer. In *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 1433–1443.
- [43] D Yu Turdakov, Arutyun Ishkhanovich Avetisyan, Konstantin Vladimirovich Arkhipenko, Anastasiya Vsevolodovna Antsiferova, Dmitry Sergeevich Vatolin, SS Volkov, Alexander Vladimirovich Gasmikov, Dmitry Alekseevich Devyatkin, MD Drobyshevsky, AP Kovalenko, et al. 2022. Trusted artificial intelligence: challenges and promising solutions. In *Doklady Mathematics*, Vol. 106. Springer, S9–S13.
- [44] James Vincent. 2023. Stack Overflow survey finds developers are ready to use AI tools – even if they don’t fully trust them. *VOX MEDIA* (2023). <https://www.theverge.com/2023/6/13/23759101/stack-overflow-developers-survey-ai-coding-tools-moderators-strike>
- [45] Xin Wang, Yasheng Wang, Yao Wan, Fei Mi, Yitong Li, Pingyi Zhou, Jin Liu, Hao Wu, Xin Jiang, and Qun Liu. 2022. Compileable neural code generation with compiler feedback. *arXiv preprint arXiv:2203.05132* (2022).
- [46] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems* 35 (2022), 24824–24837.
- [47] Marco A Wiering and Martijn Van Otterlo. 2012. Reinforcement learning. *Adaptation, learning, and optimization* 12, 3 (2012), 729.
- [48] Ronald J Williams. 1992. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning* 8 (1992), 229–256.
- [49] Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. 2024. Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural Information Processing Systems* 36 (2024).
- [50] Danyang Zhang, Lu Chen, Situo Zhang, Hongshen Xu, Zihan Zhao, and Kai Yu. 2024. Large Language Models Are Semi-Parametric Reinforcement Learning Agents. *Advances in Neural Information Processing Systems* 36 (2024).
- [51] Tianyi Zhang, Tao Yu, Tatsunori Hashimoto, Mike Lewis, Wen-tau Yih, Daniel Fried, and Sida Wang. 2023. Coder reviewer reranking for code generation. In *International Conference on Machine Learning*. PMLR, 41832–41846.
- [52] Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. 2022. Automatic chain of thought prompting in large language models. *arXiv preprint arXiv:2210.03493* (2022).
- [53] Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc Le, et al. 2022. Least-to-most prompting enables complex reasoning in large language models. *arXiv preprint arXiv:2205.10625* (2022).
- [54] Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. 2019. Fine-tuning language models from human preferences. *arXiv preprint arXiv:1909.08593* (2019).