
PERFCODEGEN: Improving Performance of LLM Generated Code with Execution Feedback

Yun Peng^{*1}, Akhilesh Deepak Gotmare^{†2}, Michael Lyu¹, Caiming Xiong²
Silvio Savarese², Doyen Sahoo²

¹The Chinese University of Hong Kong

²Salesforce Research

Abstract

Large Language Models (LLMs) are widely adopted for assisting in software development tasks, yet their performance evaluations have narrowly focused on the functional correctness of generated code. Human programmers, however, require LLM-generated code to be not only correct but also optimally efficient. We propose PERFCODEGEN, a training-free framework that enhances the performance of LLM-generated code by incorporating feedback based on runtime during test case execution into the self-refinement iterations. With PERFCODEGEN, we achieve speedups for a significantly higher proportion of problems compared to using the base LLM with sophisticated prompting techniques. Applied to open language models like Phi-3-mini, PERFCODEGEN achieves runtime efficiency comparable to prompting powerful closed models like GPT-4. We achieve state-of-the-art runtime efficiency on benchmarks such as HumanEval, MBPP, and APPS, frequently surpassing the ground truth reference solutions with PERFCODEGEN using GPT-3.5 and GPT-4. Additionally, we demonstrate the effectiveness of our approach in enhancing code quality across a range of open LLMs of varying sizes including Phi-3-mini, Llama 3 8B, Mixtral 8x7B, Command R, and Llama 3 70B.

1 Introduction

Language Models (LMs) are now widely used for code completion (Chen et al. [2021], Austin et al. [2021], Nijkamp et al. [2023]), as well as for tasks like unit test case generation (Chen et al. [2022]), bug fixing (Yang et al. [2024]), feature addition (Zhang et al. [2024]), and other stages of software development (Hong et al. [2023], Qian et al. [2023]). A recent Github survey (Shani [2023]) underscores this rapid proliferation, estimating that 92% of U.S. based developers are already using AI coding tools. However, despite this widespread adoption of Code LMs, evaluation has almost exclusively focused on functional correctness of generated code (Hendrycks et al. [2021], Li et al. [2022], Puri et al. [2021], Liu et al. [2024]), often overlooking key aspects of code quality. Runtime efficiency of a program, in particular, is a central consideration in software design decisions, due to its significant impact on user experience, serving costs and carbon footprint of a software application (Landes and Studer [1995], Chung et al. [2012]). Singhal et al. [2024] provide a distinct evaluation of Code LMs by focusing on non-functional requirements such as efficiency, security and maintainability, revealing that current models generally falter on these key quality metrics. This is particularly concerning as less experienced developers are more likely to accept AI-suggested code (Dohmke et al. [2023]), often neglecting quality aspects, which subsequently places a significant burden on the code review process (Harding and Kloster [2023]).

^{*}Work performed while interning at Salesforce Research

[†]Corresponding author: akhilesh.gotmare@salesforce.com

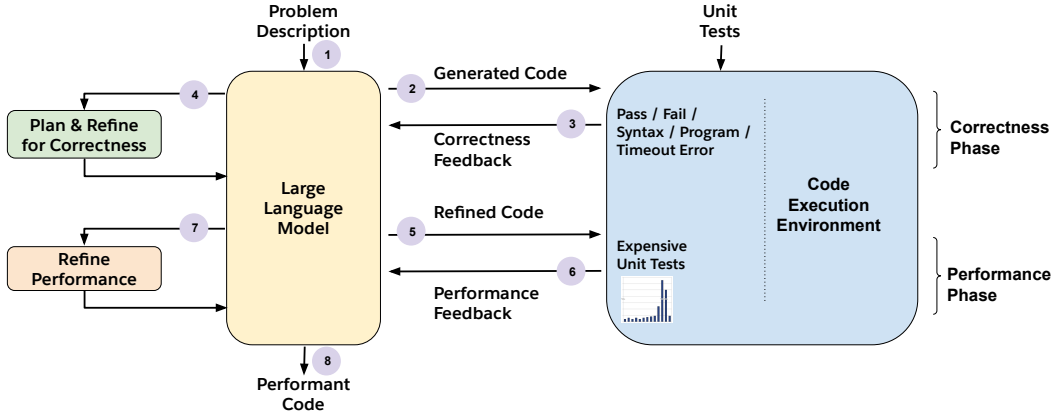


Figure 1: PERFCODEGEN Given a problem description (1), we prompt the LLM to generate a candidate solution (2), which is passed to an execution environment to collect feedback on correctness (3). The LLM is then prompted to self-reflect on this feedback in a planning stage, and accordingly generate a refinement using this context (4). This process is iterated over till correctness (4, 2, 3). Correct code obtained from this phase is self-refined for runtime-efficiency (7), and then passed to the environment to be executed (5) and the most time consuming unit test(s) is identified and passed as performance feedback to the LLM (6), that acts on it by generating a self-refinement to make the correct solution more efficient (7). This refinement is tested for correctness (2, 3) and passed as the final code solution to the problem (8) if correct, else we fall back to correct program from (3) if any.

Previous works (Madaan et al. [2023] and Garg et al. [2022]) have proposed fine-tuning LLMs to generate performance improvements for a given working program. However, this approach is challenging to scale due to the need for parallel training data in the code domain, which can be significantly more expensive to collect and manually validate than natural language data. Additionally, these methods require the availability of a functionally correct input program to optimize, which is not typically the case when writing code from scratch. Moreover, they do not leverage execution feedback from unit tests, which has been shown to be crucial in improving code correctness (Le et al. [2022]). This lack of execution feedback utilization is a notable limitation, as unit tests provide valuable insights into runtime behavior and performance characteristics.

Prompting advances such as Chain of Thought (Wei et al. [2022b], Nye et al. [2021]) and Self-Refine (Madaan et al. [2024]) have enhanced LLM output quality without additional training, but limitations remain, including issues with hallucinations, sycophancies, and unreliability in refining initial responses (Huang et al. [2023], Laban et al. [2023], Sharma et al. [2023]). To address these, several recent works (Chen et al. [2023], Shinn et al. [2024], Gou et al. [2023], Welleck et al. [2022], Summers et al. [2023]) have proposed providing verbal feedback or grounding from an automatic tool or environment to the LLM during the refinement stage. However, these approaches often fail to assess or improve on quality aspects beyond mere functionality when generating code.

Inspired by this progress, we propose PERFCODEGEN, a framework that leverages code execution feedback in the iterative self-refinement stages of an LLM to improve the performance of generated code beyond merely ensuring functional correctness. First, we use environment feedback based on unit test execution to self-refine code generated from the base LLM for correctness. This yields a larger set of correct solutions for our framework to optimize, increasing the likelihood of generating an optimally efficient program in the subsequent optimization phase. For performance refinement, we first assess the runtime required by correct solutions for each unit test, and then provide verbalised performance feedback indicating the most expensive unit test encountered during execution for the LLM to optimize its program. We evaluate PERFCODEGEN’s effectiveness in generating runtime-efficient code on tasks from three widely used Python programming benchmarks: HumanEval (Chen et al. [2021]), MBPP (Austin et al. [2021]) and APPS Hendrycks et al. [2021]. We use five open and two closed language models of varying sizes (Phi-3-mini 3.8B, Llama 3 8B, Mixtral 8x7B (13B), Command R 35B, Llama 3 70B, GPT-3.5 and GPT-4), and witness consistent and significant gains in the correctness and runtime efficiency of LLM-generated programs by applying PERFCODEGEN.

The remainder of this work is organized as follows: In Section 2, we provide a detailed explanation of the proposed PERFCODEGEN framework. Section 3 presents experimental results demonstrating its effectiveness, accompanied by ablations comparing it with several other prompting approaches in Section 3.3. In Section 4, we discuss related work, followed by a brief conclusion in Section 5.

2 PERFCODEGEN Methodology

We begin by prompting a given LLM with a problem description x in natural language that details the requirements of the program to be generated. We sample K candidate generations $C = \{y_x^i\}_{i=1\dots K}$ with nucleus sampling (using the base prompt from Figure 2). As discussed in Section 1, we use the execution feedback to refine the incorrect programs within C in order to increase the number of correct programs in this initial seed set for further optimization. We describe the correctness enhancements in Section 2.1, and the performance refinement phase in Section 2.2. The resulting PERFCODEGEN framework with its correctness and performance phases is illustrated in Figure 1.

2.1 Generating Correct Code

We assess the correctness of LLM generated programs using a set of J unit tests $U(x) = \{u_x^j\}_{j=1}^J$ corresponding to a problem x . After assessing correctness of all candidates, we iteratively refine incorrect ones based on execution feedback from the unit tests. For any $y_x^i \in C$ that fails any of the unit tests, we prompt the LLM again, this time incorporating the environment feedback for correctness verbalised as part of the prompt along with the failed solution and one of the failing unit test. The LLM is instructed to reflect and plan, and then generate a refined correct solution based on the reflection and planning result. The specific prompts we use are provided in Figure 3.

The inclusion of a reflection-planning phase, as suggested by prior work Wei et al. [2022b], Nye et al. [2021], increases the likelihood of the LLM generating a correct solution. The final refinement of y_x^i generated by the LLM is then tested for correctness, and is passed for the performance optimisation stage if it passes all the unit tests. Otherwise, the problem x is considered as unsolved for correctness by the given LLM. Subsequently, a set C_{correct} is constructed by removing incorrect samples from C and including their refined versions, if any of them achieve correctness. We could continue this iterative approach to further improve the correctness rate of a given LLM, but to minimize computational costs, we pause after this one iteration. Besides, we observe that the gains in correctness significantly diminish after the first iteration of refinement. Note that this stage is shared across all the different performance refinement methods studied in this work. Having a larger number of correct candidates as seeds for performance refinements benefits the framework’s effectiveness, as this implies higher likelihood of PERFCODEGEN generating an optimally efficient program.

2.2 Optimising Correct Code using Execution Feedback

For all correct solutions $y_x^i \in C_{\text{correct}}$ that are constructed using the samples and refinements as described in Section 2.1, we prompt the model to refine its solution to optimise performance while preserving functional equivalence. If this refinement breaks correctness, we stop and return the fastest program from C_{correct} . Otherwise, we measure the execution time consumed by this initial refinement (still denoted by y_x^i for simplicity) to pass each available unit test u_x^j corresponding to the given problem x . This involves conducting E independent executions for each solution-test pair in identical compute environments. After sorting this set of E observations, let $t(y_x^i, u_x^j)[e]$ be the e -th smallest execution time consumed by y_x^i on the j -th unit test u_x^j . We then calculate the empirical estimate of the expected execution time of a solution on a unit test, excluding the two outliers (smallest and largest execution times) to minimize the impact of potential measurement noise as follows:

$$\hat{t}(y_x^i, u_x^j) = \frac{1}{(E-2)} \sum_{e=2}^{E-1} t(y_x^i, u_x^j)[e]; \quad u_x^f = \operatorname{argmax}_{u_x^j} \hat{t}(y_x^i, u_x^j) \quad (1)$$

We hypothesize that the f -th unit test u_x^f of a problem x , that corresponds to the largest execution time, as defined above, can be highly informative in optimising the performance of y_x^i if included in the feedback to the LLM for generating a revision. This approach mirrors how developers would identify the most time consuming pieces (hot spots) in their code to come up with runtime improving code changes. Therefore, we re-prompt the model with its latest generation y_x^i , its most time consuming

unit test u_x^f and an instruction to optimise the performance given this feedback (detailed in Figure 7). This leads to a refinement denoted by \tilde{y}_x^i . The fastest amongst the refined correct outputs ($\{\tilde{y}_x^i \mid \tilde{y}_x^i \text{ passes correctness}\}_{i=1\dots K}$) is considered as the final performant solution for x . We employ the greedy decoding algorithm (sampling temperature set to 0) in this stage of performance refinement, as we intend to collect only one refinement per correct code piece to minimize LLM inference costs. If none of the refinement \tilde{y}_x^i pass correctness, we fall back to the fastest correct base solution from C_{correct} , where fastest from C_{correct} is also found using E executions on all unit tests. Similar to optimising for correctness, we could continue refining for performance with more iterations, but we pause after this one iteration for algorithmic simplicity and lower inference costs. Note that the LLM self-refinement step for runtime-efficiency is planned both before and after environment execution. To maximize the utility of execution feedback, we aim to expose to the environment a program that has undergone sufficient correctness and performance based self-refinement by the LLM. Ideally one could also include a planning substep in this self-refinement stage. However, to stay within the token limits of most LLMs, we omit this intermediate step.

3 Experiments

We describe the experiments demonstrating the effectiveness of PERFCODEGEN for generating runtime-efficient programs in this Section. Section 3.1 outlines our experimental setup, Section 3.2 provides the main results of with all the models on the three benchmarks. In Section 3.3 we compare PERFCODEGEN’s execution feedback scheme with alternative prompting strategies, and in Section 3.4 we validate the effectiveness of the planning phase in improving correctness.

3.1 Setup: Metrics, Datasets and Models

To evaluate the correctness and runtime-efficiency of LLM generated solutions using different approaches, we follow prior work (Madaan et al. [2023]) to compute and report the below metrics using the fastest (Best@ k) LLM generated correct program out of k samples.

- **Percent Optimized [%Opt]:** The proportion of problems in the test set where the fastest correct LLM generated program y_x is more runtime-efficient (at least 10% faster) than the ground truth g_x . $\frac{100}{N} \cdot \sum_x \mathbb{1}_{\sum_j \hat{t}(y_x, u_x^j) < 0.9 \cdot \sum_j \hat{t}(g_x, u_x^j)}$ where N is the total number of test set problems.
- **Percent Correct [%Correct]:** The proportion of problems in the test set where the LLM generates at least one correct solution out of k candidates.
- **Speedup:** For problems where we obtain atleast one correct LLM-generated program y_x , we calculate speedup as the absolute ratio between the execution time required by g_x and the execution time required by the fastest correct y_x . $\frac{\sum_x \sum_j \hat{t}(g_x, u_x^j)}{\sum_x \sum_j \hat{t}(y_x, u_x^j)}$

Following prior work (Singhal et al. [2024]), we rely on an empirical estimation of the execution time of Python programs, despite its drawbacks and challenges like high compute requirements. While tools like the gem5 simulator (Binkert et al. [2011]) for reliably and efficiently determine CPU cycles of a program, adapting them for Python is non-trivial. Nevertheless, our qualitative analysis (Listing 1) confirms that the differences observed in execution time (\hat{t}) correspond to clear differences in coding patterns. To estimate the execution time of a candidate solution, we use $E = 12$ executions for each unit test as described in Equation 1. We then compute the above three metrics using the fastest correct program y_x obtained from k (Best@ k) candidates. If there are multiple ground truth solutions (like in APPS), we only use the fastest one as the reference for computing all metrics. We study the impact of sampling budget on the effectiveness of our framework by using $k \in \{1, 8, 20\}$.

We perform our analysis by treating %Opt as the preferred metric over speedup when comparing different methods. A method achieving higher %Opt would be generally more desirable to users than one with lower %Opt, irrespective of speedups observed. This is because users often prefer a larger number of tasks solved optimally rather than a few tasks solved more optimally. However, this preference can be adjusted based on user requirements in different contexts. Speedup should only be analyzed in conjunction with %Opt and %Correct, not in isolation, as it is defined using problems with a correctly generated program by the LLM. Note that the %Correct metric is equivalent to the commonly reported pass@ k , when $n = k$. However, since our approach leverages unit tests in

Model	Method	Best@1			Best@8		
		%Opt	%Correct	Speedup	%Opt	%Correct	Speedup
HumanEval							
Phi-3-mini	Base	14.63	51.83	1.24	35.98	78.05	1.44
	PERFCODEGEN	18.29(+3.66)	57.32(+5.49)	1.23	40.85(+4.87)	85.37(+7.32)	1.68
Mixtral-8x7B	Base	9.43	27.04	1.31	19.5	63.52	1.37
	PERFCODEGEN	11.32(+1.89)	32.70(+5.66)	1.31	27.67(+8.17)	75.47(+11.95)	1.71
Command R	Base	19.02	54.6	1.37	25.15	71.17	1.43
	PERFCODEGEN	20.25(+1.23)	57.06(+2.46)	1.37	32.52(+7.37)	79.75(+8.58)	1.46
Llama 3 8B	Base	15.85	56.71	1.35	29.88	76.22	1.39
	PERFCODEGEN	17.07(+1.22)	62.80(+6.09)	1.29	31.10(+1.22)	81.71(+5.49)	1.62
Llama 3 70B	Base	24.39	75.61	1.39	33.54	84.15	1.42
	PERFCODEGEN	25.61(+1.22)	84.76(+9.15)	1.62	39.02(+5.48)	93.29(+9.14)	1.71
GPT-3.5	Base	16.05	54.94	1.37	29.63	78.4	2.34
	PERFCODEGEN	22.22(+6.17)	67.28(+12.34)	1.34	38.89(+9.26)	90.12(+11.72)	2.33
GPT-4	Base	24.54	72.39	1.81	39.26	88.96	1.82
	PERFCODEGEN	28.83(+4.29)	87.12(+14.73)	1.68	46.63(+7.37)	94.48(+5.52)	1.91
MBPP (test)							
Phi-3-mini	Base	20.26	61.21	2.61	38.36	75.86	3.16
	PERFCODEGEN	20.26(+0.00)	69.40(+8.19)	2.50	44.40(+6.04)	84.48(+8.62)	3.03
Mixtral-8x7B	Base	11.26	45.95	1.46	22.07	66.67	1.63
	PERFCODEGEN	13.06(+1.8)	53.15(+7.2)	1.34	38.29(+16.22)	78.38(+11.71)	2.83
Command R	Base	15.86	55.51	1.68	25.11	69.16	2.27
	PERFCODEGEN	17.18(+1.32)	56.83(+1.32)	1.73	31.72(+6.61)	75.33(+6.17)	2.78
Llama 3 8B	Base	16.02	59.74	2.21	28.14	71.86	2.23
	PERFCODEGEN	17.75(+1.73)	68.40(+8.66)	1.91	32.90(+4.76)	82.68(+10.82)	3.03
Llama 3 70B	Base	18.92	65.32	1.44	26.13	77.93	1.63
	PERFCODEGEN	17.12(-1.8)	68.47(+3.15)	1.68	34.68(+8.55)	88.29(+10.36)	2.21
GPT-3.5	Base	44.1	66.38	3.42	57.21	76.86	3.69
	PERFCODEGEN	47.16(+3.06)	73.36(+6.98)	4.19	71.18(+13.97)	88.21(+11.35)	4.13
GPT-4	Base	20.87	72.17	2.76	43.48	82.17	2.85
	PERFCODEGEN	30.87(+10)	86.52(+14.35)	2.70	56.52(+13.14)	93.91(+11.74)	4.01

Table 1: %Correct, %Opt, and Speedup results on HumanEval and MBPP with k of 1 and 8. Gains in %Opt and %Correct observed by applying our PERFCODEGEN framework are indicated in the relevant cells (+ δ , - δ) for each model. PERFCODEGEN generates correct and runtime-efficient programs for more problems than the base LLM. Despite the higher correctness rate, PERFCODEGEN maintains similar speedup, i.e. speedup over more samples, demonstrating its ability to generate runtime-efficient solutions for more challenging problems that the base LLM could not solve optimally.

execution feedback, it is not fair to compare our correctness metric with those from previous works (Shinn et al. [2024]) that do not assume access to unit tests, and instead reserve them for correctness evaluation. Our study presents an alternative formulation where we seek to generate a program with optimal runtime efficiency, given all unit tests for a task. Our proposed framework can nonetheless be applied in settings we do not have any unit tests, by generating synthetic unit tests for a problem using LLMs (Chen et al. [2022], Shinn et al. [2024]). For simplicity, we instead re-purpose the HumanEval, MBPP and APPS datasets which include tests commonly used to evaluate LLMs on programming tasks. Appendix B.1 details our pre-processing and dataset sanitisation steps. Given our focus on Python datasets, we are unable to compare with Madaan et al. [2024]’s CodeLlama 13B fine-tuned for generating C++ performance improvements. We include open LLMs of varying sizes: Phi-3-mini 3.8B (Abdin et al. [2024]), Llama 3 8B (Meta), Mixtral 8x7B (13B active params, Jiang et al. [2024]), Command R (Cohere), Llama 3 70B (Meta). We also include the closed and commercial GPT-3.5 (Ouyang et al. [2022]) and GPT-4 models via OpenAI APIs. We provide the details of our OpenAI API usage, LLM inference, and code execution environments in Appendix B.2.

3.2 PERFCODEGEN Results

We report the performance of different LLMs on problems from the the HumanEval and MBPP benchmarks in Table 1 using our PERFCODEGEN framework and the aforementioned metrics, with a sampling budget of $k = 1$ and 8 samples. For comparison, we also list the base LLM performance without using PERFCODEGEN. Performance results with a k of 20 are provided in Table 6 in the Appendix. On both the benchmarks, we witness that our framework leads to significant improvements in %Opt and %Correct for all base LLMs, particularly with the higher k of 8.

With PERFCODEGEN we notably enhance the runtime-efficiency of programs generated by open models like Phi-3-mini (%Opt of 40.85 @ $k = 8$) and Llama 3 70B (39.02) making them comparable to GPT-4 in the base setting, which attains the highest base %Opt of 39.26 ($k = 8$). Similarly, with PERFCODEGEN, open models like Mixtral (27.67), Command R (32.52) and Llama 3 8B (31.10) can achieve comparable %Opt to GPT-3.5 (29.63). While we can elevate the performance of open models to match that of closed commercial models, we witness even higher gains in %Opt and %Correct when using PERFCODEGEN on the closed GPT-3.5 and GPT-4 models. This finding can be attributed to the differences in the reasoning capabilities of these model categories, as the effectiveness of PERFCODEGEN heavily relies on the reasoning skills of the given LLM.

On MBPP, we continue to witness significant gains in %Opt when using PERFCODEGEN in most cases when sampling 8 candidates. An exception is Llama 3 70B, whose performance marginally drops on MBPP with our method at $k = 1$, likely due to the high variance in estimating %Opt with a single sample. This drop can be mitigated in practice by leveraging execution time evaluation to fall back to the base LLM output if our refinement is correct but suboptimal. However, we avoid doing so here for a stricter evaluation of our scheme. We provide an example in Listing 1 where PERFCODEGEN generates a faster program than the ground truth. As shown in Table 1, most LLMs generate solutions that are faster than the ground truth for significant portions of the test set, raising questions about their optimality.

```
1 '''
2 Problem: Write a python function to find the average of cubes of first
      n natural numbers.
3 '''
4 # Optimal solution generated by \tool based on GPT-4:
5 def solution(n):
6     sum_of_cubes = (n * (n + 1) / 2.0) ** 2
7     return sum_of_cubes / n
8
9 # Original ground truth solution in MBPP:
10 def solution(n):
11     sum = 0
12     for i in range(1, n + 1):
13         sum += i * i * i
14     return round(sum / n, 6)
```

Listing 1: An optimal solution generated by PERFCODEGEN in MBPP.

Results on the APPS dataset are reported in Table 2. We observe a significantly lower correctness rate compared to HumanEval and MBPP with all LLMs, which is expected given the higher difficulty of APPS. Additionally, since each problem in this dataset has on average more than 25 ground truth solutions, and we use the best one as the reference for comparison, most models fail to generate a solution more optimal than the best ground truth reference for over 1% of the dataset with a k of 1. Using PERFCODEGEN’s execution feedback, we observe gains in %Opt, %Correct and speedup of GPT-3.5 and GPT-4 generations on APPS-test problems, whereas open models benefit less due to task difficulty and their limited reasoning skills. Notably, the gap in correctness rates between commercial GPT models and open models is significantly wider on APPS than on simpler problems from HumanEval and MBPP, highlighting the capability differences.

3.3 Evaluating Alternatives to PERFCODEGEN’s Execution Feedback

Given a correct solution, we evaluate some common prompting techniques for the performance improvement phase. We provide the specific prompts in verbatim in Appendix B.7. Table 3 lists the results with all these methods on the HumanEval and MBPP datasets using GPT-3.5. We exclude

Model	Method	Best@1			Best@8		
		%Opt	%Correct	Speedup	%Opt%	%Correct	Speedup
APPS (test)							
Phi-3-mini	Base	0.09	6.51	1.03	0.28	13.77	1.86
	PERFCODEGEN	0.09(+0.0)	8.40(+1.89)	1.00	0.40(+0.12)	16.73(+2.96)	2.02
Mixtral-8x7B	Base	0.12	6.33	0.97	0.31	11.02	1.08
	PERFCODEGEN	0.19(+0.07)	7.08(+0.75)	1.95	0.40(+0.09)	14.50(+3.48)	1.81
Command R	Base	0.31	14.27	1.00	0.58	24.86	1.00
	PERFCODEGEN	0.43(+0.12)	16.48(+2.21)	1.00	0.86(+0.28)	28.82(+3.96)	1.00
Llama 3 8B	Base	0.18	9.38	2.29	0.52	15.8	3.07
	PERFCODEGEN	0.25(+0.07)	10.21(+0.83)	2.62	0.61(+0.09)	17.71(+1.91)	3.07
Llama 3 70B	Base	0.37	18.65	1.04	0.65	26.12	1.04
	PERFCODEGEN	0.31(-0.06)	19.30(+0.65)	1.04	0.86(+0.21)	27.94(+1.82)	1.04
GPT-3.5	Base	0.49	25.18	1.33	1.11	39.92	1.25
	PERFCODEGEN	0.58(+0.09)	30.56(+5.38)	1.32	1.48(+0.37)	46.26(+6.34)	1.24
GPT-4	Base	0.31	36.63	0.98	1.14	57.75	1.15
	PERFCODEGEN	0.61(+0.30)	45.62(+8.99)	1.26	1.96(+0.82)	65.80(+8.05)	1.90

Table 2: %Correct, %Opt, and Speedup results on APPS with k of 1 and 8. Gains in %Opt and %Correct observed by applying our PERFCODEGEN framework are indicated in the relevant cells (+ δ , - δ) for each model. PERFCODEGEN generates correct and runtime-efficient programs for more problems than the base LLM. Despite the higher relative correctness rate, PERFCODEGEN maintains similar speedup, i.e. speedup over more samples. While the absolute %Opt and %Correct on APPS problems are lower due to the higher difficulty level of problems, PERFCODEGEN offers sizeable relative gains on both these metrics.

other models from this analysis to avoid the high costs associated with LLM inference. Due to the significantly larger size of the APPS-test set compared to HumanEval and MBPP, we exclude it from this comparative analysis to limit inference costs. APPS-test contains roughly 20x more problems, each with approximately three times as many test cases on average, making it prohibitively expensive to perform a comprehensive analysis with multiple LLMs.

We start with three single-round prompting methods. First, we use vanilla prompting for performance improvement, instructing the LLM to optimize the correct code while ensuring functional equivalence (Appendix Figure 4(a)). Next, we evaluate a 3-Shot In-context learning baseline, which includes three examples of program refinements along with optimization instructions (Appendix Figure 4(b)). We also assess an improved prompt that includes common Python optimization tricks as predefined strategies, along with the usual optimization instruction (Appendix Figure 4(c)). We then consider two multi-round approaches similar to those in Madaan et al. [2024], Nye et al. [2021], which include a planning or analysis stage, as shown in Figure 5 in the Appendix. In the Plan and Refine approach, we prompt the LLM to generate a plan for performance refinement, then prompt it again with this output to implement the proposed refinement. In Analyze and Refine, we prompt the LLM to first analyze the time complexity, then prompt it again to refine the code based on this analysis.

We also evaluate multi-agent prompting. First, we implement a multi-agent coder-reviewer setup for performance refinement (Figure 8), where a coder refines the base solution and a reviewer provides feedback, followed by another refinement attempt based on this feedback. Additionally, we implement a more elaborate variant with leader-coder-reviewer roles, where the three agents take turns planning, refining, and reviewing code (Figures 9 and 10 in the Appendix). Finally, we implement two variants leveraging execution feedback for performance improvements after the LLM attempts to refine the correct code solution. In the first variant (Figure 6), we execute and evaluate the effectiveness of the refinement and verbalize the result (positive if the refinement is faster than the base, negative otherwise), feeding this feedback to the LLM for another refinement attempt. In PERFCODEGEN, as described in Section 2, we provide the most time-consuming unit test as feedback to the LLM.

From Table 3, we observe that the base model generations offer significant performance optimizations over the ground truth. However, the gains with performance-improvement prompting, multi-round, and multi-agent techniques are insignificant and often negative. This can be attributed to cascading LLM reasoning errors, as discussed in Section 1. Direct execution feedback produces mixed results,

Prompting Method Summarized Instruction(s)	HumanEval		MBPP (test)	
	Speedup	%Opt	Speedup	%Opt
Base LLM Generation GPT-3.5 prompted to solve a problem	1.37	16.05	3.42	44.1
Perf Improvement Prompt Optimize given code, maintaining equivalence	1.39	19.14(+3.09)	3.45	44.1(+0.00)
In-context learning + Here’s an example of optimisation: <i>{demo}</i>	1.38	19.14(+3.09)	3.20	44.98(+0.88)
Pre-defined Strategies + Here are common ways to optimize: <i>{strategies}</i>	1.27	16.67(+0.62)	2.86	32.75(-11.35)
Plan and Refine (a) Generate optim. plan (b) Optimize w.r.t. plan	1.34	19.14(+3.09)	3.30	45.85(+1.75)
Analyze and Refine (a) Analyze \mathcal{O} time (b) Optimize given a’s analysis	1.35	18.52(+2.47)	3.32	46.29(+2.19)
Multi-Agent w/ Reviewer (a) Coder: Optimize (b) Reviewer: Suggest changes	1.25	18.52(+2.47)	3.57	41.48(-2.62)
Multi-Agent w/ Team (a) Leader: Plan optim (b) Coder Optimize w.r.t. plan (c) Reviewer: Suggest changes to b	1.28	20.37(+4.32)	3.20	42.79(-1.31)
Direct Execution Feedback (a) Optimize (b) It worked/didn’t, try again	1.38	21.60(+5.55)	3.58	42.79(-1.31)
PERFCODEGEN (a) Optimize given code (b) Your costliest unit test is <i>{test}</i> , optimize accordingly	1.34	22.22(+6.17)	4.19	47.16(+3.06)

Table 3: Comparison of PERFCODEGEN Alternatives in the Performance-Refinement Stage ($k = 1$) using GPT-3.5. Gains in %Opt are denoted by $+\delta$ or $-\delta$ in the respective cells. PERFCODEGEN’s execution feedback demonstrates the highest gains while maintaining a similar speedup.

Prompt	GPT-3.5			Llama 3 70B		
	Best@1	Best@8	Best@20	Best@1	Best@8	Best@20
HumanEval						
Base	54.94	78.40	84.57	75.61	84.76	89.02
Testcase Feedback	72.84	85.80	90.74	84.76	90.24	92.07
PERFCODEGEN	68.52	90.12	93.83	85.37	93.90	95.12
MBPP (test)						
Base	66.38	76.86	79.48	65.32	77.93	81.53
Testcase Feedback	78.60	90.83	91.27	72.97	87.39	91.89
PERFCODEGEN	73.36	88.21	92.58	68.47	88.29	93.24

Table 4: %Correct using GPT-3.5 and Llama 3 70B at k of 1, 8 and 20.

with a %Opt gain of 5.55% on HumanEval and a drop of 1.31% on MBPP. In contrast, PERFCODEGEN results in substantial gains in optimization rate on HumanEval (6.17% gain in %Opt) and MBPP (3.06% gain in %Opt) problems, validating the higher performance-improvement effectiveness of PERFCODEGEN’s execution feedback, in the form of the most expensive unit test.

3.4 Role of Planning in Correctness Refinement

We evaluate the effectiveness of PERFCODEGEN’s planning step in achieving higher correctness compared to directly using execution feedback. As discussed in Section 2.1, a high correctness rate is

essential for generating optimized solutions effectively across a larger proportion of test set problems. PERFCODEGEN is more likely to produce maximally optimal solutions by refining from a larger pool of correct candidate solutions, benefiting from the greater diversity within the seed set.

We implement the (direct) Testcase Feedback approach as follows: starting with the base LLM generation, we evaluate correctness based on available unit tests and instruct the LLM to refine its solution according to the environment output (Appendix Figure 3(b)). This process mirrors how a developer would verify correctness by executing a program to ensure it passes the test suite. In contrast, PERFCODEGEN incorporates an additional planning step based on verbalized environment output. The generated plan is then included in the prompt for refinement in the subsequent step.

Results with these two approaches are shown in Table 4 on the HumanEval and MBPP datasets with GPT-3.5 and Llama 3 70B. Both approaches achieve higher correctness than the base LLM. With a k of 1, we observe that the additional planning step of PERFCODEGEN often leads to slightly lower correctness gain compared to the direct approach. However, with a k of 8 and 20, we generally observe higher correctness rate with the planning step. Notably, PERFCODEGEN achieves the highest correctness rates on both benchmarks, underscoring the utility of its additional planning step, suggesting this should also be included in the performance phase given additional token budget.

4 Related Work

Many software engineering works have proposed a code-to-code editing formulation for improving code quality in the form of tasks like fixing bugs (Gupta et al. [2017]), performance improving edits (Madaan et al. [2023], Garg et al. [2022]), improving maintainability (Loriot et al. [2022], Al Madi [2022]), and security enhancing edits (He and Vechev [2023], Perry et al. [2023], Tony et al. [2023], Pearce et al. [2022], Bhatt et al. [2023]). Contrary to this approach, we formulate a text-to-code task for our work on runtime efficiency aspect of quality improvements. As programmers continue to rely on prompting LLMs for generating programs for repetitive tasks in software engineering (Dohmke et al. [2023], Feng and Chen [2024], White et al. [2023], Denny et al. [2023]), we opine that it is critical for research on code quality to focus on the prompting stage by studying natural language inputs that describe developer intent or program specifications.

To improve the general LLM output quality (Chiang et al. [2024]) post the pre-training and supervised instruction fine-tuning (Wei et al. [2022a]) stages, recently proposed algorithms like Reinforcement Learning from Human Feedback (Ouyang et al. [2022], Stiennon et al. [2020]) and Direct Preference Optimization (Rafailov et al. [2024]) that use human preference data have become industry standard (Touvron et al. [2023]). While one could continue to scale these approaches for improving LLM generated code quality, this would require gathering large-scale preference data for code, which is arguably more difficult and expensive than collecting natural language response preferences. Besides needing extensive number of samples, RL techniques also involve expensive model fine-tuning and are known to be notoriously prone to training instabilities (Casper et al. [2023], Wang et al. [2024]).

Our work builds upon the effectiveness of prior work like Madaan et al. [2024]’s Self-Refine, Nye et al. [2021]’s Scratchpads for LLMs and Chen et al. [2023]’s Self-Debug who propose LLM based self-refinements to improve output quality by adding intermediate planning or analysis stages. Our framework is also closely related to Shinn et al. [2024]’s Reflexion who use environment or tool feedback to improve LLM output quality, but focus only on functional correctness in the context of code generation. With PERFCODEGEN, we extend these ideas to improve program runtime efficiency, an aspect that has been largely ignored in favor of functional correctness by prior work.

5 Conclusion

We introduced PERFCODEGEN, a novel framework that leverages code execution feedback during the iterative self-refinement stages of LLMs to enhance the runtime-efficiency of generated code. We show that using our approach, open LLMs like Phi-3-mini can achieve code quality comparable to closed LLMs like GPT-4. Our evaluation of PERFCODEGEN on three widely used Python programming benchmarks using both open and closed language models of varying sizes, demonstrates consistent and significant gains in correctness and runtime efficiency. On a sizeable fraction of the test set problems from HumanEval and MBPP, we achieve programs with state-of-the-art runtime efficiency, far exceeding the ground truth reference solutions in several cases with PERFCODEGEN using GPT-4.

These findings underscore the importance of integrating execution feedback into the code generation process, highlighting a path forward for more robust and reliable AI-driven software development.

References

- Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat Behl, et al. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*, 2024.
- Naser Al Madi. How readable is model-generated code? examining readability and visual inspection of github copilot. In *Proceedings of the 37th IEEE/ACM International Conference on Automated Software Engineering*, pages 1–5, 2022.
- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*, 2021.
- Manish Bhatt, Sahana Chennabasappa, Cyrus Nikolaidis, Shengye Wan, Ivan Evtimov, Dominik Gabi, Daniel Song, Faizan Ahmad, Cornelius Aschermann, Lorenzo Fontana, et al. Purple llama cyberseceval: A secure coding benchmark for language models. *arXiv preprint arXiv:2312.04724*, 2023.
- Nathan Binkert, Bradford Beckmann, Gabriel Black, Steven K Reinhardt, Ali Saidi, Arkaprava Basu, Joel Hestness, Derek R Hower, Tushar Krishna, Somayeh Sardashti, et al. The gem5 simulator. *ACM SIGARCH computer architecture news*, 39(2):1–7, 2011.
- Stephen Casper, Xander Davies, Claudia Shi, Thomas Krendl Gilbert, Jérémy Scheurer, Javier Rando, Rachel Freedman, Tomasz Korbak, David Lindner, Pedro Freire, et al. Open problems and fundamental limitations of reinforcement learning from human feedback. *arXiv preprint arXiv:2307.15217*, 2023.
- Bei Chen, Fengji Zhang, Anh Nguyen, Daoguang Zan, Zeqi Lin, Jian-Guang Lou, and Weizhu Chen. Codet: Code generation with generated tests, 2022. URL <https://arxiv.org/abs/2207.10397>.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Pondé de Oliveira Pinto, Jared Kaplan, Harrison Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Joshua Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large language models trained on code. *CoRR*, abs/2107.03374, 2021. URL <https://arxiv.org/abs/2107.03374>.
- Xinyun Chen, Maxwell Lin, Nathanael Schärli, and Denny Zhou. Teaching large language models to self-debug. *arXiv preprint arXiv:2304.05128*, 2023.
- Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li, Dacheng Li, Hao Zhang, Banghua Zhu, Michael Jordan, Joseph E. Gonzalez, and Ion Stoica. Chatbot arena: An open platform for evaluating llms by human preference, 2024.
- Lawrence Chung, Brian A Nixon, Eric Yu, and John Mylopoulos. *Non-functional requirements in software engineering*, volume 5. Springer Science & Business Media, 2012.
- Paul Denny, Viraj Kumar, and Nasser Giacaman. Conversing with copilot: Exploring prompt engineering for solving cs1 problems using natural language. In *Proceedings of the 54th ACM Technical Symposium on Computer Science Education V. 1*, pages 1136–1142, 2023.

- Thomas Dohmke, Marco Iansiti, and Greg Richards. Sea change in software development: Economic and productivity analysis of the ai-powered developer lifecycle. *arXiv preprint arXiv:2306.15033*, 2023.
- Shihan Dou, Yan Liu, Haoxiang Jia, Limao Xiong, Enyu Zhou, Wei Shen, Junjie Shan, Caishuang Huang, Xiao Wang, Xiaoran Fan, Zhiheng Xi, Yuhao Zhou, Tao Ji, Rui Zheng, Qi Zhang, Xuanjing Huang, and Tao Gui. StepCoder: Improve code generation with reinforcement learning from compiler feedback, 2024.
- Sidong Feng and Chunyang Chen. Prompting is all you need: Automated android bug replay with large language models. In *Proceedings of the 46th IEEE/ACM International Conference on Software Engineering*, pages 1–13, 2024.
- Spandan Garg, Roshanak Zilouchian Moghaddam, Colin B Clement, Neel Sundaresan, and Chen Wu. Deepperf: A deep learning-based approach for improving software performance. *arXiv preprint arXiv:2206.13619*, 2022.
- Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen, Yujiu Yang, Nan Duan, and Weizhu Chen. Critic: Large language models can self-correct with tool-interactive critiquing. *arXiv preprint arXiv:2305.11738*, 2023.
- Rahul Gupta, Soham Pal, Aditya Kanade, and Shirish Shevade. Deepfix: Fixing common c language errors by deep learning. In *Proceedings of the aaai conference on artificial intelligence*, 2017.
- William Harding and Matthew Kloster. Coding on copilot: 2023 data suggests downward pressure on code quality. https://www.gitclear.com/coding_on_copilot_data_shows_ais_downward_pressure_on_code_quality, 2023. Accessed: 2024-05-20.
- Jingxuan He and Martin Vechev. Large language models for code: Security hardening and adversarial testing. In *Proceedings of the 2023 ACM SIGSAC Conference on Computer and Communications Security*, pages 1865–1879, 2023.
- Dan Hendrycks, Steven Basart, Saurav Kadavath, Mantas Mazeika, Akul Arora, Ethan Guo, Collin Burns, Samir Puranik, Horace He, Dawn Song, and Jacob Steinhardt. Measuring coding challenge competence with apps. *NeurIPS*, 2021.
- Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiawu Zheng, Yuheng Cheng, Ceyao Zhang, Jinlin Wang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, Chenyu Ran, Lingfeng Xiao, Chenglin Wu, and Jürgen Schmidhuber. MetaGPT: Meta programming for a multi-agent collaborative framework, 2023.
- Jie Huang, Xinyun Chen, Swaroop Mishra, Huaixiu Steven Zheng, Adams Wei Yu, Xinying Song, and Denny Zhou. Large language models cannot self-correct reasoning yet. *arXiv preprint arXiv:2310.01798*, 2023.
- Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*, 2023.
- Philippe Laban, Lidiya Murakhovska, Caiming Xiong, and Chien-Sheng Wu. Are you sure? challenging llms leads to performance drops in the flipflop experiment. *arXiv preprint arXiv:2311.08596*, 2023.
- Dieter Landes and Rudi Studer. The treatment of non-functional requirements in mike. In *European software engineering conference*, pages 294–306. Springer, 1995.
- Hung Le, Yue Wang, Akhilesh Deepak Gotmare, Silvio Savarese, and Steven Chu Hong Hoi. Coderl: Mastering code generation through pretrained models and deep reinforcement learning. *Advances in Neural Information Processing Systems*, 35:21314–21328, 2022.

- Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, et al. Competition-level code generation with alphacode. *Science*, 378(6624):1092–1097, 2022.
- Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. Is your code generated by chatgpt really correct? rigorous evaluation of large language models for code generation. *Advances in Neural Information Processing Systems*, 36, 2024.
- Benjamin Lorient, Fernanda Madeiral, and Martin Monperrus. Styler: learning formatting conventions to repair checkstyle violations. *Empirical Software Engineering*, 27(6):149, 2022.
- Aman Madaan, Alexander Shypula, Uri Alon, Milad Hashemi, Parthasarathy Ranganathan, Yiming Yang, Graham Neubig, and Amir Yazdanbakhsh. Learning performance-improving code edits. *CoRR*, abs/2302.07867, 2023. doi: 10.48550/ARXIV.2302.07867. URL <https://doi.org/10.48550/arXiv.2302.07867>.
- Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. Self-refine: Iterative refinement with self-feedback. *Advances in Neural Information Processing Systems*, 36, 2024.
- Erik Nijkamp, Bo Pang, Hiroaki Hayashi, Lifu Tu, Huan Wang, Yingbo Zhou, Silvio Savarese, and Caiming Xiong. Codegen: An open large language model for code with multi-turn program synthesis. *ICLR*, 2023.
- Maxwell Nye, Anders Johan Andreassen, Guy Gur-Ari, Henryk Michalewski, Jacob Austin, David Bieber, David Dohan, Aitor Lewkowycz, Maarten Bosma, David Luan, et al. Show your work: Scratchpads for intermediate computation with language models. *arXiv preprint arXiv:2112.00114*, 2021.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.
- Hammond Pearce, Baleegh Ahmad, Benjamin Tan, Brendan Dolan-Gavitt, and Ramesh Karri. Asleep at the keyboard? assessing the security of github copilot’s code contributions. In *2022 IEEE Symposium on Security and Privacy (SP)*, pages 754–768. IEEE, 2022.
- Neil Perry, Megha Srivastava, Deepak Kumar, and Dan Boneh. Do users write more insecure code with ai assistants? In *Proceedings of the 2023 ACM SIGSAC Conference on Computer and Communications Security*, pages 2785–2799, 2023.
- Ruchir Puri, David S Kung, Geert Janssen, Wei Zhang, Giacomo Domeniconi, Vladimir Zolotov, Julian Dolby, Jie Chen, Mihir Choudhury, Lindsey Decker, et al. Codenet: A large-scale ai for code dataset for learning a diversity of coding tasks. *arXiv preprint arXiv:2105.12655*, 2021.
- Chen Qian, Xin Cong, Wei Liu, Cheng Yang, Weize Chen, Yusheng Su, Yufan Dang, Jiahao Li, Juyuan Xu, Dahai Li, Zhiyuan Liu, and Maosong Sun. Communicative agents for software development, 2023.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36, 2024.
- Bernardino Romera-Paredes, Mohammadamin Barekatin, Alexander Novikov, Matej Balog, M Pawan Kumar, Emilien Dupont, Francisco JR Ruiz, Jordan S Ellenberg, Pengming Wang, Omar Fawzi, et al. Mathematical discoveries from program search with large language models. *Nature*, 625(7995):468–475, 2024.
- Inbal Shani. Survey reveals AI’s impact on the developer experience, 2023. URL <https://github.blog/2023-06-13-survey-reveals-ais-impact-on-the-developer-experience/>.

- Mrinank Sharma, Meg Tong, Tomasz Korbak, David Duvenaud, Amanda Askell, Samuel R Bowman, Newton Cheng, Esin Durmus, Zac Hatfield-Dodds, Scott R Johnston, et al. Towards understanding sycophancy in language models. *arXiv preprint arXiv:2310.13548*, 2023.
- Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning. *Advances in Neural Information Processing Systems*, 36, 2024.
- Disha Shrivastava, Hugo Larochelle, and Daniel Tarlow. Repository-level prompt generation for large language models of code. In *International Conference on Machine Learning*, pages 31693–31715. PMLR, 2023.
- Manav Singhal, Tushar Aggarwal, Abhijeet Awasthi, Nagarajan Natarajan, and Aditya Kanade. Nofuneval: Funny how code lms falter on requirements beyond functional correctness. *arXiv preprint arXiv:2401.15963*, 2024.
- Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. *Advances in Neural Information Processing Systems*, 33:3008–3021, 2020.
- Theodore R Sumers, Shunyu Yao, Karthik Narasimhan, and Thomas L Griffiths. Cognitive architectures for language agents. *arXiv preprint arXiv:2309.02427*, 2023.
- Catherine Tony, Markus Mutas, Nicolás E Díaz Ferreyra, and Riccardo Scandariato. Llmseceval: A dataset of natural language prompts for security evaluations. In *2023 IEEE/ACM 20th International Conference on Mining Software Repositories (MSR)*, pages 588–592. IEEE, 2023.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- Yuanhao Wang, Qinghua Liu, and Chi Jin. Is rlhf more difficult than standard rl? a theoretical perspective. *Advances in Neural Information Processing Systems*, 36, 2024.
- Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. Finetuned language models are zero-shot learners. In *International Conference on Learning Representations*, 2022a.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022b.
- Sean Welleck, Ximing Lu, Peter West, Faeze Brahman, Tianxiao Shen, Daniel Khashabi, and Yejin Choi. Generating sequences by learning to self-correct. *arXiv preprint arXiv:2211.00053*, 2022.
- Jules White, Sam Hays, Quchen Fu, Jesse Spencer-Smith, and Douglas C Schmidt. Chatgpt prompt patterns for improving code quality, refactoring, requirements elicitation, and software design. *arXiv preprint arXiv:2303.07839*, 2023.
- John Yang, Carlos E. Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik Narasimhan, and Ofir Press. Swe-agent: Agent-computer interfaces enable automated software engineering, 2024.
- Yuntong Zhang, Haifeng Ruan, Zhiyu Fan, and Abhik Roychoudhury. Autocoderover: Autonomous program improvement. *arXiv preprint arXiv:2404.05427*, 2024.

A Limitations

The challenge of writing efficient and high-quality software with LLMs spans various levels of granularity, from line-level optimizations to multi-class project repositories Shrivastava et al. [2023]. In our current scope, we focus on generating performant modules or Python functions, which are typically small components of real-world systems. However, addressing this problem comprehensively should ideally involve ensuring architectural design patterns such as minimal redundancy or wasteful computation across the entire scope of a project.

Another limitation is our focus solely on measuring the runtime performance of LLM-generated code, disregarding memory consumption, which can be crucial in many applications. Future extensions of PERFCODEGEN could prioritize optimizing for both factors or allow users to specify preferences for optimization. Additionally, beyond performance, developers desire attributes like readability, ease of maintenance, security, and harmlessness Bhatt et al. [2023], Singhal et al. [2024], which are not within the scope of our current work. While PERFCODEGEN could be adapted to incorporate feedback from different environments or tools evaluating these attributes, achieving a balance in optimizing code generation across these dimensions is non-trivial.

As emphasized by prior research, reliably measuring the runtime performance of code poses challenges Madaan et al. [2023]. A piece of code may exhibit varying execution times across different compute environments, even with identical underlying hardware. Unfortunately, tools like the gem5 simulator Binkert et al. [2011] do not support the execution of Python programs at the time of our study. To mitigate this, we ensure identical compute environments for each candidate code snippet and run only a single Python program at any given time to minimize effects from concurrent execution. However, averaging execution time measurements from 10 independent runs of each program significantly adds to our execution costs. Future work could explore more efficient methods for reliably measuring runtime, such as determining the instruction count of LLM-generated programs deterministically.

B Appendix

B.1 Sanitized Benchmarks

Benchmark	Original			Sanitized		
	#Problem	#Groundtruth	#Testcase	#Problem	#Groundtruth	#Testcase
HumanEval	164	1.0	9.6	164	1.0	9.6
MBPP-test	257	1.0	3.0	257	1.0	3.0
APPS-test	5,000	30.4	21.2	3,249	26.9	27.4

Table 5: Details of sanitized benchmarks used in the evaluation. “#Problem” indicates the number of problems in the benchmark, “#Groundtruth” indicates the average number of ground truth programs for each problem, and “#Testcase” indicates the average number of test cases for each problem. We use the sanitized benchmarks for correctness evaluation and the common subsets for time efficiency evaluation.

We utilize three benchmarks for our evaluation: the HumanEval benchmark developed by Chen et al. [2021], the sanitized version of the MBPP benchmark created by Austin et al. [2021], and the APPS benchmark established by Hendrycks et al. [2021]. Detailed information on these benchmarks is presented in Table 5, which we directly obtain from the HuggingFace Hub. In our evaluation process, we initially verify whether the provided ground truth programs within each benchmark can successfully pass the associated test cases. Notably, we identify 1,511 instances in the APPS benchmark where the ground truth solutions fail to meet the test case criteria. This discrepancy arises because our evaluation requires exact matches between program outputs and expected outputs in test cases, whereas the APPS benchmark allows for similar matches and may contain incorrect ground truth programs. This observation aligns with previous findings reported by Dou et al. [2024]. To ensure a fair comparison among the different benchmarks, we exclude the aforementioned 1,751 instances from our analysis, retaining the remaining 3,249 instances for correctness evaluation. Additionally,

for assessing time efficiency, we construct subsets from the three benchmarks, comprising problems for which all baseline methods can generate at least one functionally correct program. Time-related data are computed solely on these subsets, rather than the entire benchmark, to maintain consistency in evaluation conditions.

B.2 Compute

LLM inference: We use the vLLM Kwon et al. [2023] library on a node with 16 Nvidia A100 GPUs for approximately three weeks to complete all the experiments in this work.

OpenAI API: We use the gpt-4-0613 and gpt-3.5-turbo-0125 model endpoints from OpenAI. In total, we required nearly 411k GPT-4 and 1.5M GPT-3 requests for all our experiments, contributing to the major costs of this study.

Code Execution: We use 40 instances of virtual machines (n1-highmem-16 GCP instances), each with 16 CPUs and 104 GB RAM for executing all the LLM generated programs generated in our experiments. We employ these instances for roughly four weeks to complete the execution of all the LLM generated programs using different frameworks in our study. Interestingly, unlike most LLM research, gathering this environment feedback tends to be the much costlier bottleneck in our experiments compared to the LLM inference costs. We implement the safeguards prescribed in Section 2.3 of Chen et al. [2021] to mitigate the security risks in executing untrusted programs in our environment.

B.3 Statistical Significance

Our main results are based on findings in Table 1 and Table 2, where we report that using PERFCODEGEN leads to significant gains in the %Opt metric compared to the base LLM. For results with the GPT-4 model, we compute the Z-scores to compare PERFCODEGEN’s output with that of the base model: $-0.878(k = 1)$ and $-1.34(k = 8)$ on HumanEval, $-2.588(k = 1)$ and $-2.947(k = 8)$ on MBPP and $-1.8(k = 1)$ and $-2.675(k = 8)$ on APPS. The improvement obtained with PERFCODEGEN is thus statistically significant with $\alpha < 0.05$ on the MBPP and APPS problems, and with a lower confidence ($\alpha < 0.2$) on HumanEval which has a smaller number of problems (164).

B.4 Broader Impact

By enabling LLMs to generate code that is not only functionally correct but also efficient, our work with PERFCODEGEN can significantly accelerate the software development process. This can lead to faster creation of new applications, reduced development costs, and increased innovation across various industries. Frameworks like PERFCODEGEN can potentially empower individuals with less coding experience to leverage LLMs for basic programming tasks. This could lead to a wider pool of software developers and a more inclusive tech landscape. More efficient code translates to lower energy consumption during program execution. This can contribute to a more sustainable software development ecosystem and reduce the environmental impact of the tech industry. Our work demonstrates the ability of PERFCODEGEN to enhance the performance of open LLMs, making them more competitive with closed models. This can foster a more open and accessible environment for LLM development and research. Our work could also offer a promising route to discover novel algorithms to solve long standing problems more efficiently (Romera-Paredes et al. [2024]).

While PERFCODEGEN aims to improve code efficiency, it could be misused to automate the generation of efficient harmful or malicious code. For instance, cybercriminals could use optimized code to create more efficient malware or exploit software vulnerabilities more effectively. Incorporating content filtering and malicious code detection algorithms to identify and block harmful code generation can help reduce the risk of such misuse. Optimized code may sometimes introduce new types of vulnerabilities that are difficult to detect. If PERFCODEGEN generates code that is highly efficient but less readable or maintainable, it could lead to challenges in debugging and maintaining software, potentially resulting in unexpected failures or security issues. Addressing this risk requires implementing comprehensive testing and validation, including code reviews, to ensure the generated code maintains high reliability and security standards. More efficient code generation could lead to further automation in software industry, potentially displacing some human programmers. This necessitates discussions on re-skilling initiatives and the evolving nature of jobs in the tech industry.

Model	Speedup	Opt%	Correct%
HumanEval			
Phi3	1.81	47.56	92.68
Mixtral-8x7B	1.85	42.14	87.42
Command R	2.01	46.63	87.12
Llama3 8B	2.27	43.9	87.2
Llama3 70B	1.84	45.73	94.51
GPT-3.5	2.33	43.21	93.83
GPT-4	1.98	55.21	95.71
MBPP (test)			
Phi3	3.40	59.91	89.22
Mixtral-8x7B	3.52	47.3	87.39
Command R	3.21	46.26	85.9
Llama3 8B	3.30	43.72	86.15
Llama3 70B	2.88	44.59	93.24
GPT-3.5	4.25	75.98	92.58
GPT-4	4.16	66.96	95.65
APPS (test)			
Phi3	2.83	0.65	21.46
Mixtral-8x7B	1.75	0.59	21.42
Command R	1.05	1.26	34.53
Llama3 8B	2.86	0.95	22.44
Llama3 70B	1.06	1.39	37.08
GPT-3.5	1.24	2.06	51.34
GPT-4	2.57	2.95	72.18

Table 6: PERFCODEGEN results on HumanEval, MBPP and APPS with a k of 20.

B.5 Best@20 Results of PERFCODEGEN for All Models

B.6 Examples of Optimal Solutions Generated by LLMs

```
1 '''
2 Problem:
3 The FibFib number sequence is a sequence similar to the Fibonacci
   sequence that's defined as follows:
4 fibfib(0) == 0
5 fibfib(1) == 0
6 fibfib(2) == 1
7 fibfib(n) == fibfib(n-1) + fibfib(n-2) + fibfib(n-3).
8 Please write a function to efficiently compute the n-th element of the
   fibfib number sequence.
9 >>> fibfib(1)
10 0
11 >>> fibfib(5)
12 4
13 >>> fibfib(8)
14 24
15 '''
16 # Optimal solution generated by PerfCodeGen based on GPT-3.5
17 def fibfib(n: int):
18     if n == 0 or n == 1:
19         return 0
20     if n == 2:
21         return 1
22     (a, b, c) = (0, 0, 1)
23     for _ in range(3, n + 1):
24         (a, b, c) = (b, c, a + b + c)
25     return c
26 # Original ground truth solution in HumanEval:
27 def fibfib(n):
28     if n == 0:
29         return 0
30     if n == 1:
31         return 0
32     if n == 2:
33         return 1
34     return fibfib(n - 1) + fibfib(n - 2) + fibfib(n - 3)
```

Listing 2: An optimal solution generated by PERFCODEGEN in HumanEval.

B.7 Prompts

We design multiple prompts to evaluate the base performance of LLM and generate runtime-efficient code solutions.

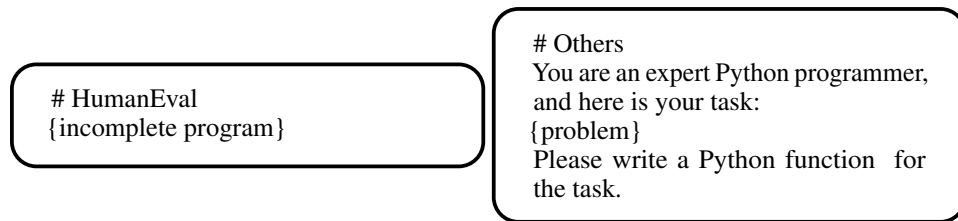


Figure 2: Base prompt for HumanEval and other benchmarks.

Round 1:
Your generated solution for the problem is not correct and cannot pass the following test case:
{testcase}

The error message is as follows:
“*{error}*”

Could you please analyze the reason of failure and propose a strategy to modify your solution so that it can pass the above test case?

Round 2 :
Could you please modify your solution so that it can fulfill the requirements in the problem and pass the test case?
Give your solution as follows. Wrap it with “python”.

(a) Reflection and Test Case Feedback

Your generated solution for the problem is not correct and cannot pass the following test case:
{testcase}

The error message is as follows:
“*{error}*”

Could you please modify your solution so that it can fulfill the requirements in the problem and do not have any syntax error?
Give your solution as follows. Wrap it with “python”.

(b) Test Case Feedback

Figure 3: The two correctness prompts discussed in the paper. PERFCODEGEN uses the reflection and test case feedback prompt in (a).

Good job! You generated the correct solution for the problem! Now let's step further and optimize the time performance of the solution.
Based on the correctly generated solution, could you please refine it so that it consumes less time in the execution?
Please make sure your refined solution is functionally equivalent with the original solution and do not change the input-output format and the name of the major components.
Give your solution as follows. Wrap it with `python`.

(a) Perf Improvement Prompt

Good job! You generated the correct solution for the problem! Now let's step further and optimize the time performance of the solution.
Here are some examples of optimization:
{demo}
Based on the correctly generated solution and the above examples, could you please refine it so that it consumes less time in the execution?
Please make sure your refined solution is functionally equivalent with the original solution and do not change the input-output format and the name of the major components.
Give your solution as follows. Wrap it with `python`.

(b) In-context learning (few-shot)

Good job! You generated the correct solution for the problem! Now let's step further and optimize the time performance of the solution.
Here are some commonly used strategies for optimization:
1. Do not generate useless testing code.
2. Use builtin functions and libraries instead of importing third-party libraries.
3. Use list and dict comprehension to avoid loops.
4. Eliminate unnecessary variable definitions, function definitions and print statements.
5. Optimize the loops to avoid unnecessary iterations.
6. Use local variables instead of global variables.
7. Use multiple assignments in one statement instead of in multiple statements.
8. Make use of generators.
Based on the correctly generated solution and the above strategies, could you please refine it so that it consumes less time in the execution?
Please make sure your refined solution is functionally equivalent with the original solution and do not change the input-output format and the name of the major components.
Give your solution as follows. Wrap it with `python`.

(c) Pre-defined Strategies

Figure 4: The single-round runtime performance improving prompts.

Round 1:
Good job! You generated the correct solution for the problem! Now let's step further and optimize the time performance of the solution.
Could you please propose a strategy to optimize the generated solution so that it consumes less time in the execution?

Round2:
Please implement the strategy you proposed for the previous solution and generate the optimized solution.
Please make sure your refined solution is functionally equivalent with the original solution and do not change the input-output format and the name of the major components.
Give your solution as follows. Wrap it with `“python”`.

(a) Plan and Refine

Round 1:
Good job! You generated the correct solution for the problem! Now let's step further and optimize the time performance of the solution.
Could you please analyze the time complexity of the solution?

Round 2:
Now given the time complexity you estimated, could you please optimize your solution by reducing its time complexity?
Please make sure your refined solution is functionally equivalent with the original solution and do not change the input-output format and the name of the major components.
Give your solution as follows. Wrap it with `“python”`.

(b) Analyze and Refine

Figure 5: The multi-round runtime performance improving prompts.

Round 1: Good job! You generated the correct solution for the problem! Now let's step further and optimize the time performance of the solution.
Based on the correctly generated solution, could you please refine it so that it consumes less time in the execution?
Please make sure your refined solution is functionally equivalent with the original solution and do not change the input-output format and the name of the major components.
Give your solution as follows. Wrap it with “python”.

Round 2 (Negative Feedback):
We tested your optimized solution and found that your optimized solution runs slower than the previous version. The following is the execution time for both solutions:
Original solution:
{ori_time}
Your optimized solution:
{opt_time}
Based on the execution time and the original solution, could you please propose another optimized solution so that it consumes less time than the original solution?
Please make sure your refined solution is functionally equivalent with the original solution and do not change the input-output format and the name of the major components.
Give your solution as follows. Wrap it with “python”.

Round 2 (Positive Feedback):
We test your optimized program. Great! It runs faster than the original program.
The following is the execution time for both solutions:
Original solution:
{ori_time}
Your optimized solution:
{opt_time}
Can you refine the current optimized version as follows and provide a more efficient one?
Please make sure your refined solution is functionally equivalent with the original solution and do not change the input-output format and the name of the major components.
Give your solution as follows. Wrap it with “python”.

Figure 6: Direct Execution Feedback prompt.

Round 1: Good job! You generated the correct solution for the problem! Now let's step further and optimize the time performance of the solution. Based on the correctly generated solution, could you please refine it so that it consumes less time in the execution?
Please make sure your refined solution is functionally equivalent with the original solution and do not change the input-output format and the name of the major components. Give your solution as follows. Wrap it with “python”.

Round 2:
We tested your optimized program and found that the following test case costs the most time in execution.

{testcase}

Could you please refine your optimized program according to the test case below?
Please make sure your refined solution is functionally equivalent with the original solution and do not change the input-output format and the name of the major components. Give your solution as follows. Wrap it with “python”.

Figure 7: PERFCODEGEN Testcase Feedback prompt used.

Round 1 (Coder):
Good job! You generated the correct solution for the problem! Now let's step further and optimize the time performance of the solution.
Based on the correctly generated solution, could you please refine it so that it consumes less time in the execution?
Please make sure your refined solution is functionally equivalent with the original solution and do not change the input-output format and the name of the major components.
Give your solution as follows. Wrap it with “python”.

Round 2 (Reviewer): You are in a discussion group, aiming to optimize a Python program so that it runs faster.
The original Python program is:

{program}

An optimized Python program given by your group member:

{opt_program}

You are a reviewer. Based on your knowledge, can you check whether the optimized program given by your group member runs faster than the original program?
Your response should follow the following rules:
1. Begin your response with [Agree] if you think the optimized program really runs faster than the original one and [Disagree] if not or there exists some syntax errors in the optimized program.
2. Provide your comments after the keyword “Comment:” to explain your decision.
3. Give suggestions for further improvements in comments.
Give your response here:

Round 3 (Coder):
A code reviewer carefully reviewed your optimized code solution. He thinks that your optimization is *{decision}*.
Please refer to his comments as follows to further refine your optimized program.
Comments:

{comment}

Please make sure your refined solution is functionally equivalent with the original solution and do not change the input-output format and the name of the major components.
Give your solution as follows. Wrap it with “python”.

Figure 8: The Multi-Agent (Coder - Reviewer) prompt.

Round 1 (Leader):
You are a team leader and your team is working on the code of the following problem:

{problem}

Here is the correct solution your team wrote before:

{program}

Given the Python program, can you make a plan about how to optimize it step by step?
Give your plan here:

Round 2 (Coder):
Given the previous Python program, your leader has proposed the following plan to optimize it:

{plan}

Following the above plan, can you write an optimized version of the original Python program?
Please make sure your refined solution is functionally equivalent with the original solution and do not change the input-output format and the name of the major components.
Give your solution as follows. Wrap it with “python”.

Round 3 (Reviewer):
You are a code reviewer in a team and your team is trying to optimize the following Python program so that it runs faster.
Original Python program:

{program}

Given the Python program, your leader has proposed the following plan to optimize it:

{plan}

The coder in your team implemented an optimized version of program as follows.

{opt_program}

Based on your knowledge, can you check whether the optimized program given by your group member runs faster than the original program?
Your response should follow the following rules:

1. Begin your response with [Agree] if you think the optimized program really runs faster than the original one and [Disagree] if not or there exists some syntax errors in the optimized program.
2. Provide your comments after the keyword `Comment:` to explain your decision.
3. Give suggestions for further improvements in comments.

Give your response here:

Figure 9: The first part of Multi-Agent w/ Team (Leader - Coder - Reviewer) prompt.

Round 4 (Leader):
A coder in your team implement an optimized version of program according to your plan:
{opt_program}

A code reviewer reviewed the optimized version and thinks that the optimization is *{decision}*.
He provides the following comments:
{comment}

Could you please refine your plan so that we can further improve the efficiency of the program?
Give you plan here:

Round 5 (Coder):
After your optimized program is reviewed by another group member, your leader modifies his plan and gives a new plan here:
{plan}

Could you please write a new optimized version of the program based on the new plan?
Please make sure your refined solution is functionally equivalent with the original solution and do not change the input-output format and the name of the major components.
Give your solution as follows. Wrap it with “python”.

Figure 10: The second part of Multi-Agent w/ Team (Leader - Coder - Reviewer) prompt.