

# Rethinking Code Refinement: Learning to Judge Code Efficiency

Anonymous ACL submission

## Abstract

Large Language Models (LLMs) have demonstrated impressive capabilities in understanding and generating codes. Due to these capabilities, many recent methods are proposed to automatically refine the codes with LLMs. However, we should rethink that the refined codes (from LLMs and even humans) are not always more efficient than their original versions. On the other hand, running two different versions of codes and comparing them every time is not ideal and time-consuming. Therefore, in this work, we propose a novel method based on the code language model that is trained to judge the efficiency between two different codes (generated across humans and machines) by either classifying the superior one or predicting the relative improvement. We validate our method on multiple programming languages with multiple refinement steps, demonstrating that the proposed method can effectively distinguish between more and less efficient versions of code.

## 1 Introduction

Large Language Models (LLMs) have shown significant success across a wide range of tasks, extending from natural language understanding to programming-related activities (Brown et al., 2020; Chen et al., 2021; Achiam et al., 2023; Touvron et al., 2023; Rozière et al., 2023; Abdin et al., 2024). Specifically, thanks to their capabilities in understanding and generating codes, LLMs are able to allow developers to save time, reduce errors, and boost their productivity (Shen et al., 2022). For example, several recent studies have utilized LLMs to optimize and refine the existing code bases (Madaan et al., 2023; Chen et al., 2023b). Also, more recent work has been proposed to iteratively refine the generated codes from LLMs by judging them with LLMs (Zelikman et al., 2023). Another noteworthy approach involves using LLMs to check the functional correctness of the generated codes from LLMs (Dong et al., 2023a).

However, despite huge advancements made in the field of LLMs for code generation, the aforementioned studies assume that the codes generated and refined from LLMs are more efficient than their originals. However, as shown in Figure 1, our observations contradict this assumption, showing that LLM-generated and -refined codes do not always perform better. To handle this issue, while one may calculate the efficiencies of codes both before after modifications by actually executing them, this process introduces unnecessary inefficiencies, and may be very costly and time-consuming.

In this work, to overcome those challenges, we introduce a new task of judging the efficiency of the refined code over its original version. In addition, not only LLMs but also human coders may degrade the efficiency of codes during refinement; thus, our task of judging the efficiency between two codes involves all the possible pairs of code modification sources, including human-human, human-machine, and machine-machine. Then, to address this new task, we propose a new model (based on the code LM) that is trained to compare efficiencies between two codes. Specifically, given a code pair (one is original and the other is refined from it), the code LM is trained to classify which one is more efficient or predict how much the refined code is efficient.

We experimentally validate the effectiveness of our efficiency judgement model on multiple code refinement scenarios with multiple programming languages. The results show that our judgement model can identify the more efficient code among two different versions, substantially outperforming baselines. Our contributions are as follows:

- We point out that the refined codes from LLMs or humans do not always improve their efficiency.
- We introduce a novel approach that judges the efficiency of two different versions of codes.
- We validate our model on multiple code refinement scenarios, demonstrating its effectiveness.

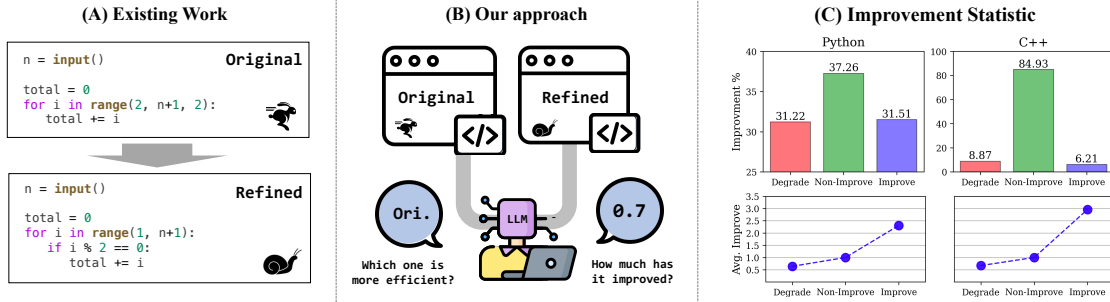


Figure 1: (A) Existing code refinement approaches sometimes generate the code which has inferior efficiency to the original code. (B) Our proposed approach identifies the efficient code among two different versions of codes (before and after modifications), and further predicts its relative improvement. (C) We categorize the refined code according to its efficiency gain (%) compared to the original into three classes: Degradation (less than 0.9), Non-Improvement (0.9 to 1.1), and Improvement (greater than 1.1).

## 2 Related Work

**Large Language Models for Code** Large Language Models (LLMs), trained on extensive corpora with multi-billion parameters, exhibit remarkable performance across a broad spectrum of tasks involving both text and code (Brown et al., 2020; Li et al., 2022; Rozière et al., 2023; Li et al., 2023; Achiam et al., 2023; Abdin et al., 2024; Guo et al., 2024). These models, particularly those trained on code-specific datasets, have opened up a new era in software development by not only assisting with basic programming tasks but also enabling more complex activities such as code generation (Chen et al., 2023a; Nijkamp et al., 2023), translation (Yin and Neubig, 2018; Rozière et al., 2020; Cassano et al., 2023), and refinement (Yu et al., 2023; Shirafuji et al., 2023). As such, these models are increasingly integrated into development environments, optimizing workflows, and reducing the time for development (Shen et al., 2022; Dong et al., 2023b).

**LLMs for Code Refinement** Beyond basic code generation, LLMs are widely used to refine the existing code bases. One of the early work in code refinement aims to detect and fix bugs in the codes by pre-training a transformer-based model on English and source code, and then fine-tuning it on commits (relevant to the part fixing bugs and improving performance) (Garg et al., 2022). In a similar vein, another work proposes to refine the code with a sequence-to-sequence model that is trained to transform the original code to its optimized version of the code (Chen et al., 2023b). Additionally, recent work showcases that LLMs are able to recursively improve their own generated codes, progressively enhancing their outputs (Zelikman et al., 2023). Further, Madaan et al. (2023) demonstrate that LLMs with sophisticated prompting strategies (to adapt LLM for code optimization) can surpass human-level performance in code optimization tasks. However, despite these substantial

achievements, prior studies have primarily focused on code enhancement with limited attention to the actual efficiency of the refined code. Meanwhile, we focus on a different angle, proposing to judge the efficiency of the refined codes in advance (without executing them) based on our observation that not all the refined codes have better efficiency.

**LLM-Powered Code Evaluation** The objective of our work which aims to evaluate the efficiency between two codes based on LLMs has a similarity to work on LLMs for code evaluation. Early work on it uses either a term-matching-based approach (similar to BLEU) or an embedding-based approach (whose representations are obtained from language models), to compare two codes (Ren et al., 2020; Zhou et al., 2023). However, as collecting ground-truth answers for every evaluation is difficult, recent work has shifted towards using LLMs to judge the quality of the generated code, such as its utility or functional correctness, without the need for comparisons to the reference code (Dong et al., 2023a; Zhuo, 2024). Yet, unlike these approaches evaluating the single instance of the code other than the efficiency, we aim to compare efficiency in the setting where the code pair is given.

## 3 Method

We first provide a general description of code refinement, and then introduce our approach.

### 3.1 Code Refinement

Let us assume that the existing code base is defined as  $c$ , which consists of a sequence of tokens as follows:  $c = \{c_1, \dots, c_n\}$ . Then, the objective of code refinement (in this work) is to transform the existing code base  $c$  into its improved version  $c' = \{c'_1, \dots, c'_m\}$ , where the execution time for  $c'$  should be faster than  $c$ , formalized as follows:  $\text{Exec}(c) > \text{Exec}(c')$ . Here,  $\text{Exec}$  is the code execution function that returns its runtime.

Table 1: Main results on the task of judging the code efficiency, where easy denotes the dataset containing only the code pairs whose efficiency difference is more than 10%.

	All			Easy		
	Python	C++	All	Python	C++	All
Zero-shot	50.87	45.30	46.55	47.04	51.21	49.70
Few-shot	51.21	51.17	51.18	49.16	48.22	48.56
Zero-shot CoT	50.35	51.69	51.39	48.94	52.10	50.95
Few-shot CoT	50.52	50.87	50.79	49.50	48.09	48.60
<b>Ours</b>	<b>72.49</b>	<b>62.08</b>	<b>64.42</b>	<b>77.65</b>	<b>70.14</b>	<b>72.86</b>

It is worth noting that, in this work, we consider three different scenarios of code refinement. The first scenario involves a human-human interaction, where one developer revises the code originally authored by another. The second scenario, termed the human-machine scenario, consists of collaborative efforts between humans and machines. This practice has become increasingly prevalent in real-world software development environments, thanks to CodeLLMs (Chen et al., 2021; Rozière et al., 2023). Lastly, the machine-machine scenario, involves autonomous code refinement by machines, a process that has shown promise in various studies (Zelikman et al., 2023; Madaan et al., 2023).

Note that, despite the huge advancements made in the field of code refinement, we find that modified codes from machines can occasionally reduce the efficiency of the original codes. Similarly, human developers may diminish the efficiency of the codes during refinement. On the other hand, executing the pair of original and modified codes at every refinement step is inefficient and time-consuming.

### 3.2 Judging Code Efficiency

To overcome the aforementioned limitation, we aim to predict the efficiency of the modified code over its original code, without actually executing them. This can be formulated by either the classification problem (where we classify the superior code) or the regression problem (where we predict the relative improvement of the modified code over the original code), given a pair of original and modified codes. Also, note that we operationalize classification and regression problems with CodeLLMs, due to their capabilities in understanding codes.

Specifically, given the code pair (e.g.,  $c$  and  $c'$ ), we concatenate it and provide it with the CodeLLM, formalized as follows:  $o = \text{CodeLLM}([c, c'])$  where  $[\cdot]$  is the concatenation operation, and  $o$  is the prediction output. Then, for the classification problem, we formulate it as the next token prediction task, as follows:  $L_C = -\sum_{t=1}^T \log p(o|[c, c'], o \in \{A, B\})$ , where  $A$  represents the improvement over the original code and vice versa for  $B$ . Simi-

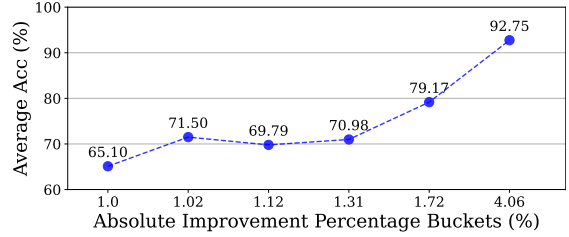


Figure 2: Results with bucketing the code pairs according to their absolute relative improvement in efficiency, on Python.

larly, for the regression problem, we train the model to predict the relative improvement of the refined code over its original by minimizing this prediction value with the actual relative improvement.

## 4 Experiment

We now describe experimental setups and results.

### 4.1 Experimental Setups

**Dataset** To validate the efficacy of our approach to judge code efficiency, we should collect pairs of two different versions of codes before and after modifications. Here, we consider three different scenarios of code editing, and, for the cases where humans refine the code, we use a dataset of code edits made by humans from Madaan et al. (2023). For the other scenarios where the machine improves the human- or machine-generated codes, we prompt the Code LLM (namely DeepSeek-Coder-Instruct-7B) (Guo et al., 2024) to refine the given codes for better efficiency. Specifically, starting with the codes generated by humans from the existing dataset, we generate the machine-refined codes with the Code LLM. In addition, from those machine-generated codes, we similarly prompt the Code LLM to improve them. Through these steps, we can obtain pairs of human-human, human-machine, and machine-machine code versions.

**Baselines and Our Model** In this work, as we tackle a novel problem of judging code efficiency, there are no direct baselines available to compare. Therefore, we turn to compare our approach against the basic models powered by LLMs. Specifically, given a code pair, we perform zero-shot and few-shot prompting with LLMs, to decide which one is more efficient. In addition, we also enhance those strategies with Chain-of-Thought prompting (Wei et al., 2022), to elicit the reasoning ability of LLMs with the instruction: "Let's think step by step". For our model, we use the classifier (predicting the class of the efficient code) for main experiments, and provide the performance of the other (predicting the relative improvement) during analysis. We use DeepSeek-Coder-Instruct-1.3B for all models.

Table 2: Breakdown results for varying the code refinement scenarios. 'H' indicates Human and 'M' indicates Machine.

Scenarios	Statistics			Breakdown Acc		
	Avg. Improve	Degrade %	Improve %	Python	C++	All
H-H	1.08	37.19	21.85	80.43	64.43	67.88
M-M	1.32	30.69	32.25	60.77	55.84	57.03
H-M	1.29	31.26	30.76	67.27	61.55	62.87

Table 3: Results with predicting the relative difference of the modified code over its original code in efficiency. Corr denotes the Spearman’s rank correlation coefficient, and Acc denotes the accuracy where we convert prediction values into classes.

	All		Easy	
	Corr	Acc	Corr	Acc
Python	0.66	76.38	0.64	82.88
C++	0.50	66.69	0.63	80.16
Python & C++	0.56	68.87	0.66	81.17

## 4.2 Experimental Results

**Main Results** We report the main results in Table 1, and, from this, we observe that our method consistently outperforms all baseline models across all settings. In addition, our model is particularly effective in scenarios where there is a clear difference in code efficiency — specifically, a difference exceeding 10% (the easy setting). To examine the performance of our model more granularly based on varying degrees of efficiency differences between code pairs, we bucketize the code pairs based on their efficiency differences. As shown in Figure 2, the performance of our model increases when the difference between two codes becomes larger.

**Results with Varying Refinement Scenarios** It is worth noting that our code refinement scenario is categorized as human-human, human-machine, and machine-machine, and we report their breakdown results in Table 2. From this, we first observe that the percentage of efficiency improvement in code refinement scenarios is low, which is around 20% to 30%, and it is similar to the percentage of efficiency degradation. In addition, the average improvement made by machines is 30%, meanwhile, the improvement by humans is around 10%. On the other hand, as shown in the Breakdown Acc column, our model can more effectively identify the code improved by humans (rather than machines).

**Relative Improvement Prediction Results** In addition to classifying the efficient code given the code pair, we can further predict how much it is improved. For this task, we measure the performance of our model, by ranking all the code pairs based on their relative improvements and comparing them with predicted improvements. Also, if the prediction score exceeds 1.00, we classify this case that there is an improvement during refinement. As

Table 4: Generalization results by varying the training data.

Training Datasets	Python	C++	All
Python	<b>72.75</b>	58.26	61.52
C++	57.53	60.32	57.90
Python & C++	72.49	<b>62.08</b>	<b>64.42</b>

Table 5: With different Code LLMs, we report their average relative improvement (Avg), as well as the percentage of their degradation and improvement in efficiency. Acc (H) and (M) denote the accuracy on human- and machine-generated codes.

LLM	Avg	Degrade %	Improve %	Acc (H)	Acc (M)
DSC	1.3	31.22	31.51	80.43	65.49
CodeQwen	1.33	28.01	29.12	80.60	58.86
Granite	1.13	21.46	24.55	76.69	48.97
gpt-3.5	1.00	25.34	21.45	78.29	53.72

shown in Table 3, we observe both the high rank correlation coefficient and high accuracy, demonstrating the effectiveness of our approach even in this actual improvement prediction setting.

**Generalization on Different Languages** To see the generalization ability of our approach to different programming languages, we train the model with Python, C++ or both, and measure the performance on Python and C++. As shown in Table 4, we observe that the model trained on one language can be generalizable to the other language, perhaps due to the algorithmic similarities in their codes.

**Results with Different Code LLMs** To see the performance of different Code LLMs in refining the given codes (in terms of efficiency), and to see the performance of our efficiency judgement model trained with the code pairs constructed by different LLMs, we change the code refinement model from DeepSeekCoder (DSC) to recent CodeQwen (Bai et al., 2023), Granite (Mishra et al., 2024), and GPT-3.5. As shown in Table 5, we find that DSC and CodeQwen are superior in improving the efficiency of codes when refining them. Yet, the percentage of improvement made by each model is comparable to the percentage of degradation, which supports again our motivation that we should rethink code refinement. Lastly, our model is able to identify the more efficient code among two different versions of models made across different LLMs.

## 5 Conclusion

In this work, we pointed out that the codes refined by humans or machines are sometimes inferior than originals, and to tackle this, we introduced a novel approach to identify the more efficient code given a pair of codes before and after modifications. We validated our method on multiple code editing scenarios involving both humans and machines, showcasing its substantial efficacy despite its simplicity.

## 320 Limitation

321 In this work, we aim to judge the efficiency of two  
322 different versions of codes, demonstrating that our  
323 approach can distinguish the more efficient code.  
324 However, there are some areas that future work may  
325 improve upon. First, we perform experiments with  
326 two widely used programming languages, such as  
327 Python and C++, and it may be promising to con-  
328 sider other languages, particularly those used less  
329 frequently. In addition, in terms of measuring the  
330 code efficiency, we consider its execution time,  
331 meanwhile, there may be additional factors to con-  
332 sider, such as memory usage, I/O operations, and  
333 the underlying OS environment. Future work may  
334 incorporate these factors to perform a more holistic  
335 assessment of program efficiency. Lastly, beyond  
336 predicting the efficient code, future work may ex-  
337 plore its interpretability, providing the reasons why  
338 certain codes are more efficient than others at a  
339 more fine-grained level (e.g., lines of codes).

## 340 Ethic Statement

341 We believe this work does not have particular con-  
342 cerns about ethics. This is because, it strictly fo-  
343 cuses on the technological aspect of comparing  
344 code efficiency, which does not engage in the uneth-  
345 ical use of LLMs for manipulating software codes,  
346 user data, or any other sensitive information.

## 347 References

348 Marah Abdin, Sam Ade Jacobs, Ammar Ahmad  
349 Awan, Jyoti Aneja, Ahmed Awadallah, Hany Has-  
350 san Awadalla, Nguyen Bach, Amit Bahree, Arash  
351 Bakhtiari, Harkirat Singh Behl, Alon Benhaim,  
352 Misha Bilenko, Johan Bjorck, Sébastien Bubeck,  
353 Martin Cai, Caio C’esar Teodoro Mendes, Weizhu  
354 Chen, Vishrav Chaudhary, Parul Chopra, Allison Del  
355 Giorno, Gustavo de Rosa, Matthew Dixon, Ronen  
356 Eldan, Dan Iter, Abhishek Goswami, Suriya Gu-  
357 nasekar, Emman Haider, Junheng Hao, Russell J.  
358 Hewett, Jamie Huynh, Mojan Javaheripi, Xin Jin,  
359 Piero Kauffmann, Nikos Karampatziakis, Dongwoo  
360 Kim, Mahoud Khademi, Lev Kurilenko, James R.  
361 Lee, Yin Tat Lee, Yuanzhi Li, Chen Liang, Weishung  
362 Liu, Eric Lin, Zeqi Lin, Piyush Madan, Arindam  
363 Mitra, Hardik Modi, Anh Nguyen, Brandon Norick,  
364 Barun Patra, Daniel Perez-Becker, Thomas Portet,  
365 Reid Pryzant, Heyang Qin, Marko Radmilac, Corby  
366 Rosset, Sambudha Roy, Olli Saarikivi, Amin Saied,  
367 Adil Salim, Michael Santacroce, Shital Shah, Ning  
368 Shang, Hiteshi Sharma, Xianmin Song, Olatunji  
369 Ruwase, Xin Wang, Rachel Ward, Guanhua Wang,  
370 Philipp Witte, Michael Wyatt, Can Xu, Jiahang Xu,  
371 Sonali Yadav, Fan Yang, Ziyi Yang, Donghan Yu,

Cheng-Yuan Zhang, Cyril Zhang, Jianwen Zhang, 372  
Li Lyna Zhang, Yi Zhang, Yunan Zhang, and Xiren 373  
Zhou. 2024. *Phi-3 technical report: A highly capa- 374*  
*able language model locally on your phone.* *ArXiv,* 375  
*abs/2404.14219.* 376

OpenAI Josh Achiam, Steven Adler, Sandhini Agarwal, 377  
Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale- 378  
man, Diogo Almeida, Janko Altenschmidt, Sam Alt- 379  
man, Shyamal Anadkat, Red Avila, Igor Babuschkin, 380  
Suchir Balaji, Valerie Balcom, Paul Baltescu, Haim- 381  
ing Bao, Mo Bavarian, Jeff Belgum, Irwan Bello, 382  
Jake Berdine, Gabriel Bernadett-Shapiro, Christo- 383  
pher Berner, Lenny Bogdonoff, Oleg Boiko, Made- 384  
laine Boyd, Anna-Luisa Brakman, Greg Brockman, 385  
Tim Brooks, Miles Brundage, Kevin Button, Trevor 386  
Cai, Rosie Campbell, Andrew Cann, Brittany Carey, 387  
Chelsea Carlson, Rory Carmichael, Brooke Chan, 388  
Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, 389  
Ruby Chen, Jason Chen, Mark Chen, Benjamin 390  
Chess, Chester Cho, Casey Chu, Hyung Won Chung, 391  
Dave Cummings, Jeremiah Currier, Yunxing Dai, 392  
Cory Decareaux, Thomas Degry, Noah Deutsch, 393  
Damien Deville, Arka Dhar, David Dohan, Steve 394  
Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, 395  
Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, 396  
Sim’on Posada Fishman, Juston Forte, Isabella Ful- 397  
ford, Leo Gao, Elie Georges, Christian Gibson, Vik 398  
Goel, Tarun Gogineni, Gabriel Goh, Raphael Gontijo- 399  
Lopes, Jonathan Gordon, Morgan Grafstein, Scott 400  
Gray, Ryan Greene, Joshua Gross, Shixiang Shane 401  
Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, 402  
Yuchen He, Mike Heaton, Johannes Heidecke, Chris 403  
Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, 404  
Brandon Houghton, Kenny Hsu, Shengli Hu, Xin 405  
Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, 406  
Joanne Jang, Angela Jiang, Roger Jiang, Haozhun 407  
Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo 408  
Jun, Tomer Kaftan, Lukasz Kaiser, Ali Kamali, In- 409  
gmar Kanitscheider, Nitish Shirish Keskar, Tabarak 410  
Khan, Logan Kilpatrick, Jong Wook Kim, Christina 411  
Kim, Yongjik Kim, Hendrik Kirchner, Jamie Ryan 412  
Kiros, Matthew Knight, Daniel Kokotajlo, Lukasz 413  
Kondraciuk, Andrew Kondrich, Aris Konstantini- 414  
dis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, 415  
Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, 416  
Jade Leung, Daniel Levy, Chak Ming Li, Rachel 417  
Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, 418  
Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Ade- 419  
ola Makanju, Kim Malfacini, Sam Manning, Todor 420  
Markov, Yaniv Markovski, Bianca Martin, Katie 421  
Mayer, Andrew Mayne, Bob McGrew, Scott Mayer 422  
McKinney, Christine McLeavey, Paul McMillan, 423  
Jake McNeil, David Medina, Aalok Mehta, Jacob 424  
Menick, Luke Metz, Andrey Mishchenko, Pamela 425  
Mishkin, Vinnie Monaco, Evan Morikawa, Daniel P. 426  
Mossing, Tong Mu, Mira Murati, Oleg Murk, David 427  
M’ely, Ashvin Nair, Reiichiro Nakano, Rajeev 428  
Nayak, Arvind Neelakantan, Richard Ngo, Hyeon- 429  
woo Noh, Ouyang Long, Cullen O’Keefe, Jakub W. 430  
Pachocki, Alex Paino, Joe Palermo, Ashley Pantu- 431  
liano, Giambattista Parascandolo, Joel Parish, Emy 432  
Parparita, Alexandre Passos, Mikhail Pavlov, Andrew 433  
Peng, Adam Perelman, Filipe de Avila Belbute Peres, 434

435	Michael Petrov, Henrique Pondé de Oliveira Pinto,	Greenberg, Abhinav Jangda, and Arjun Guha. 2023.	496
436	Michael Pokorny, Michelle Pokrass, Vitchyr H. Pong,	<a href="#">Knowledge transfer from high-resource to low-</a>	497
437	Tolly Powell, Alethea Power, Boris Power, Elizabeth	<a href="#">resource programming languages for code llms.</a>	498
438	Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya	<i>ArXiv</i> , abs/2308.09895.	499
439	Ramesh, Cameron Raymond, Francis Real, Kendra		
440	Rimbach, Carl Ross, Bob Rotsted, Henri Roussez,	Bei Chen, Fengji Zhang, Anh Nguyen, Daoguang Zan,	500
441	Nick Ryder, Mario D. Saltarelli, Ted Sanders, Shibani	Zeqi Lin, Jian-Guang Lou, and Weizhu Chen. 2023a.	501
442	Santurkar, Girish Sastry, Heather Schmidt, David	<a href="#">Codet: Code generation with generated tests.</a>	502
443	Schnurr, John Schulman, Daniel Selsam, Kyla Shep-	In <i>The Eleventh International Conference on Learning</i>	503
444	ppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker,	<i>Representations, ICLR 2023, Kigali, Rwanda, May</i>	504
445	Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie	<i>1-5, 2023.</i> OpenReview.net.	505
446	Simens, Jordan Sitkin, Katarina Slama, Ian Sohl,		
447	Benjamin D. Sokolowsky, Yang Song, Natalie Stau-	Mark Chen, Jerry Tworek, Heewoo Jun, Qiming	506
448	dacher, Felipe Petroski Such, Natalie Summers, Ilya	Yuan, Henrique Ponde, Jared Kaplan, Harrison Ed-	507
449	Sutskever, Jie Tang, Nikolas A. Tezak, Madeleine	wards, Yura Burda, Nicholas Joseph, Greg Brockman,	508
450	Thompson, Phil Tillet, Amin Tootoonchian, Eliz-	Alex Ray, Raul Puri, Gretchen Krueger, Michael	509
451	abeth Tseng, Preston Tuggle, Nick Turley, Jerry	Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin,	510
452	Tworek, Juan Felipe Cer'on Uribe, Andrea Val-	Brooke Chan, Scott Gray, Nick Ryder, Mikhail	511
453	lone, Arun Vijayvergiya, Chelsea Voss, Carroll L.	Pavlov, Alethea Power, Lukasz Kaiser, Moham-	512
454	Wainwright, Justin Jay Wang, Alvin Wang, Ben	mad Bavarian, Clemens Winter, Philippe Tillet, Fe-	513
455	Wang, Jonathan Ward, Jason Wei, CJ Weinmann,	lipo Petroski Such, David W. Cummings, Matthias	514
456	Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian	Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel	515
457	Weng, Matt Wiethoff, Dave Willner, Clemens Win-	Herbert-Voss, William H. Guss, Alex Nichol, Igor	516
458	ter, Samuel Wolrich, Hannah Wong, Lauren Work-	Babuschkin, Suchir Balaji, Shantanu Jain, Andrew	517
459	man, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao,	Carr, Jan Leike, Joshua Achiam, Vedant Misra, Evan	518
460	Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Woj-	Morikawa, Alec Radford, Matthew M. Knight, Miles	519
461	ciech Zaremba, Rowan Zellers, Chong Zhang, Mar-	Brundage, Mira Murati, Katie Mayer, Peter Welinder,	520
462	vin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang	Bob McGrew, Dario Amodei, Sam McCandlish, Ilya	521
463	Zhuang, William Zhuk, and Barret Zoph. 2023. <a href="#">Gpt-</a>	Sutskever, and Wojciech Zaremba. 2021. <a href="#">Evaluat-</a>	522
464	<a href="#">4 technical report.</a>	<a href="#">ing large language models trained on code.</a> <i>ArXiv</i> ,	523
		abs/2107.03374.	524
465	Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang,	Zimin Chen, Sen Fang, and Monperrus Martin. 2023b.	525
466	Xiaodong Deng, Yang Fan, Wenhang Ge, Yu Han,	<a href="#">Supersonic: Learning to generate source code opti-</a>	526
467	Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang	<a href="#">mizations in c/c++.</a> <i>ArXiv</i> , abs/2309.14846.	527
468	Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang		
469	Lu, K. Lu, Jianxin Ma, Rui Men, Xingzhang Ren,	Yihong Dong, Ji Ding, Xue Jiang, Zhuo Li, Ge Li,	528
470	Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong	and Zhi Jin. 2023a. <a href="#">Codescore: Evaluating code</a>	529
471	Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang	<a href="#">generation by learning code execution.</a> <i>ArXiv</i> ,	530
472	Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian	abs/2301.09043.	531
473	Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen		
474	Yu, Yu Bowen, Hongyi Yuan, Zheng Yuan, Jianwei	Yihong Dong, Xue Jiang, Zhi Jin, and Ge Li.	532
475	Zhang, Xing Zhang, Yichang Zhang, Zhenru Zhang,	2023b. <a href="#">Self-collaboration code generation via chat-</a>	533
476	Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and	<a href="#">gpt.</a> <i>ArXiv</i> , abs/2304.07590.	534
477	Tianhang Zhu. 2023. <a href="#">Qwen technical report.</a> <i>ArXiv</i> ,		
478	abs/2309.16609.	Spandan Garg, Roshanak Zilouchian Moghaddam,	535
		Colin B. Clement, Neel Sundaresan, and Chen	536
479	Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie	Wu. 2022. <a href="#">Deeppperf: A deep learning-based ap-</a>	537
480	Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind	<a href="#">proach for improving software performance.</a> <i>ArXiv</i> ,	538
481	Neelakantan, Pranav Shyam, Girish Sastry, Amanda	abs/2206.13619.	539
482	Askell, Sandhini Agarwal, Ariel Herbert-Voss,		
483	Gretchen Krueger, Tom Henighan, Rewon Child,	Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai	540
484	Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,	Dong, Wentao Zhang, Guanting Chen, Xiao Bi,	541
485	Clemens Winter, Christopher Hesse, Mark Chen, Eric	Yu Wu, Y. K. Li, Fuli Luo, Yingfei Xiong, and Wen-	542
486	Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess,	feng Liang. 2024. <a href="#">Deepseek-coder: When the large</a>	543
487	Jack Clark, Christopher Berner, Sam McCandlish,	<a href="#">language model meets programming - the rise of code</a>	544
488	Alec Radford, Ilya Sutskever, and Dario Amodei.	<a href="#">intelligence.</a> <i>ArXiv</i> , abs/2401.14196.	545
489	2020. <a href="#">Language models are few-shot learners.</a> In <i>Ad-</i>		
490	<i>vances in Neural Information Processing Systems 33:</i>	Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas	546
491	<i>Annual Conference on Neural Information Process-</i>	Muennighoff, Denis Kocetkov, Chenghao Mou, Marc	547
492	<i>ing Systems 2020, NeurIPS 2020, December 6-12,</i>	Marone, Christopher Akiki, Jia Li, Jenny Chim,	548
493	<i>2020, virtual.</i>	Qian Liu, Evgenii Zheltonozhskii, Terry Yue Zhuo,	549
		Thomas Wang, Olivier Dehaene, Mishig Davaadorj,	550
494	Federico Cassano, John Gouwar, Francesca Lucchetti,	Joel Lamy-Poirier, João Monteiro, Oleh Shliazhko,	551
495	Claire Schlesinger, Carolyn Jane Anderson, Michael	Nicolas Gontier, Nicholas Meade, Armel Zebaze,	552

553	Ming-Ho Yee, Logesh Kumar Umapathi, Jian Zhu,	Shuai Ma. 2020. <a href="#">Codebleu: a method for automatic</a>	612
554	Benjamin Lipkin, Muhtasham Oblokulov, Zhiruo	<a href="#">evaluation of code synthesis</a> . <i>ArXiv</i> , abs/2009.10297.	613
555	Wang, Rudra Murthy, Jason Stillerman, Siva Sankalp		
556	Patel, Dmitry Abulkhanov, Marco Zocca, Manan	Baptiste Rozière, Jonas Gehring, Fabian Gloeckle,	614
557	Dey, Zhihan Zhang, Nourhan Fahmy, Urvashi Bhat-	Sten Sootla, Itai Gat, Xiaoqing Tan, Yossi Adi,	615
558	tacharyya, W. Yu, Swayam Singh, Sasha Luccioni,	Jingyu Liu, Tal Remez, Jérémy Rabin, Artyom	616
559	Paulo Villegas, Maxim Kunakov, Fedor Zhdanov,	Kozhevnikov, I. Evtimov, Joanna Bitton, Manish P	617
560	Manuel Romero, Tony Lee, Nadav Timor, Jennifer	Bhatt, Cristian Cantón Ferrer, Aaron Grattafiori, Wen-	618
561	Ding, Claire Schlesinger, Hailey Schoelkopf, Jana	han Xiong, Alexandre D’efossez, Jade Copet, Faisal	619
562	Ebert, Tri Dao, Mayank Mishra, Alexander Gu,	Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier,	620
563	Jennifer Robinson, Carolyn Jane Anderson, Bren-	Thomas Scialom, and Gabriel Synnaeve. 2023. <a href="#">Code</a>	621
564	dan Dolan-Gavitt, Danish Contractor, Siva Reddy,	<a href="#">llama: Open foundation models for code</a> . <i>ArXiv</i> ,	622
565	Daniel Fried, Dzmitry Bahdanau, Yacine Jernite, Car-	abs/2308.12950.	623
566	los Muñoz Ferrandis, Sean M. Hughes, Thomas Wolf,		
567	Arjun Guha, Leandro von Werra, and Harm de Vries.	Baptiste Rozière, Marie-Anne Lachaux, Lowik Chanus-	624
568	2023. <a href="#">Starcoder: may the source be with you!</a> <i>ArXiv</i> ,	sot, and Guillaume Lample. 2020. <a href="#">Unsupervised</a>	625
569	abs/2305.06161.	<a href="#">translation of programming languages</a> . In <i>Advances</i>	626
		<i>in Neural Information Processing Systems 33: Annual</i>	627
570	Yujia Li, David Choi, Junyoung Chung, Nate Kush-	<i>Conference on Neural Information Processing</i>	628
571	man, Julian Schrittwieser, Rémi Leblond, Tom, Ec-	<i>Systems 2020, NeurIPS 2020, December 6-12, 2020,</i>	629
572	cles, James Keeling, Felix Gimeno, Agustin Dal	<i>virtual</i> .	630
573	Lago, Thomas Hubert, Peter Choy, Cyprien de,		
574	Masson d’Autume, Igor Babuschkin, Xinyun Chen,	Sijie Shen, Xiang Zhu, Yihong Dong, Qizhi Guo,	631
575	Po-Sen Huang, Johannes Welbl, Sven Gowal,	Yankun Zhen, and Ge Li. 2022. <a href="#">Incorporating do-</a>	632
576	Alexey, Cherepanov, James Molloy, Daniel Jaymin	<a href="#">main knowledge through task augmentation for front-</a>	633
577	Mankowitz, Esme Sutherland Robson, Pushmeet	<a href="#">end javascript code generation</a> . In <i>Proceedings of</i>	634
578	Kohli, Nando de, Freitas, Koray Kavukcuoglu, and	<i>the 30th ACM Joint European Software Engineering</i>	635
579	Oriol Vinyals. 2022. <a href="#">Competition-level code genera-</a>	<i>Conference and Symposium on the Foundations of</i>	636
580	<a href="#">tion with alphacode</a> . <i>Science</i> , 378:1092 – 1097.	<i>Software Engineering, ESEC/FSE 2022, Singapore,</i>	637
		<i>Singapore, November 14-18, 2022</i> , pages 1533–1543.	638
		ACM.	639
581	Aman Madaan, Alex Shypula, Uri Alon, Milad	Atsushi Shirafuji, Md. Mostafizer Rahman, Md.	640
582	Hashemi, Parthasarathy Ranganathan, Yiming Yang,	Faizul Ibne Amin, and Yutaka Watanobe. 2023. <a href="#">Pro-</a>	641
583	Graham Neubig, and Amir Yazdanbakhsh. 2023.	<a href="#">gram repair with minimal edits using codet5</a> . In <i>12th</i>	642
584	<a href="#">Learning performance-improving code edits</a> . <i>ArXiv</i> ,	<i>International Conference on Awareness Science and</i>	643
585	abs/2302.07867.	<i>Technology, iCAST 2023, Taichung, Taiwan, Novem-</i>	644
		<i>ber 9-11, 2023</i> , pages 178–184. IEEE.	645
586	Mayank Mishra, Matt Stallone, Gaoyuan Zhang, Yikang	Hugo Touvron, Louis Martin, Kevin R. Stone, Peter	646
587	Shen, Aditya Prasad, Adriana Meza Soria, Michele	Albert, Amjad Almahairi, Yasmine Babaei, Niko-	647
588	Merler, Parameswaran Selvam, Saptha Surendran,	lay Bashlykov, Soumya Batra, Prajjwal Bhargava,	648
589	Shivdeep Singh, Manish Sethi, Xuan-Hong Dang,	Shruti Bhosale, Daniel M. Bikel, Lukas Blecher, Cris-	649
590	Pengyuan Li, Kun-Lung Wu, Syed Zawad, Andrew	tian Cantón Ferrer, Moya Chen, Guillem Cucurull,	650
591	Coleman, Matthew White, Mark Lewis, Raju Pavu-	David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin	651
592	luri, Yan Koyfman, Boris Lublinsky, Maximilien	Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami,	652
593	de Bayser, Ibrahim Abdelaziz, Kinjal Basu, Mayank	Naman Goyal, Anthony S. Hartshorn, Saghar Hos-	653
594	Agarwal, Yi Zhou, Chris Johnson, Aanchal Goyal,	seini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor	654
595	Hima Patel, Yousaf Shah, Petros Zefos, Heiko Lud-	Kerkez, Madian Khabsa, Isabel M. Kloumann, A. V.	655
596	wig, Asim Munawar, Maxwell Crouse, Pavan Ka-	Korenev, Punit Singh Koura, Marie-Anne Lachaux,	656
597	panipathi, Shweta Salaria, Bob Calio, Sophia Wen,	Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai	657
598	Seetharami R. Seelam, Brian M. Belgodere, Carlos	Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov,	658
599	Fonseca, Amith Singhee, Nirmal Desai, David Cox,	Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew	659
600	Ruchir Puri, and Rameswar Panda. 2024. <a href="#">Granite</a>	Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan	660
601	<a href="#">code models: A family of open foundation models</a>	Saladi, Alan Schelten, Ruan Silva, Eric Michael	661
602	<a href="#">for code intelligence</a> .	Smith, R. Subramanian, Xia Tan, Binh Tang, Ross	662
		Taylor, Adina Williams, Jian Xiang Kuan, Puxin	663
603	Erik Nijkamp, Bo Pang, Hiroaki Hayashi, Lifu Tu, Huan	Xu, Zhengxu Yan, Iliyan Zarov, Yuchen Zhang, An-	664
604	Wang, Yingbo Zhou, Silvio Savarese, and Caiming	gela Fan, Melanie Kambadur, Sharan Narang, Aure-	665
605	Xiong. 2023. <a href="#">Codegen: An open large language</a>	lien Rodriguez, Robert Stojnic, Sergey Edunov, and	666
606	<a href="#">model for code with multi-turn program synthesis</a> . In	Thomas Scialom. 2023. <a href="#">Llama 2: Open foundation</a>	667
607	<i>The Eleventh International Conference on Learning</i>	<a href="#">and fine-tuned chat models</a> . <i>ArXiv</i> , abs/2307.09288.	668
608	<i>Representations, ICLR 2023, Kigali, Rwanda, May</i>		
609	<i>1-5, 2023</i> . OpenReview.net.		
610	Shuo Ren, Daya Guo, Shuai Lu, Long Zhou, Shujie	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	669
611	Liu, Duyu Tang, M. Zhou, Ambrosio Blanco, and	Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le,	670

- 671 and Denny Zhou. 2022. [Chain-of-thought prompting](#)  
672 [elicits reasoning in large language models](#). In *Ad-*  
673 *vances in Neural Information Processing Systems 35:*  
674 *Annual Conference on Neural Information Process-*  
675 *ing Systems 2022, NeurIPS 2022, New Orleans, LA,*  
676 *USA, November 28 - December 9, 2022*.
- 677 Pengcheng Yin and Graham Neubig. 2018. [TRANX:](#)  
678 [A transition-based neural abstract syntax parser for](#)  
679 [semantic parsing and code generation](#). In *Proceed-*  
680 *ings of the 2018 Conference on Empirical Methods*  
681 *in Natural Language Processing, EMNLP 2018: Sys-*  
682 *tem Demonstrations, Brussels, Belgium, October 31 -*  
683 *November 4, 2018*, pages 7–12. Association for Com-  
684 putational Linguistics.
- 685 Zhaojian Yu, Xin Zhang, Ning Shang, Yangyu Huang,  
686 Can Xu, Yishujie Zhao, Wenxiang Hu, and Qiufeng  
687 Yin. 2023. [Wavecoder: Widespread and versatile](#)  
688 [enhanced instruction tuning with refined data genera-](#)  
689 [tion](#). *ArXiv*, abs/2312.14187.
- 690 E. Zelikman, Eliana Lorch, Lester Mackey, and  
691 Adam Tauman Kalai. 2023. [Self-taught optimizer](#)  
692 [\(stop\): Recursively self-improving code generation](#).  
693 *ArXiv*, abs/2310.02304.
- 694 Shuyan Zhou, Uri Alon, Sumit Agarwal, and Graham  
695 Neubig. 2023. [Codebertscore: Evaluating code gen-](#)  
696 [eration with pretrained models of code](#). In *Proceed-*  
697 *ings of the 2023 Conference on Empirical Methods*  
698 *in Natural Language Processing, EMNLP 2023, Sin-*  
699 *gapore, December 6-10, 2023*, pages 13921–13937.  
700 Association for Computational Linguistics.
- 701 Terry Yue Zhuo. 2024. [Ice-score: Instructing large lan-](#)  
702 [guage models to evaluate code](#). In *Findings of the*  
703 *Association for Computational Linguistics: EACL*  
704 *2024, St. Julian’s, Malta, March 17-22, 2024*, pages  
705 2232–2242. Association for Computational Linguis-  
706 tics.



Table 6: Dataset statistics. The first two rows represent the code statistics made by humans and the other rows are the ones made by machines. PL denotes the programming language.

	PL	Training	Validation	Test
Human	Python	14936	761	562
Human	C++	27005	1316	2038
DSC	Python	17079	1046	594
DSC	C++	35666	1290	1949
gpt-3.5	Python	10674	1306	724
CodeQwen	Python	14854	227	207
Granite	Python	20783	650	315

Table 7: Results on the DeepSeek-Coder-Instruct-7B model.

	All			Easy		
	Python	C++	All	Python	C++	All
Zero-shot	52.50	47.91	48.88	51.62	49.05	49.98
Few-shot	54.50	49.86	50.90	55.87	49.81	52.00
Zero-shot CoT	54.50	49.86	50.90	55.53	52.73	53.75
Few-shot CoT	53.11	49.99	50.69	47.04	52.16	50.3
<b>Ours</b>	<b>73.44</b>	<b>62.90</b>	<b>65.27</b>	<b>78.88</b>	<b>72.68</b>	<b>74.93</b>

## A Additional Experimental Setups

**Dataset Details** We report the dataset statistics in Table 6. Note that, in order to obtain the stable code execution result to decide which code is more efficient than the other, we run every refined code with its original code three times and then select the one whose results are consistent across those three runs. In addition, the code execution is performed following the existing setup (Madaan et al., 2023).

**Fine-tuning Details** We provide details on fine-tuning the efficiency judgment model: we fine-tune the Code LLM (namely, DeepSeek-Coder-Instruct-1.3B) over 10 epochs with a batch size of 16 and a learning rate of  $2e-5$ , and we select the best epoch based on performance on the validation set.

**Prompts** In Table 9, we provide the prompts used to elicit the Code LLM to refine the code and to predict the code efficiency (in classification and regression settings). For the efficiency prediction problem, we randomly shuffle the sequence of the original and its refined codes.

## B Additional Experimental Results

Here, we provide additional experimental results.

**Results with Larger Models** We conduct an auxiliary analysis to see how the performance of different methods changes if a model larger than DeepSeek-Coder-Instruct-1.3B (that we use for main experiments) is used. Specifically, we use its 7B model as the base code LLM and then classify the efficient code given code pairs. As shown

Table 8: Average accuracy results for code improvement classification with order perturbation across multiple different runs, where we report the variance in parentheses.

	All			Easy		
	Python	C++	All	Python	C++	All
Zero-shot	50.26 (0.7)	47.02 (5.9)	47.75 (2.9)	48.49 (4.2)	49.94 (3.2)	49.42 (0.2)
Few-shot	50.43 (1.2)	50.00 (2.7)	50.10 (2.3)	49.61 (0.4)	48.32 (0.0)	48.79 (0.1)
Zero-shot CoT	49.70 (0.9)	50.53 (2.7)	50.34 (2.2)	48.94 (0.0)	50.29 (6.6)	49.80 (2.7)
Few-shot CoT	49.27 (3.2)	49.79 (2.3)	49.67 (2.5)	50.62 (2.5)	48.03 (0.0)	48.97 (0.3)
<b>Ours</b>	<b>71.67 (1.3)</b>	<b>62.09 (0.0)</b>	<b>64.25 (0.1)</b>	<b>77.60 (0.0)</b>	<b>70.52 (0.3)</b>	<b>73.09 (0.1)</b>

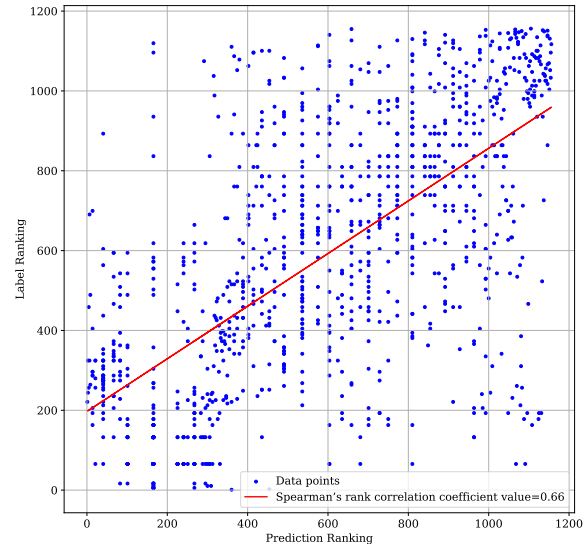


Figure 3: Visualization of the Spearman’s rank correlation between the ranks of the actual relative improvements and the predicted relative improvements of code pairs, for our model.

in Table 7, we observe results similar to those obtained from the smaller 1.3B model, where our model is consistently superior to other baselines.

**Analysis on Bias for Code Sequence** In our code efficiency judgment task, we put a sequence of two codes in the input of Code LLMs, and the Code LLMs may have a bias in this sequence (e.g., predicting the code at the last more often). To see whether they have such a bias, we conduct an additional experiment, flipping the order of the code pairs in the input. We report the results in Table 8, and, from this, we observe that there are no such the notable bias in the sequence of codes.

**Visualization of Rank Correlation** To visualize how accurate the predicted results of relative improvements of code pairs from our model are, we compare their ranks with the ground-truth ranks calculated by actual relative improvements of code pairs. As shown in Figure 3 where we present a scatter plot of rank correlations along with their coefficient value, we observe that the results from our approach have a positive correlation with the ground truth, demonstrating its effectiveness in predicting the relative improvement of code pairs.

761 **Qualitative Analysis** We provide some example  
762 codes in Python and C++ in Figures 4 and 5. From  
763 these two examples, we observe that, despite the  
764 difference in grammar across different program-  
765 ming languages, code pairs from them can share  
766 the same underlying algorithms. This result sup-  
767 ports our finding on generalization ability that our  
768 model trained on one programming language can  
769 be generalizable to other languages (See Table 4).

Table 9: A list of prompts that we used for code refinement and efficiency predictions. It is worth noting that the variable inside the parentheses `{}` is replaced with its actual code.

Types	Prompts
<b>Code Refinement</b>	Update the given code to make it more efficient. {Original code}
<b>Efficiency Classification</b>	Given a selection of code, determine which one is the most efficient in computing. A: {Original code or Refined code} B: {Refined code or Original code}
<b>Efficiency Regression</b>	Given two sets of code, assess how much Code B has improved compared to Code A. A. {Original code} B. {Refined code}

```
1 # Python Example
2 N = int(eval(input()))
3 print(((N*(N-1))/2))
```

```
1 // C++ Example
2 #include<iostream>
3 using namespace std;
4 int main() {
5     long long int n;
6     cin>>n;
7     cout<<n*(n-1)/2<< endl;
8 }
```

Figure 4: Generated Python and C++ samples for the question "For an integer  $N$ , we will choose a permutation  $\{P_1, P_2, \dots, P_N\}$  of  $\{1, 2, \dots, N\}$ . Then, for each  $i = 1, 2, \dots, N$ , let  $M_i$  be the remainder when  $i$  is divided by  $P_i$ . Find the maximum possible value of  $M_1 + M_2 + \dots + M_N$ . Constraints  $N$  is an integer satisfying  $1 \leq N \leq 10^9$ ".

```
1 # Python Example
2 from math import floor, ceil
3
4 X = int(eval(input()))
5 cash = 100
6 count = 0
7 while cash < X:
8     cash=floor(cash*1.01)
9     count += 1
10
11 print(count)
```

```
1 // C++ Example
2 #include <bits/stdc++.h>
3 using namespace std;
4 int main() {
5     long long X;
6     cin >> X;
7
8     int year=0;
9     long long s=100;
10
11     while(s<X){
12         s=s*1.01;
13         year++;
14     }
15     cout << year << endl;
16 }
```

Figure 5: Generated Python and C++ samples for the question "Takahashi has a deposit of 100 yen (the currency of Japan) in AtCoder Bank. The bank pays an annual interest rate of 1% compounded annually. (A fraction of less than one yen is discarded.) Assuming that nothing other than the interest affects Takahashi's balance, in how many years does the balance reach X yen or above for the first time?".