# BEYOND CORRECTNESS: BENCHMARKING MULTI DIMENSIONAL CODE GENERATION FOR LARGE LAN GUAGE MODELS

Anonymous authors

Paper under double-blind review

#### ABSTRACT

In recent years, researchers have proposed numerous benchmarks to evaluate the impressive coding capabilities of large language models (LLMs). However, current benchmarks primarily assess the accuracy of LLM-generated code, while neglecting other critical dimensions that also significantly impact code quality in real-world development. Moreover, relying exclusively on correctness as the guiding metric renders LLMs susceptible to data contamination. Therefore, this paper proposes the **RACE** benchmark, which comprehensively evaluates the quality of code generated by LLMs across 4 dimensions: Readability, mAintainability, Correctness, and Efficiency. Specifically, considering the *demand-dependent* nature of dimensions beyond correctness, we design various types of user requirements for each dimension to assess the model's ability to generate correct code that also meets user demands. We analyze 28 representative LLMs based on RACE and find that: 1) current correctness-centric benchmarks fail to capture the multifaceted requirements of code in real-world scenarios, while RACE provides a comprehensive evaluation that reveals the defects of LLMs across multiple dimensions; 2) the RACE benchmark serves as an effective tool for resisting the risk of data contamination; 3) even the most advanced code LLMs still encounter significant challenges in customized requirements involving complex instructions; 4) most LLMs exhibit an inherent preference for specific coding style. These findings highlight the need for a multidimensional evaluation of code LLMs, emphasizing metrics beyond correctness for real-world applications. Future efforts should aim to develop novel learning algorithms to enhance code generation under varied constraints and improve coverage and usability for diverse user needs<sup>1</sup>.

034

006

008 009 010

011

013

014

015

016

017

018

019

021

023

025

026

027

028

029

031

032

#### 1 INTRODUCTION

037

The impressive coding capabilities demonstrated by Large Language Models (LLMs) are reshaping
the landscape of software development (Zheng et al., 2023c;b; Fan et al., 2023), attracting significant
attention from researchers. To accurately measure and compare the coding capabilities of various
LLMs, numerous benchmarks have been proposed to evaluate the code generation (Chen et al., 2021;
Austin et al., 2021; Hendrycks et al., 2021), completion (Gong et al., 2024), and execution (Jain et al., 2024a) abilities of LLMs.

However, current benchmarks primarily focus on evaluating the correctness of LLM-generated code,
while neglecting other critical dimensions that also significantly impact code quality in real-world
development scenarios. For example, Börstler et al. (2023) investigate various aspects of code quality and find that code readability is the most decisive property for high-quality code (Dantas et al.,
2023; Oliveira et al., 2020). Additionally, code maintainability is crucial for ensuring the software
remains adaptable and easy to update, ultimately reducing long-term costs and technical debt (Hegedus, 2013). Code efficiency is essential for optimizing performance, reducing resource consumption,
and ensuring scalability in software applications (Curtis et al., 2022; Börstler et al., 2023). As shown
in Figure 1, current benchmarks lack evaluation on these critical dimensions that impact code qual-

<sup>&</sup>lt;sup>1</sup>We provide our benchmark and source code at this anonymous repository https://anonymous. 4open.science/r/RACE-4AF1.

069



Figure 1: Current benchmarks perform single-dimension evaluations and mostly focus only on code correctness (upper right); our proposed RACE benchmark performs multi-dimensional code evaluations to identify truly high-quality code beyond correctness (lower right).

071 ity, making it challenging to distinguish genuinely high-quality code from merely correct code. Such deficiency in evaluation could lead to incomplete assessments of the coding capabilities of different LLMs in real-world development scenarios. Furthermore, if these correctness-based benchmarks 073 serve as guiding indicators and correctness alone becomes the sole criterion for driving LLM devel-074 opment, these models might end up memorizing the exact solutions from the training data instead of 075 understanding the underlying principles or patterns. This overfitting implies the model may repro-076 duce code that is highly similar to the training data during inference, leading to data leakage. Con-077 sequently, this singular focus on correctness can render LLMs susceptible to data contamination, which has been proven to be quite prevalent due to the exponential scaling of pre-training data (Rid-079 dell et al., 2024; Cao et al., 2024). Therefore, there is an urgent need for a multidimensional code evaluation benchmark that transcends correctness, addressing the gap between LLM-generated code 081 and real-world scenarios, and steering code LLMs towards comprehensive development.

082 To this end, we propose the **RACE** benchmark, designed to comprehensively evaluate the code gen-083 erated by LLMs across multiple dimensions including Readability, mAintainability, Correctness, 084 and Efficiency. However, it is not trivial to develop a reliable multi-dimensional benchmark for 085 code generation. The first challenge is to design a quantifiable evaluation framework with corresponding metrics for each dimension. Unlike correctness, other dimensions are typically difficult to 087 quantify with a single metric (e.g., accuracy). To address this, we refer to the definition of readability, maintainability, and efficiency in various quality models (Curtis et al., 2022; Nistala et al., 2019; Sadeghzadeh Hemayati & Rashidi, 2017), and summarize multiple representative factors for each dimension of code quality. Furthermore, dimensions beyond correctness cannot directly use the pass 090 rate of test cases as the performance metric. Therefore, we develop specific evaluation metrics for 091 each factor within these dimensions, which can be objectively and automatically calculated based 092 on static analysis and runtime monitoring methods. As illustrated in Figure 1 and 2, by integrating performance across multiple factors, we can comprehensively assess the quality of LLM-generated 094 code in each dimension. Another more critical challenge is that dimensions other than correctness are *demand-dependent*. This means a fixed and uniform standard cannot be used to assess what 096 constitutes better code. Instead, different application scenarios could have varying requirements for code generation. For instance, various projects require unique coding styles and interface standards 098 for adaptability and scalability. Additionally, balancing time and space efficiency based on hardware conditions ensures code operates efficiently. Therefore, a genuinely practical model should generate 099 correct, customizable code that meets multiple dimensional requirements. To achieve this, we design 100 various demands for each factor and incorporate them into the task descriptions, requiring the model 101 to generate code that is both correct and meets the specified requirements. For example, we design 102 multiple instructions that direct the model to generate multiple versions of code, each optimized 103 differently for time and space efficiency. By incorporating the aforementioned quantifiable evalua-104 tion framework, we can accurately and efficiently quantify the extent to which LLM-generated code 105 fulfills the corresponding customized requirements across each dimension. 106

107 Based on the RACE benchmark, we conduct a comprehensive evaluation and systematic analysis of 28 LLMs across various scales, encompassing the most advanced open-source (e,g, Qwen2.5108 72B (Yang et al., 2024) and DeepSeek-V2.5 (Zhu et al., 2024)) and closed-source models (e.g., 109 GPT-40 (OpenAI, 2024a), o1-mini (OpenAI, 2024b) and Claude-3.5-Sonnet) in terms of coding ca-110 pabilities. Our findings reveal that using correctness as the only guiding indicator is insufficient 111 for code benchmarks to steer code LLMs towards comprehensive advancement. 1) Current 112 benchmarks fail to capture the multifaceted requirements of code in real-world scenarios. Therefore, current code LLMs, developed with a primary focus on correctness, exhibit significant room 113 for improvement in other critical dimensions, including code readability, maintainability, and effi-114 ciency (§4.2). 2) We present concrete evidence that current code benchmarks are susceptible to data 115 contamination, compromising the fairness and reliability of the evaluation conclusions. In contrast, 116 a contaminated model may simply reproduce memorized code without the ability to generate diverse 117 solutions that meet various user requirements. Therefore, the RACE benchmark can robustly pro-118 vide stable assessment results even under data contamination settings (§4.3). Moreover, a deeper 119 analysis based on RACE reveals notable deficiencies in current code LLMs: 3) Even the most 120 advanced code LLMs severely struggle to understand and follow complex instructions that include 121 several customization requirements, with performance deteriorating significantly as the number of 122 requirements increases (§4.4). 4) Most LLMs exhibit an inherent preference for specific coding 123 styles, making it difficult for them to follow user instructions that are inconsistent with their preference (§4.5). The findings above highlight the importance of a multidimensional evaluation of code 124 LLMs, while also revealing the necessity for metrics that extend beyond correctness to guide the 125 development in real-world scenarios. In the future, code LLMs will require the design of novel 126 learning algorithms to acquire the ability to generate high-quality code subject to multidimensional 127 constraints, as well as to enhance their coverage and usability concerning diverse user requirements. 128

The main contributions of this paper can be summarized as follows:

- We propose a novel multi-dimensional evaluation framework for code generation.
- Based on the framework, we construct a comprehensive benchmark named RACE, featuring data construction, customized requirement instructions, and specific evaluation metrics.
- We evaluate and analyze 28 LLMs on the RACE benchmark, and obtain valuable conclusions that reveal the limitations of current benchmarks and models.

## 2 RELATED WORK

2.1 CODE LLMS

141 The outstanding code generation capabilities exhibited by LLMs have attracted considerable atten-142 tion from researchers (Wang et al., 2021; Li et al., 2022; Fried et al., 2022; Xu et al., 2022; Roziere 143 et al., 2023; Zheng et al., 2023a). Some representative code LLMs, such as CodeX (Chen et al., 144 2021), CodeGen (Nijkamp et al., 2022), and AlphaCode (Li et al., 2022), have achieved notable 145 performance in code generation, program repair, and code translation. Currently, research on LLMs for code primarily focuses on data and pretraining methods. For training data collection, Wizard-146 Coder (Luo et al., 2024) introduces code instruction-following training constructed by Evol-Instruct 147 to enhance the capabilities of code LLMs. For pretraining methods, StarCoder (Li et al., 2023a) 148 and DeepSeek-Coder (Guo et al., 2024) incorporate fill-in-the-middle training task to enhance the 149 model's capability to handle various structural arrangements in code. With the rapid advancement 150 of code LLM capabilities, there is an increasing demand for reliable and comprehensive code eval-151 uation benchmarks.

152 153 154

129

130

131 132

133 134

135

136 137

138 139

140

## 2.2 CODING BENCHMARK FOR LLMS

The existing benchmarks for LLM-based code (Ni et al., 2023), such as HumanEval (Chen et al., 2021), APPS (Hendrycks et al., 2021), MBPP (Austin et al., 2021), CodeContests (Li et al., 2022), and DS-1000 (Lai et al., 2023), focusing on the correctness of generated code in scenarios such as code exercises, data science, and competitions (Yan et al., 2023; Li et al., 2023b; Shinn et al., 2024). However, these efforts only focus on the correctness of the generated code, using the pass rate of test cases as the sole evaluation metric. Meanwhile, there has been a recent trend in considering other dimensions (Li et al., 2024; Jain et al., 2024b; Tian et al., 2024); for example, Huang et al. (2024) evaluate the efficiency of the generated code, while Dillmann et al. (2024) bridge the con-



Figure 2: The overall evaluation pipeline in RACE benchmark.

nection between cross-entropy and logical lines of code. Nevertheless, these studies neither account for the demand-dependent nature of these dimensions nor systematically evaluate the LLM's code capabilities across multiple dimensions.

#### 3 **RACE BENCHMARK CONSTRUCTION**

The philosophy of our framework design comes from the demands for code quality in software 182 engineering (Börstler et al., 2023). Firstly, we summarize multiple representative factors for each 183 dimension based on their respective quality definitions (Curtis et al., 2022; Nistala et al., 2019; Sadeghzadeh Hemayati & Rashidi, 2017). Secondly, we design several reasonable customized re-185 quirements for each factor and integrate them into task descriptions, requiring the model to generate 186 code that is both correct and meets these requirements. Information on the detailed evaluation data is presented in Table 1. Finally, leveraging static analysis and runtime monitoring techniques, we 188 develop evaluation metrics tailored to each factor to achieve accurate and efficient evaluation. The 189 specific designs of each instruction refer to Appendix B.2.

190 191

192

187

174 175

176

177

178 179

180 181

#### 3.1 CORRECTNESS

193 Correctness is the core and foundation for evaluating whether the functionality of code generated by mod-194 els meets expectations. Therefore, to comprehensively 195 assess the capability of LLMs in generating function-196 ally correct code across various task scenarios, we se-197 lect 4 datasets with different distributions: HumanEval+ and MBPP+ (Liu et al., 2024) for code exercise prob-199 lems, ClassEval (Du et al., 2023) for class-level code 200 generation, and LeetCode (Guo et al., 2024) for coding 201 competition problems. To mitigate bias from extrane-202 ous information in the original dataset affecting the cus-203 tomized requirements, we exclude such information from 204 the datasets. We use the macro accuracy across 4 datasets as the metric for correctness. 205

Table 1: The sources and number of evaluation cases for each factor in the **RACE** benchmark

Factors	Data Source	# Cases
	Correctness	
Correctness	HumanEval+, MBPP+, ClassEval, LeetCode	923
	Readability	
Code Length Name Convention Comments	HumanEval+	492 984 328
М	aintainability	
Maintainability Index Modularity	ClassEval LeetCode	100 540
	Efficiency	
Time Complexity Space Complexity	LeetCode	101

206 Furthermore, to investigate the impact of adding cus-207 tomized demands on code correctness, we also calculate 208

the accuracy of the generated code when instructions with customized requirements are provided.

209 210

#### 3.2 READABILITY 211

212 In real-world development scenarios, maintaining a consistent coding style is essential for enhancing 213 comprehensibility and minimizing the time required for code maintenance, often referred to as code readability (Börstler et al., 2023). One of the most fundamental aspects of coding style is line length; 214 excessively long lines can lead to truncation on screens. Additionally, adopting clear and consistent 215 naming conventions enables developers to quickly grasp the functionality of interfaces, while wellplaced comments facilitate a rapid understanding of the implementation logic. Consequently, we condense code readability into three key factors: Length, Naming Convention, and Comment. In response to real-world development needs, we collect a set of customizable options for each factor.

219 For the Length factor, readability requirements for code length can vary depending on display scales 220 in different user scenarios. To address this, we refer to PEP8 style guidelines for Python and define 221 the following user requirements regarding code length: (60, 20), (70, 30), and (79, 40), where the 222 parentheses represent the maximum line length and the maximum number of lines in functions, re-223 spectively. For the Naming Convention factor, camel-case and snake-case are widely used naming 224 conventions in programming, with specific preferences varying by project. Therefore, we provide 225 the option to choose between camel-case and snake-case based on the conventions employed for 226 functions and variables. For the **Comment** factor, different levels of granularity serve distinct purposes. Line-level comments clarify implementation details and are beneficial for novice program-227 mers, while function-level comments enhance understanding of functionality and usage. Thus, we 228 offer customization options for both comment types. Please kindly note that although a few read-229 ability requirements can be addressed using formatting tools, we believe that assessing whether a 230 model can directly generate code that meets readability standards provides valuable insights into the 231 model's ability to follow user instructions. Moreover, we further find that the model's capability in 232 readability serves as a significant indicator of its overall coding proficiency. Refer to Appendix A 233 for corresponding experiment results due to page limitations. 234

To align with real-world scenarios that require readability, we conduct experiments on HumanEval+ (Liu et al., 2024) dataset, which consists of coding exercise tasks. We incorporate the aforementioned customized requirements into the problem descriptions to evaluate the model's coding capabilities in terms of readability. To measure code readability, we analyze the various components of the generated code using abstract syntax trees. We then develop corresponding rule-based and regular expression-based methods to measure code length, detect naming conventions, and differentiate between different levels of comment granularity.

241

260 261 262

# 242 3.3 MAINTAINABILITY 243

The maintainability of code plays a vital role in the long-term health of software and the efficiency
of development teams. Numerous quality models propose empirical quantitative measures to assess
maintainability. Additionally, the single responsibility principle is essential in code design, helping
to prevent excessive functional coupling. Based on these principles, we identify two key factors
influencing code maintainability: Maintainability Metric and Modularity.

249 For the Maintainability Metric factor, we use the Maintainability Index (MI) (Coleman et al., 1994) 250 to measure how maintainable the code is, which is widely used in the Microsoft Visual Studio 2010 251 development environment. To address concerns regarding the unreliability of evaluations related to the MI metric when faced with differences in the volume and organizational structure of the code 252 under assessment (Heitlager et al., 2007), we conduct assessments solely on fixed evaluation data. 253 We calculate the maintainability index values of the generated code, thereby enabling a horizontal 254 comparison of the complexity of code generated by different models when faced with the same task. 255 This comparison reflects the variations in their ability to produce maintainable code and substantially 256 alleviates the aforementioned concerns. 257

258 Specifically, MI is a four-metric polynomial equation, resulting in a value between 0 and 100, with 259 higher values indicating greater maintainability. The formulation is as follows:

$$MI = \max\left[0,100 \cdot \frac{171 - 5.2\ln V - 0.23G - 16.2\ln L + 50\sin(\sqrt{2.4C})}{171}\right]$$
(1)

where V is Halstead Volume to identify measurable properties of the code, G is Cyclomatic Complexity corresponding to the number of decisions a block of code contains plus 1, L is the number of source lines of code, and C is the percent of comment lines. To comprehensively assess the maintainability requirements satisfaction of LLM-generated code, we conduct experiments on ClassEval (Du et al., 2023) dataset, to ensure the complexity of the code problems.

For the **Modularity** factor, different requirements determine the varying levels of code modularization. Achieving compactness often requires implementing functionality within a single function, while maximizing code reusability typically involves the use of multiple functions. Accordingly, we define several customization options: implementing functionality using 1, 2, or 3 functions. To
assess the modularity of the generated code, we conduct experiments on LeetCode (Guo et al., 2024)
dataset, which presents a greater challenge and thus enhances discriminative capability. Additionally, we employ the abstract syntax tree to extract all function nodes from the generated code to
verify whether the number of function nodes aligns with the defined level of modularity.

#### 276 3.4 EFFICIENCY

275

277

278

279

280

281

282

283

284

291 292 293

295 296

297 298

299

300

301

302

303 304

305

309

In most applications, code efficiency is closely tied to user experience and business process effectiveness. Typically, efficiency is assessed through time complexity and space complexity. Given the varying hardware conditions of users, it is common practice to balance execution time and memory usage or to optimize one of these aspects to the extreme to ensure code efficiency. To address these scenarios, we gather 101 cases from LeetCode programming problems designed to simulate such conditions. These cases are customized with specific requirements for time complexity, space complexity, or both, to evaluate how well the LLM-generated code meets the efficiency standards.

To measure code efficiency, we propose the Normalized Index (NI), i.e., to measure the degree to which the generated code satisfies the complexity requirement. Given two standard solution with time and space complexity  $C_1^T, C_1^S$  and  $C_2^T, C_2^S$ , respectively, where  $C_1^T$  and  $C_2^S$  are better, and given their total running time  $T_1, T_2$  ( $T_1 < T_2$ ) and memory usage  $S_1, S_2$  ( $S_1 > S_2$ ) on all test cases. Now there is a code  $\hat{C}$  to be evaluated, which has a running time  $\hat{T}$  and memory usage  $\hat{S}$ , with requirements  $C_1^T, C_1^S$ , then the normalized index is:

$$NI_T = 100 \cdot Clip\left(1 - \frac{\hat{T} - T_1}{T_2 - T_1}, 0, 1\right), \ NI_S = 100 \cdot Clip\left(1 - \frac{\hat{S} - S_1}{S_1 - S_2}, 0, 1\right)$$
(2)

 $NI_T$  indicates the degree of time complexity toward  $C_1^T$ , and  $NI_S$  indicates the degree of space complexity toward  $C_2^S$ .

## 4 EXPERIMENTS

In this section, we conduct a detailed evaluation of 28 LLMs and obtain several valuable findings. We first introduce the input formats and inference configurations for code generation tasks, along with the selection of LLMs. Subsequently, we present the overall experimental findings and conduct further analysis of the results to derive meaningful conclusions. The detailed experimental results are shown in Appendix C.

4.1 Settings

Task formats We construct the different prompts based on the completion style and chat style, to
 better induce the LLMs to accomplish the corresponding tasks, see details in Appendix B.2. In the
 inference process, we use a greedy strategy and set the temperature to 0.

- 320 321
- 321 4.2 OVERALL RESULTS322
- The overall evaluation results on all 4 dimensions of each LLM are demonstrated in Table 2, and Figure 3 provides a more intuitive comparison of the capabilities across various dimensions for rep-

Models We select 28 widely-used closed-source and open-source LLMs ranging in different 310 sizes, including state-of-the-art code LLMs. For closed-source models, our experiments include 311 the GPT series (GPT-3.5-turbo-0125, GPT-4o-2024-05-13, and GPT-4o-mini), o1-mini-2024-09-312 12 (OpenAI, 2024b), and Claude-3.5-Sonnet. For open-source models, our experiments include 313 several series: DeepSeek (Guo et al., 2024; Zhu et al., 2024) (DeepSeek-Coder-Ins-7B/13B/34B and 314 DeepSeek-V2.5), CodeLlama (Roziere et al., 2023) (CodeLlama-Ins-7B/13B/34B and CodeLlama-315 Python-7B/13B/34B), WizardCoder (Luo et al., 2024) (WizardCoder-15B/33B and WizardCoder-316 Python-7B/13B), Qwen (Bai et al., 2023; Yang et al., 2024; Hui et al., 2024) (CodeQwen1.5-317 7B-Chat, Qwen2.5-Coder-7B-Ins, Qwen2-72B-Ins, and Qwen2.5-72B-Ins), Llama3 (Dubey et al., 318 2024) (Llama3-Ins-8B/70B), Mixtral-8x22B (Jiang et al., 2024) and StarCoder2-15B (Lozhkov 319 et al., 2024). The full list of models is shown in Appendix B.1.

324 Table 2: Based on the RACE benchmark, the performance results for each LLM in code correctness 325 (C), readability (R), maintainability (M), and efficiency (E). RN, RL, RC, and EC denote the Name 326 Convention, Length, Comments, and Complexity factor. MI denotes the Maintainability Index. MC denotes the Modularity factor.  $NI_T$  and  $NI_S$  are metrics for code efficiency. RACE Score represents 327 the overall metrics at the dimension level. The symbol (\*) indicates that the metric is a scalar from 0 328 to 100, and the rest are percentages (%). The symbol  $(\dagger)$  indicates that the results are obtained from 329 a randomly sampled 30% of the evaluation data, in order to optimize budget efficiency. 330

		RACE	С		R		N	1	I	3		
Models		Overall	C	RN	RL	RC	MI*	MC	$NI_T^*$	NI <sub>S</sub> *		
			Ir	nstruct-T	vne							
	2024 00 12	(25	70.1	00.7	47.5		64.4	(( 1†	(0.2 <sup>†</sup>	40.0 <sup>†</sup>		
Ol-mini- Claude	2024-09-12 3.5. Sonnet	63.5 62.3	70.1 64.6	80.7 74 4	47.5	65.5	64.4 75.3	50.1	60.31 56.8	40.0 <sup>+</sup> 49.7		
GPT-40	.5-30met	<u>57.2</u>	<u>59.9</u>	78.6	63.2	70.4	75.1	$\frac{39.8}{35.2}$	$\frac{30.8}{44.0}$	42.0		
GPT-4o-	mini	52.5	56.4	67.6	55.7	72.9	73.5	23.3	40.3	39.5		
GPT-3.5	-turbo-0125	43.6	44.7	51.4	46.1	47.5	80.2	18.5	27.5	36.5		
CL-7B-I	ns	23.2	23.9	17.8	23.4	22.2	71.8	7.2	8.2	8.8		
CL-13B	-Ins	26.9	24.4	22.9	23.6	29.0	82.1	7.6	10.4	16.1		
CL-34B	-Ins	24.4	26.0	21.9	17.5	10.7	73.2	8.5	14.4	13.8		
DS-Cod	er-6.7B-Ins	39.8	39.2	45.8	46.6	50.0	79.3	8.2	27.1	30.0		
DS-Cod	er-/B-Ins	38.9	39.9	30.8	46.0	53.7	79.6 75.7	8.9	25.1	26.8		
DS-Cod	21-33D-1118 pr_V2_16B_Ins	44.0	44.7 50.9	39.0 41.8	55.5 57.7	54.0 47.5	78.2	19.8	33.3 40.2	30.1 47.7		
DS-V2 4	-236B	57.1	59.0	72.2	66.1	65.8	72.9	33.9	46.4	49.5		
CodeOw	en1 5-7B-Chat	45.2	46.3	48.8	$\frac{00.1}{47.0}$	62.2	82.3	13.0	30.7	$\frac{47.5}{37.7}$		
Owen2.5	5-Coder-7B-Ins	49.0	57.1	53.0	51.8	61.3	78.6	17.6	37.0	33.7		
Qwen2-'	72B-Ins	50.1	53.1	73.6	47.6	60.1	79.4	22.8	32.3	39.4		
Qwen2.5	5-72B-Ins	61.3	64.1	77.2	72.1	72.8	76.7	40.4	47.9	49.4		
Mixtral-	8x22B	42.2	42.0	56.2	47.8	56.1	79.6	9.1	24.7	33.2		
Llama3-	8B-Ins	35.2	35.6	44.3	3 23.6	40.0	79.8	8.1	23.5	26.9		
lama3-	70B-Ins	47.2	44.4	66.0	47.8	54.2	79.8	25.2	29.2	42.8		
			Cor	npletion	Туре							
CL-7B-I	Py	24.0	20.4	20.9	25.8	12.5	79.4	3.7	14.3	14.4		
CL-13B	Py	25.6	21.7	23.1	30.9	24.4	78.6	2.4	13.8	14.7		
CL-34B	-Py	23.6	19.2	18.8	26.7	8.6	85.3	2.2	12.0	14.4		
WC-Py-	7B	26.2	25.2	22.8	28.0	10.1	79.3	7.2	15.3	16.7		
WC-Py-	13B	29.3	26.3	23.9	33.1	30.5	78.8	8.5	16.2	19.8		
WC 33E	5	30.4 40.8	28.0	24.0 40.0	27.8 47.6	28.1	80.0	/.8	21.8	24.2		
StarCod	er2-15B	29.2	28.5	25.8	27.9	44.0 22.0	74.2	9.5 6.1	20.6	25.1		
Claude-3	3.5-Sonnet C	GPT-4o-2024-05-13	De	epSeek-V2.5-	236B	Qwen2.5	-72B-Instruct	Co	deLlama-34B-	Instruct		
RC	RL RN RC.	RL 612 RN	RC 🗸	RL	RN	RC	RL 721 RM	RC.	RL	RN		
		24 72			72	(The second seco	- AN					
11 22	64.6 C MI 731	((()))399	С МІ 720	$((\otimes))$	59.0 C MI	nd ((((	<b>64.1</b>	) с мі (732	11-521-8 10-5 12-521-8 25-12-9	c		
	394	100		340 AGA AS		hot /	82.0 Bg/		X			
MC	NIS MC	NIs	мс	Ŧ	NIs	мс	NI	Is MC NIs				

Figure 3: Performance radar charts of several representative LLMs on the RACE benchmark, visually illustrating the capability distribution across 8 detailed factors.

361

362

363 364

365

366

367

368

369

370

371

372

373

374

375

376

377

#### resentative models. We find that compared to previous correctness-centric benchmarks, RACE can provide a multi-dimensional comprehensive evaluation for code LLMs, offering valuable insights for their application in real-world scenarios.

• From an overview perspective, current code LLMs still have considerable room for improvement in generating correct and user-compliant code across multiple dimensions. For instance, even the most advanced model, o1-mini, achieves only a score of 60.3 in time complexity, with most models below 50. Additionally, apart from o1-mini and Claude-3.5-Sonnet, all other models exhibit performance below 45% in modularity. Further, we find that incorporating different user demands into instructions has varying impacts on code accuracy. For example, increasing the requirement to add comments in appropriate sections can enhance accuracy, which we hypothesize is due to comments facilitating an implicit chain-of-thought. In comparison, adding requirements related to code length tends to decrease accuracy, which may be attributed to the model's inherent preferences towards code of varying lengths (see detailed results in Appendix C.1). These findings offer valuable insights for designing better prompting methods and future optimization directions.

			-													
		2 Epochs			4 Epochs			6 Epochs			8 Epochs		10 Epochs			
Benchmark	Clean	w/ Test	$\Delta\downarrow$	Clean	w/ Test	$\Delta \downarrow$										
HumanEval+	22.0	40.9	+18.9	23.8	62.2	+38.4	20.1	89.6	+69.5	22.0	95.7	+73.7	22.6	97.6	+75.0	
MBPP+	35.7	52.4	+16.7	36.8	77.2	+40.4	33.1	91.5	+58.4	32.3	98.1	+65.8	33.1	98.4	+65.3	
ClassEval	14.0	22.0	+8.0	11.0	46.0	+35.0	13.0	74.0	+61.0	14.0	86.0	+72.0	13.0	89.0	+76.0	
LeetCode	3.9	3.9	+0.0	7.8	20.0	+12.2	6.1	70.0	+63.9	7.2	95.6	+88.4	7.2	97.8	+90.6	
RACE - Overall	20.3	15.6	-4.7	18.8	20.7	+1.9	19.3	35.8	+16.5	19.8	42.9	+23.1	19.1	43.8	+24.7	
RACE - IF rate	51.1	38.3	-12.8	53.9	37.4	-16.5	53.0	40.7	-12.3	53.4	41.9	-11.5	53.2	43.2	-10.0	

Table 3: The performance comparison of LLMs trained on clean data versus contamination data over the same number of epochs.

 Since current benchmarks use correctness as the sole guiding indicator, some LLMs perform well only on correctness but exhibit significant deficiencies in other dimensions. For example, Qwen2.5-Coder-7B-Ins demonstrates comparable levels of code correctness to GPT-4o-mini; however, GPT-4o-mini outperforms it by at least 5 percentage points regarding comments, modularity, and space complexity. These findings highlight the shortcomings of previous benchmarks and suggest that such deficiencies may be related to potential data leakage (see analysis in Section 4.3).

• The evaluation results reveal the importance of preserving the general instruction-following and language-understanding capabilities of code LLMs during training. This aligns with the direction of recent advancements and provides valuable insights for guiding further development. Specifically, on the one hand, current LLMs with the best coding abilities are often general-purpose LLMs, such as Claude-3.5-Sonnet, GPT-4o, and Qwen2.5-72B-Ins, all with overall scores exceeding 57. On the other hand, in the technical reports for Qwen2.5 (Qwen, 2024) and DeepSeek-Coder-V2 (Zhu et al., 2024), it is mentioned that a significant proportion of natural language corpora and general instruction data are included in training data. This approach not only enhances coding capabilities but also preserves general-purpose abilities. Moreover, both models achieve an overall score exceeding 57, outperforming most code LLMs.

These findings indicate that future research should prioritize improving instruction-following capabilities in terms of code readability, maintainability, and efficiency, while ensuring code accuracy remains uncompromised. This approach seeks to develop code LLMs that consistently meet realworld development requirements across multiple dimensions.

407 408 409

38 38 38

38

386

387

388

389

390

391

392

393

394

396

397

398

399

400

401

402 403

#### 4.3 ROBUSTNESS OF RACE ON DATA CONTAMINATION

410 Data contamination refers to the mixing of evaluation data into the training dataset, resulting in an 411 overly optimistic estimation of the model's true performance and leading to erroneous conclusions. 412 With the rapid expansion in training datasets for LLMs and the opacity surrounding critical information such as data sources and detailed data processing methods, addressing data contamination 413 is crucial for obtaining accurate and trustworthy evaluation results for LLMs. In this regard, RACE 414 plays a vital role, as it requires LLMs to generate code that is correct and meets user-specific cus-415 tomization demands across several dimensions. A contaminated model would only fit the data itself, 416 merely enhancing accuracy without improving the model's ability to follow user instructions. There-417 fore, it is intuitive to suggest that RACE can effectively mitigate the impact of data contamination. 418

To validate the RACE benchmark's robust-419 ness against data contamination, we select 420 starcoderbase-7b (Liet al., 2023a) as our 421 baseline. This model has been carefully cu-422 rated to exclude data from HumanEval (Chen 423 et al., 2021) and MBPP (Austin et al., 2021) 424 during its training, and ClassEval (Du et al., 425 2023) and LeetCode (Guo et al., 2024) are 426 not within the temporal coverage of its train-427 ing data. Furthermore, we compare the perfor-428 mance of models under clean data conditions 429 and varying levels of data contamination on the existing benchmarks and the RACE bench-430 mark. Specifically, we employ LoRA (Hu et al., 431 2022) to train models on both clean and con-



Figure 4: The variation in instruction-following capabilities under different factors in the context of data contamination.



Figure 5: The variation in the ability of different LLMs to generate code that is correct and satisfies all requirements as the number of requirements in complex instructions gradually increases from 2 (2-req) to 5 (5-req).

taminated datasets. In this case, the contaminated dataset consists of HumanEval+, MBPP+, ClassEval, and LeetCode data, while the clean dataset consists of an equivalent number of samples
randomly selected from Magicoder-OSS-Instruct (Wei et al., 2023). The model is trained for 10
epochs on the corresponding dataset, with a batch size of 32 and a learning rate of 1e-3. Additionally, starcoderbase-7b is pre-trained for 3 epochs on Code-Alpaca (Chaudhary, 2023) to
enhance instruction-following before performance comparisons.

The results presented in Table 3 clearly illustrate the role of the RACE benchmark in mitigating the 451 impact of data contamination on evaluation: 1) The original benchmarks are significantly affected 452 by severe data contamination. As the level of contamination increases, the code accuracy on the 453 corresponding benchmark rapidly rises. For instance, when trained on the contaminated dataset for 8 454 epochs, the accuracy of each original benchmark exceeds 85%. 2) The RACE benchmark provides 455 more stable evaluation results, and therefore effectively resist the risks of data contamination. 456 When data contamination is present, the instruction-following rate (IF rate) of the model on the 457 RACE benchmark consistently remains below 10% compared to the model without contamination. 458 As the degree of contamination increases, the rate of increase in the IF rate significantly slows down. 459 This phenomenon occurs because data contamination merely guides models to fit the data itself, 460 thereby improving the accuracy of generated code, which leads to a slow increase in the proportion of generated code that is both correct and meets user requirements (*Overall score*). However, this does 461 not contribute to improving the model's ability to follow user instructions. Therefore, the RACE 462 benchmark featured multidimensional evaluation effectively mitigates the risks associated with data 463 contamination. Furthermore, Figure 4 illustrates the variation in the model's instruction-following 464 capability across different factors as the degree of data contamination increases. It is evident that 465 data contamination significantly impairs the model's instruction-following ability in all factors. 466 For instance, in the case of code comments (RC), the model experiences a dramatic drop from 47%467 to 2.4% after only one epoch of data leakage. Notably, for factors related to naming conventions 468 (RN) and code length (RC), the model's performance first declines significantly and then gradually 469 improves, ultimately remaining considerably lower than its actual instruction-following capability.

470 471

442

443

444

4.4 COMPLEX INSTRUCTION FOLLOWING ABILITIES OF CODE LLMS

473 In real-world development scenarios, the requirements involved are often multifaceted. For example, 474 when developing a real-time data processing system, it is necessary to maintain code efficiency 475 while also ensuring the readability of the code through standardized variable names and comments. 476 To investigate the performance of code LLMs under complex requirements, we construct complex instructions based on LeetCode cases. Specifically, we randomly select K (where  $2 \le K \le 5$ ) 477 factors from the five factors: naming convention, length, comment, modularity, and efficiency, to 478 construct complex instructions with K customized requirements. Subsequently, we calculate the 479 proportion of generated code that is correct and meets all customized requirements (Acc.IF), serving 480 as the performance metric for the model's adherence to complex instructions. 481

The results from representative LLMs are illustrated in Figure 5. It is evident that as the number
 of requirements increases, the performance of all models gradually declines. Even the most ad vanced LLMs struggle to adequately satisfy all customized requirements within the complex
 instructions. For instance, when the number of requirements reaches 5 (5-req), almost all models
 exhibit a significant drop in Acc.IF rate, approaching 0, revealing a notable performance gap. We

488 489 490

493

494

495

496

497 498

499

500

501

502

491 492

Figure 6: The instruction-following rates of different LLMs for different customization needs. For naming convention, we adapt camel-case or snake-case for both function names and variable names. For code length, we set single-line lengths are limited to 60 and 79 characters, and methods are limited to 20 and 40 lines, respectively. For loop structure, we adapt for or while statements.

For the Length facto

attribute this deficiency to the fact that existing code benchmarks prioritize correctness, overlooking the models' capability to follow complex instructions. It is imperative for future work to investigate the underlying mechanisms of code LLMs when confronted with complex instructions, with the aim of enhancing their applicability in practical development contexts with multifaceted demands.

## 4.5 PREFERENCE BIAS OF CODE LLMS

To investigate whether the internal code preferences of the model affect its ability to follow user instructions, we conduct a more fine-grained comparison across various factors that are likely to induce such preferences. Specifically, we designed distinct instructions for naming convention, code length, and loop structure, respectively. Our objective is to observe whether the model exhibits a stronger capability to meet certain customization requirements. Furthermore, we calculate the proportion of LLM-generated code that follows these customized requirements, referred to as the Instruction-Following (IF) rate.

512 Figure 6 demonstrates the IF rates of 15 representative LLMs across all the customized requirements 513 above. Our analysis reveals that most LLMs exhibit an inherent preference bias towards gen-514 erating code in specific styles. This bias often leads to difficulties in following user instructions 515 when the requested style diverges from that prevalent in their training data. Specifically, for naming conventions, Python typically employs snake-case for function and variable names. When instructed 516 to use camel-case, most LLMs, such as CodeLlama and WizardCoder, almost fail to comprehend 517 and fulfill this requirement, with IF rates below 30%. Regarding code length, when presented with 518 stricter length constraints, the IF rates of most instruct-type LLMs drop by nearly 15%. When 519 it comes to loop structures, certain LLMs, including CodeLlama-34B-Ins and WizardCoder-33B, 520 exhibit a strong inclination towards using for statements. These observations suggest that many 521 LLMs primarily learn the inherent patterns of token prediction from examples, lacking a compre-522 hensive understanding of code logic. Such preference bias may result in the rigidification of coding 523 styles in code LLMs, ultimately impeding their ability to meet specific real-world requirements and 524 affecting the adaptability and scalability of the generated code. This issue could be even more pro-525 nounced in programming languages like Perl, JavaScript, and PHP, where there is no strict, widely accepted standard for coding styles. 526

527 528

## 5 CONCLUSION

529 530

We introduce the RACE benchmark, a comprehensive multi-dimensional evaluation framework 531 for code generation, including correctness, readability, maintainability, and efficiency. The RACE 532 benchmark assesses the ability of LLMs to generate code that is both correct and meets customized 533 requirements across different dimensions. Through extensive experiments involving 28 representa-534 tive LLMs, we find that the current code LLMs still fall short in generating high-quality code on demand. Moreover, the RACE benchmark serves as an effective tool for mitigating data contami-536 nation. Our research underscores the importance of enhancing code LLMs to generate high-quality code across multiple dimensions beyond mere correctness. Future work should prioritize developing novel learning algorithms to improve the coverage and usability of code LLMs in addressing diverse 538 user needs. This will enable the models to better handle complex instructions, ultimately guiding code LLMs towards becoming practical software development agents.

## 540 REFERENCES

566

567

568

578

579

580

581

- 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.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
- Jürgen Börstler, Kwabena E Bennin, Sara Hooshangi, Johan Jeuring, Hieke Keuning, Carsten Kleiner, Bonnie MacKellar, Rodrigo Duran, Harald Störrle, Daniel Toll, et al. Developers talking about code quality. *Empirical Software Engineering*, 28(6):128, 2023.
- <sup>551</sup>Boxi Cao, Mengjie Ren, Hongyu Lin, Xianpei Han, Feng Zhang, Junfeng Zhan, and Le Sun. StructEval: Deepen and broaden large language model assessment via structured evaluation. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Findings of the Association for Computational Linguistics ACL 2024*, pp. 5300–5318, Bangkok, Thailand and virtual meeting, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.314. URL https://aclanthology.org/2024.findings-acl.314.
- Sahil Chaudhary. Code alpaca: An instruction-following llama model for code generation. https:
   //github.com/sahil280114/codealpaca, 2023.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- 564 Don Coleman, Dan Ash, Bruce Lowther, and Paul Oman. Using metrics to evaluate software system 565 maintainability. *Computer*, 27(8):44–49, 1994.
  - Bill Curtis, Robert A Martin, and Philippe-Emmanuel Douziech. Measuring the structural quality of software systems. *Computer*, 55(3):87–90, 2022.
- 569 Carlos Eduardo C Dantas, Adriano M Rocha, and Marcelo A Maia. How do developers improve
   570 code readability? an empirical study of pull requests. In 2023 IEEE International Conference on
   571 Software Maintenance and Evolution (ICSME), pp. 110–122. IEEE, 2023.
- Marc Dillmann, Julien Siebert, and Adam Trendowicz. Evaluation of large language models for assessing code maintainability. *arXiv preprint arXiv:2401.12714*, 2024.
- 575 Xueying Du, Mingwei Liu, Kaixin Wang, Hanlin Wang, Junwei Liu, Yixuan Chen, Jiayi Feng,
  576 Chaofeng Sha, Xin Peng, and Yiling Lou. Classeval: A manually-crafted benchmark for evaluat577 ing llms on class-level code generation. *arXiv preprint arXiv:2308.01861*, 2023.
  - Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Angela Fan, Beliz Gokkaya, Mark Harman, Mitya Lyubarskiy, Shubho Sengupta, Shin Yoo, and Jie M Zhang. Large language models for software engineering: Survey and open problems. *arXiv preprint arXiv:2310.03533*, 2023.
- Daniel Fried, Armen Aghajanyan, Jessy Lin, Sida Wang, Eric Wallace, Freda Shi, Ruiqi Zhong,
   Scott Yih, Luke Zettlemoyer, and Mike Lewis. Incoder: A generative model for code infilling and
   synthesis. In *The Eleventh International Conference on Learning Representations*, 2022.
- Linyuan Gong, Sida Wang, Mostafa Elhoushi, and Alvin Cheung. Evaluation of Ilms on syntaxaware code fill-in-the-middle tasks. *arXiv preprint arXiv:2403.04814*, 2024.
- Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao
   Bi, Y Wu, YK Li, et al. Deepseek-coder: When the large language model meets programming-the rise of code intelligence. *arXiv preprint arXiv:2401.14196*, 2024.

624

626

627

594 Péter Hegedus. Revealing the effect of coding practices on software maintainability. In 2013 ieee international conference on software maintenance, pp. 578–581. IEEE, 2013. 596

- Ilia Heitlager, Tobias Kuipers, and Joost Visser. A practical model for measuring maintainability. 597 In Ricardo Jorge Machado, Fernando Brito e Abreu, and Paulo Rupino da Cunha (eds.), *Quality* 598 of Information and Communications Technology, 6th International Conference on the Quality of Information and Communications Technology, QUATIC 2007, Lisbon, Portugal, September 12-600 14, 2007, Proceedings, pp. 30–39. IEEE Computer Society, 2007. doi: 10.1109/QUATIC.2007.8. 601 URL https://doi.org/10.1109/QUATIC.2007.8. 602
- 603 Dan Hendrycks, Steven Basart, Saurav Kadavath, Mantas Mazeika, Akul Arora, Ethan Guo, Collin Burns, Samir Puranik, Horace He, Dawn Song, et al. Measuring coding challenge competence 604 with apps. In Thirty-fifth Conference on Neural Information Processing Systems Datasets and 605 Benchmarks Track (Round 2), 2021. 606
- 607 Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, 608 and Weizhu Chen. Lora: Low-rank adaptation of large language models. In The Tenth Inter-609 national Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022. 610 OpenReview.net, 2022. URL https://openreview.net/forum?id=nZeVKeeFYf9.
- Dong Huang, Jie M Zhang, Yuhao Qing, and Heming Cui. Effibench: Benchmarking the efficiency 612 of automatically generated code. arXiv preprint arXiv:2402.02037, 2024. 613
- 614 Binyuan Hui, Jian Yang, Zeyu Cui, Jiaxi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, 615 Bowen Yu, Kai Dang, et al. Qwen2. 5-coder technical report. arXiv preprint arXiv:2409.12186, 616 2024.
- 617 Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando 618 Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free 619 evaluation of large language models for code. arXiv preprint arXiv:2403.07974, 2024a. 620
- 621 Naman Jain, Manish Shetty, Tianjun Zhang, King Han, Koushik Sen, and Ion Stoica. R2e: Turning 622 any github repository into a programming agent environment. In ICLR 2024 Workshop on Large Language Model (LLM) Agents, 2024b. 623
- Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bam-625 ford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. Mixtral of experts. arXiv preprint arXiv:2401.04088, 2024.
- Yuhang Lai, Chengxi Li, Yiming Wang, Tianyi Zhang, Ruiqi Zhong, Luke Zettlemoyer, Wen-tau 628 Yih, Daniel Fried, Sida Wang, and Tao Yu. Ds-1000: A natural and reliable benchmark for data 629 science code generation. In International Conference on Machine Learning, pp. 18319–18345. 630 PMLR, 2023. 631
- 632 Bowen Li, Wenhan Wu, Ziwei Tang, Lin Shi, John Yang, Jinyang Li, Shunyu Yao, Chen Qian, 633 Binyuan Hui, Qicheng Zhang, et al. Devbench: A comprehensive benchmark for software devel-634 opment. arXiv preprint arXiv:2403.08604, 2024.
- 635 Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou, 636 Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, et al. Starcoder: may the source be with 637 you! arXiv preprint arXiv:2305.06161, 2023a. 638
- 639 Rongao Li, Jie Fu, Bo-Wen Zhang, Tao Huang, Zhihong Sun, Chen Lyu, Guang Liu, Zhi Jin, and Ge Li. Taco: Topics in algorithmic code generation dataset. arXiv preprint arXiv:2312.14852, 640 2023b. 641
- 642 Y Li, D Choi, J Chung, N Kushman, J Schrittwieser, R Leblond, T Eccles, J Keeling, F Gimeno, 643 A Dal Lago, et al. Competition-level code generation with alphacode. Science (New York, NY), 644 378(6624):1092-1097, 2022. 645
- Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. Is your code generated by chat-646 gpt really correct? rigorous evaluation of large language models for code generation. Advances 647 in Neural Information Processing Systems, 36, 2024.

0.40	
648	Anton Lozhkov, Raymond Li, Loubna Ben Allal, Federico Cassano, Joel Lamy-Poirier, Nouamane
649	Tazi, Ao Tang, Dmytro Pykhtar, Jiawei Liu, Yuxiang Wei, et al. Starcoder 2 and the stack v2: The
650	next generation arXiv preprint arXiv:2402 19173 2024
651	
652	Zivang Luo, Can Xu, Pu Zhao, Oingfeng Sun, Xiubo Geng, Wenxiang Hu, Chongyang Tao, Jing
652	Ma Oingwei Lin, and Daxin Jiang. Wizardcoder: Empowering code large language models with
003	evolutionstruct. In The Twelfth International Conference on Learning Representations 2024 URL
654	bitns. //openreview net/forum2id=UnUwStak5W
655	https://openreview.net/iordm.id=onowsigkow.
656	Ansong Ni Pengcheng Yin, Yilun Zhao, Martin Riddell, Troy Feng, Rui Shen, Stephen Yin, Ye Liu
657	Semi Yayuz Caiming Xiong et al L2ceval: Evaluating language-to-code generation canabilities
658	of large language models, arXiv preprint arXiv:2309.17446, 2023
659	
022	Erik Nijkamp, Bo Pang, Hiroaki Havashi, Lifu Tu, Huan Wang, Yingbo Zhou, Silvio Savarese,
000	and Caiming Xiong. Codegen: An open large language model for code with multi-turn program
661	synthesis In The Eleventh International Conference on Learning Representations 2022
662	
663	Padmalata Nistala, Kesay Vithal Nori, and Raghu Reddy, Software quality models: A systematic
664	mapping study. In 2019 IEEE/ACM International Conference on Software and System Processes
665	( <i>ICSSP</i> ), pp. 125–134, IEEE, 2019.
666	( ····· /) [[····························
667	Delano Oliveira, Reydne Bruno, Fernanda Madeiral, and Fernando Castor. Evaluating code read-
862	ability and legibility: An examination of human-centric studies. In 2020 IEEE International
000	Conference on Software Maintenance and Evolution (ICSME), pp. 348–359, IEEE, 2020.
669	
670	<b>OpenAI. Hello gpt-40, 2024a. URL</b> https://openai.com/index/hello-gpt-40.
671	
672	OpenAI. Openai o1-mini, 2024b. URL https://openai.com/index/
673	openai-ol-mini-advancing-cost-efficient-reasoning.
674	
675	Qwen. Qwen2.5: A party of foundation models, September 2024. URL https://qwenlm.
676	github.io/blog/qwen2.5/.
070	
077	Martin Riddell, Ansong Ni, and Arman Cohan. Quantifying contamination in evaluating code gen-
678	eration capabilities of language models. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar
679	(eds.), Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics
680	(Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024, pp. 14116–14137.
681	Association for Computational Linguistics, 2024. doi: 10.18653/V1/2024.ACL-LONG./61. URL
682	https://doi.org/10.18653/v1/2024.acl-long.761.
683	
684	Baptiste Roziere, Jonas Genring, Fabian Gioeckie, Sten Sootia, Ital Gat, Xlaoqing Ellen Ian, Yossi
695	Adi, Jingyu Liu, Tai Remez, Jeremy Rapin, et al. Code liama: Open foundation models for code.
005	arXiv preprint arXiv:2308.12950, 2023.
080	M Sadaghzadah Hamayati and H Pashidi. Softwara quality models: A comprehensive review and
687	analysis Journal of Electrical and Computer Engineering Impositions (IECEI) 6(1):50-76
688	anarysis. Journal of Electrical and Computer Engineering Innovations (JECEI), 0(1):59-70, 2017
689	2017.
690	Rylan Schaeffer, Brando Miranda, and Sanmi Koveio. Are emergent abilities of large language
691	models a mirage? Advances in Neural Information Processing Systems 36, 2024
692	models a minage. Mavances in Neural Information Processing Systems, 50, 2024.
693	Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion:
60/	Language agents with verbal reinforcement learning. Advances in Neural Information Processing
COF	Systems, 36, 2024.
690	~,~~~~, = •, = •= ••
696	Runchu Tian, Yining Ye, Yujia Qin, Xin Cong, Yankai Lin, Zhivuan Liu, and Maosong Sun.
697	Debugbench: Evaluating debugging capability of large language models. arXiv preprint
698	arXiv:2401.04621, 2024.
699	
700	Yue Wang, Weishi Wang, Shafiq Joty, and Steven CH Hoi. Codet5: Identifier-aware unified pre-
701	trained encoder-decoder models for code understanding and generation. In <i>Proceedings of the</i> 2021 Conference on Empirical Methods in Natural Language Processing, pp. 8696–8708, 2021.

702 703 704	Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yo- gatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language models. <i>Transactions on Machine Learning Research</i> , 2022.
705 706	Yuxiang Wei, Zhe Wang, Jiawei Liu, Yifeng Ding, and Lingming Zhang. Magicoder: Source code is all you need. <i>arXiv preprint arXiv:2312.02120</i> , 2023.
707 708 709 710	<ul> <li>Frank F Xu, Uri Alon, Graham Neubig, and Vincent Josua Hellendoorn. A systematic evaluation of large language models of code. In <i>Proceedings of the 6th ACM SIGPLAN International Symposium on Machine Programming</i>, pp. 1–10, 2022.</li> </ul>
711 712 713 714	Weixiang Yan, Haitian Liu, Yunkun Wang, Yunzhe Li, Qian Chen, Wen Wang, Tingyu Lin, Weishan Zhao, Li Zhu, Shuiguang Deng, et al. Codescope: An execution-based multilingual multitask multidimensional benchmark for evaluating llms on code understanding and generation. <i>arXiv</i> preprint arXiv:2311.08588, 2023.
715 716 717 718	An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, et al. Qwen2 technical report. <i>arXiv preprint arXiv:2407.10671</i> , 2024.
719 720 721	Qinkai Zheng, Xiao Xia, Xu Zou, Yuxiao Dong, Shan Wang, Yufei Xue, Zihan Wang, Lei Shen, Andi Wang, Yang Li, et al. Codegeex: A pre-trained model for code generation with multilingual evaluations on humaneval-x. <i>arXiv preprint arXiv:2303.17568</i> , 2023a.
722 723 724 725	Zibin Zheng, Kaiwen Ning, Jiachi Chen, Yanlin Wang, Wenqing Chen, Lianghong Guo, and We- icheng Wang. Towards an understanding of large language models in software engineering tasks. <i>arXiv preprint arXiv:2308.11396</i> , 2023b.
726 727 728	Zibin Zheng, Kaiwen Ning, Yanlin Wang, Jingwen Zhang, Dewu Zheng, Mingxi Ye, and Jiachi Chen. A survey of large language models for code: Evolution, benchmarking, and future trends. <i>arXiv preprint arXiv:2311.10372</i> , 2023c.
729 730 731 732	Qihao Zhu, Daya Guo, Zhihong Shao, Dejian Yang, Peiyi Wang, Runxin Xu, Y Wu, Yukun Li, Huazuo Gao, Shirong Ma, et al. Deepseek-coder-v2: Breaking the barrier of closed-source models in code intelligence. <i>arXiv preprint arXiv:2406.11931</i> , 2024.
733 734 735	
736 737 738	
739 740	
741 742 743	
744 745	
746 747 748	
749 750 751	
752 753	
754 755	



## A CORRELATION ANALYSIS ACROSS DIMENSIONS

Figure 7: The Pearson correlation coefficient matrix among factors under the dimensions of code correctness, readability, maintainability, and efficiency. We can observe that readability is a critical indicator of overall code quality.



Figure 8: Comparison of code correctness among LLM-generated code without custom requirements, with function-level comments, and with line-level comments.

To conduct a more in-depth analysis of how different factors across various dimensions influence overall code quality, we analyze the correlations between different factors across all models. Specifically, we first compute the proportion of the generated code that is both correct and follows customized instructions across 8 factors for 28 Code LLMs. Subsequently, we calculate the Pearson correlation coefficients between these factors.

The results of the correlation analysis are presented in Figure 7, which demonstrate that **readability** serves as a critical indicator of overall code quality. Notably, significant correlations are observed between readability and nearly all the factors, with most correlation coefficients exceeding 0.8, and particularly surpassing 0.85 in relation to correctness. For example, if LLM-generated code consistently uses proper naming conventions, follows length constraints, and includes appropriate comments, it is more likely to exhibit higher overall quality. This finding aligns with the conclu810 sions of Börstler et al. (2023), which identifies readability as a decisive factor in code quality, 811 suggesting that improving the readability of LLM-generated code is a critical path for enhancement. 812 Furthermore, we analyze the comments factor within the readability dimension and compare the 813 accuracy of LLM-generated code before and after incorporating comments. As shown in Figure 8, 814 requiring models to include comments in appropriate sections of code enhances the performance of some LLMs. We hypothesize that this improvement may be attributed to an emerging ability in 815 large-scale LLMs (Wei et al., 2022; Schaeffer et al., 2024), where comments serve as an implicit 816 chain-of-thought mechanism, thereby enhancing the accuracy of the generated code. 817

818 819

820

824

**B** EXPERIMENT SETUP

## 821 B.1 MODEL SHORT NAMES 822

823 We demonstrate the details of LLMs in our experiment in Table 4.

825	Table 4: The short names of all L	LMs in the experiment.
826		
827	Model ID	Short Name
828	claude-3.5-sonnet	Claude-3.5-Sonnet
829	gpt-4o-2024-05-13	GPT-40
830	gpt-4o-mini	GPT-4o-mini
831	gpt-3.5-turbo-0125	GPT-3.5-turbo-0125
832	o1-mini-2024-09-12	o1-mini-2024-09-12
833	CodeLlama-7b-Python-hf	CL-7B-Py
834	CodeLlama-7b-Instruct-hf	CL-7B-Ins
835	CodeLlama-13b-Python-hf	CL-13B-Py
836	CodeLlama-13b-Instruct-hf	CL-13B-Ins
837	CodeLlama-34b-Python-hf	CL-34B-Py
000	CodeLlama-34b-Instruct-hf	CL-34B-Ins
838	WizardCoder-15B-V1.0	WC-15B
839	WizardCoder-33B-V1.1	WC-33B
840	WizardCoder-Python-7B-V1.0	WC-Py-7B
841	WizardCoder-Python-13B-V1.0	WC-Py-13B
842	deepseek-coder-6.7b-instruct	DS-Coder-6.7B-Ins
843	deepseek-coder-7b-instruct-v1.5	DS-Coder-7B-Ins
844	deepseek-coder-33b-instruct	DS-Coder-33B-Ins
845	DeepSeek-Coder-V2-Lite-Instruct	DS-Coder-V2-16B-Ins
846	deepseek-v2.5	DS-V2.5-236B
847	CodeQwen1.5-7B-Chat	CodeQwen1.5-7B-Chat
8/18	Qwen2.5-Coder-7B-Instruct	Qwen2.5-Coder-7B-Ins
040	Qwen2-72B-Instruct	Qwen2-72B-Ins
049	Qwen2.5-72B-Instruct	Qwen2.5-72B-Ins
850	mixtral-8x22b	Mixtral-8x22B
851	Meta-Llama-3-8B-Instruct	Llama3-8B-Ins
852	Meta-Llama-3-70B-Instruct	Llama3-70B-Ins
853	starcoder2-15b	StarCoder2-15B
054		

855 856

857

#### B.2 EVALUATION DATA AND CUSTOMIZED INSTRUCTIONS

Based on widely recognized data, we design customized requirements that are both reasonable and closely aligned with real-world application scenarios. These requirements are incorporated into the task descriptions to generate evaluation data for our RACE benchmark. Detailed customization instructions for each factor are shown in Figure 9 and Figure 10.

For code correctness, we utilize data from HumanEval+ (Liu et al., 2024), MBPP+ (Liu et al., 2024),
ClassEval (Du et al., 2023), and LeetCode (Guo et al., 2024). Code readability is evaluated using HumanEval+ (Liu et al., 2024) data, while maintainability is evaluated using ClassEval (Du et al., 2024)

2023) and LeetCode (Guo et al., 2024) data. Code efficiency is measured using a self-constructed dataset based on LeetCode problems. We follow the task settings defined in the original datasets while incorporating our customization requirements. In the case of the HumanEval+ and MBPP+ datasets (Liu et al., 2024), we modify the original prompt format by extracting the core task descriptions to serve as the final prompts. This adjustment helps prevent conflicts between function template information in the original prompts and our requirements for code readability, providing a more accurate assessment of code-related capabilities. Additionally, it mitigates the potential for data leakage, thereby increasing the difficulty and robustness of the RACE benchmark.

875 A) The templates for the correctness dimension 876 Please generate the Python code to solve the following problem. $n\nProblem:\n\n{problem}$ 877 878 B) The templates for the readability dimension 879 1) For the Naming Convention factor 881 Please generate the Python code to solve the following problem, and use camel case for both function 882 names and variable names. $n\nProblem:\n\n{problem}$ 883 Please generate the Python code to solve the following problem, and use snake case for both function 885 names and variable names. $n\nProblem:\n\n{problem}$ 886 Please generate the Python code to solve the following problem, and use camel case for function 887 names. $n\nProblem:\n\n{problem}$ 888 889 Please generate the Python code to solve the following problem, and use snake case for function 890 names. $n\nProblem:\n\problem$ 891 Please generate the Python code to solve the following problem, and use camel case for variable 892 names. $n\nProblem:\n\n{problem}$ 893 894 Please generate the Python code to solve the following problem, and use snake case for variable names. $n\nProblem:\n\problem$ 895 896 2) For the Length factor 897 Please generate the Python code to solve the following problem, where each line is less than 60 899 characters long and each function is less than 20 lines long. \n\nProblem: \n\n{problem} 900 Please generate the Python code to solve the following problem, where each line is less than 70 901 characters long and each function is less than 30 lines long. $\n\ensuremath{nnclenslaml}$ 902 903 Please generate the Python code to solve the following problem, where each line is less than 79 904 characters long and each function is less than 40 lines long.\n\nProblem:\n\n{problem} 905 3) For the Comment factor 906 907 Please generate the Python code to solve the following problem, and add the necessary docstring for 908 each function. $n\nProblem:\n\n{problem}$ 909 Please generate the Python code to solve the following problem, and add comments for each line in 910 each function. $n\nProblem:\n\n{problem}$ 911 912 913 914 Figure 9: The prompt templates for each factor in the correctness and readability dimension for the

915 RACE benchmark. 916

917

	C) The templates for the maintainability dimension
	1) For the MI factor
	Please complete the class {class_name} in the following code, and ensure that the code has
	good maintainability. Code maintainability refers to how easy it is to support and change the
	code. $\ln \ln \exp(\ln \ln $
	2) For the Modularity factor
	{problem} $n\n$ problem, and use only the given function $n \ n$
	{problem}\n\nPlease complete the code below to solve above problem, and use only the given
	function and one addition sub-function.\n\n{starter_code}
	{problem} $n$ , n/nPlease complete the code below to solve above problem, and use only the given function and two addition sub-functions $n$
	runction and two addition sub-runctions. In Infistance_code
	3) For the loop structure (Only for experiments)
	Please generate the Python code to solve the following problem, and just use the for statement to
	implement the desired loop structures.\n\nProblem:\n\n{problem}
	Please generate the Python code to solve the following problem, and just use the while statement to
	implement the desired loop structures. $n\nProblem:\n\nFroblem}$
	D) The templates for the efficiency dimension
	$problem} n please complete the code below to solve above problem, and make sure that the time$
	complexity of the code is \${complexity}\$.\n\n{starter_code}
	{problem}\n\nPlease complete the code below to solve above problem, and make sure that the space
	complexity of the code is $complexity$ .\n\n{starter_code}
	{problem}\n\nPlease complete the code below to solve above problem, and make
	sure that the time complexity is \${time_complexity}\$ and the space complexity is
	\${space_complexity}\$.\n\n{starter_code}
	)
g	ure 10: The prompt templates for each factor in the maintainability and efficiency dimension for
e	KACE DENCHMARK.
ı	EXPERIMENTAL RESULTS
	LAI ERIMENTAL RESULTS
1	OVERALL RESULTS WITH DETAILED CODE ACCUPACY
. 1	OVERALL RESULTS WITH DETAILED CODE ACCURACT
he	experimental results with the code accuracy before and after incorporating customization in-
ru	ctions are presented in Table 5.
	•

964 C.2 DETAILED RESULTS ON CODE READABILITY

The detailed experimental results under all customized instructions for the various readability factors are presented in Table 6 and Table 7. For the Naming Convention factor, we design 6 settings that require generated code to follow specified naming conventions for function names (function\_camel, function\_snake), variable names (var\_camel, var\_snake), or both (camel, snake). It is evident that most models struggle to consistently follow the camel-case naming convention. Additionally, the variance in model performance is most pronounced in scenarios requiring function names to follow camel-case conventions (function\_camel). For the Length factor, we observe that as the constraints become more stringent, ranging from maximum single-line length of 79 and maximum method line

Table 5: Based on the RACE benchmark, the performance results for each LLM in code correctness (C), readability (R), maintainability (M), and efficiency (E). The performance metrics include accu-racy (Acc) (%) and the proportion of code that is both functionally correct and follows customized instructions (Acc. IF) (%). RN, RL, RC, and EC denote the Name Convention, Length, Comments, and Complexity factor. MI denotes the Maintainability Index. MC denotes the Modularity factor.  $NI_T$  and  $NI_S$  are metrics for code efficiency. RACE Score represents the overall Acc. IF values at the dimension level. The (\*) symbol indicates that the indicator is a scalar from 0 to 100, and the rest are percentages (%). The symbol  $(\dagger)$  indicates that the results are obtained from a randomly sampled 30% of the evaluation data, in order to optimize budget efficiency. 

981		RACE	Correctness			I	Readabi	lity					Main	tainabili		Efficiency				
982		-	С	С		RN		RL		RC	С	Ν	11	С	]	MC	С		EC	
000	Models	Overall	Acc.	Acc.	Acc.	Acc. IF	Acc.	Acc. IF	Acc.	Acc. IF	Acc.	Acc.	MI*	Acc.	Acc.	Acc. IF	Acc.	Acc.	$NI_T^*$	$NI_S^*$
983								Inst	pe	e										
984	Claude-3.5-Sonnet GPT-40	62.3 57.2	<u>64.6</u> 59.9	77.4 80.5	76.3 81.2	74.4 78.6	62.2 78 9	52.0 63.2	74.1 79.8	65.5 70.4	42.0 38.0	32.0 35.0	75.3 75.1	71.7	68.5 56.3	59.8 35.2	68.3 59.4	66.3 58.4	56.8 44.0	49.7 42.0
985	GPT-4o-mini GPT 3.5 turbo 0125	52.5	56.4	78.0	76.4	67.6	70.3	55.7	74.1	72.9	37.0	27.0	73.5	51.7	49.1	23.3	52.5	46.5	40.3	39.5
986	o1-mini-2024-09-12	63.5 23.2	70.1	82.9 32.3	83.2 31.5	80.7 17.8	76.4	47.5	80.2 30.2	77.7	36.0 16.0	25.0	64.4 71.8	<b>79.6</b> <sup>†</sup>	83.3 <sup>†</sup>	66.1 <sup>†</sup>	87.1 <sup>†</sup>	77.4 <sup>†</sup>	60.3 <sup>†</sup>	40.0 <sup>†</sup>
987	CL-13B-Ins CL-34B-Ins	26.9 24.4	24.4 26.0	36.0 36.0	37.7 36.5	22.9 21.9	35.0 35.8	23.6 17.5	35.7 36.3	29.0 10.7	17.0 12.0	19.0 18.0	82.1 73.2	10.6 15.6	13.1 14.2	7.6 8.5	17.8 20.8	17.8 15.8	10.4 14.4	16.1 13.8
988	DS-Coder-6.7B-Ins DS-Coder-7B-Ins	39.8 38.9	39.2 39.9	65.2 61.0	65.5 61.5	45.8 36.8	61.2 62.6	46.6 46.0	61.2 62.8	50.0 53.7	26.0 23.0	25.0 24.0	79.3 79.6	18.9 23.3	18.7 20.9	8.2 8.9	28.7 32.7	30.7 27.7	27.1 25.1	30.0 26.8
989	DS-Coder-33B-Ins DS-Coder-V2-16B-Ins DS-V2 5-236B	44.8 48.2 57.1	44.7 50.9 59.0	65.9 72.0 72.0	64.6 71.2 74.5	59.0 41.8 72.2	65.0 66.5 72.8	53.5 57.7 66.1	66.5 67.1 74.1	54.0 47.5 65.8	28.0 26.0 41.0	30.0 30.0 36.0	75.7 78.2 72.9	22.2 44.4 61.7	27.6 44.3 59.1	11.3 19.8 33.9	45.5 49.5 57.4	38.6 55.4 54.5	35.3 40.2 46.4	36.1 47.7 49.5
990	CodeQwen1.5-7B-Chat Qwen2.5-Coder-7B-Ins	45.2 49.0	46.3 57.1	76.2 78.0	76.8 81.4	48.8 53.0	73.4 77.4	47.0 51.8	74.7 75.3	62.2 61.3	22.0 29.0	22.0 27.0	82.3 78.6	33.3 54.4	32.6 50.4	13.0 17.6	39.6 59.4	38.6 48.5	30.7 37.0	37.7 33.7
991	Qwen2-72B-Ins Qwen2.5-72B-Ins	50.1 61.3	53.1 64.1	73.2 79.3	76.8 79.6	73.6 77.2	74.8 77.4	47.6 72.1	71.1 80.5	60.1 72.8	40.0 34.0	33.0 32.0	79.4 76.7	42.8 <u>72.8</u>	37.2 <u>71.8</u>	22.8 40.4	45.5 <u>68.3</u>	40.6 <u>69.3</u>	32.3 47.9	39.4 49.4
992	Mixtral-8x22B Llama3-8B-Ins	42.2 35.2	42.0 35.6	61.0 49.4	64.4 45.5	56.2 44.3	62.4 28.7	47.8 23.6	64.9 48.1	56.1 40.0	33.0 24.0	30.0 19.0	79.6 79.8	20.0 20.6	22.6 19.1	9.1 8.1	35.6 33.7	31.7 31.7	24.7 23.5	33.2 26.9
993	Llama3-70B-Ins	47.2	44.4	65.2	67.8	66.0	56.1	47.8 Comr	64.6	54.2 Гуре	28.0	29.0	79.8	31.7	31.7	25.2	38.6	38.6	29.2	42.8
994	CL-7B-Py	24.0	20.4	29.3	29.5	20.9	30.1	25.8	24.7	12.5	11.0	10.0	79.4	5.6	6.5	3.7	14.9	15.8	14.3	14.4
995	CL-13B-Py CL-34B-Py	25.6	19.2	40.2 31.7	27.2	18.8	34.8 32.5	26.7	27.8	24.4 8.6	3.0	2.0	78.0 85.3	7.2	4.8 5.4	2.4	17.8	17.8	13.8	14.7
996	WC-Py-7B WC-Py-13B WC-15B	20.2 29.3 30.4	25.2 26.3 28.0	34.8 36.0 38.4	35.8 38.2 38.7	22.8 23.9 24.0	34.5 38.4 41.9	28.0 33.1 27.8	35.4 43.6 40.0	30.5 28.1	20.0 22.0	23.0 21.0 21.0	79.5 78.8 80.0	10.6 12.8 11.7	9.8 12.8 11.5	7.2 8.5 7.8	20.8 21.8	19.8 18.8 22.8	15.5 16.2 21.8	16.7 19.8 24.2
997	WC-33B StarCoder2-15B	40.8 29.2	44.4 28.5	58.5 36.0	58.8 39.5	40.9 25.8	62.2 40.2	47.6 27.9	58.8 35.4	44.8 22.0	34.0 24.0	34.0 25.0	71.2 74.2	26.1 16.1	25.0 13.7	9.3 6.1	38.6 26.7	35.6 25.7	33.9 20.6	34.9 25.1

count of 40 (L\_79\_40), to maximum single-line length of 60 and maximum method line count of 20 (L\_60\_20), most models exhibit a significant decline in their ability to meet requirements. For the Comment factor, models show varying responses to comment-related requirements. However, we find that several models, such as DS-Coder-33B-Ins, DS-V2.5-236B, and WC-Py-13B, improve in code correctness when they meet the code comment requirements.

										Reada	ability (Na	ning Co	onvention	1)						
		С		came	:1		snak	e	fu	nction_	camel	fı	inction_s	nake		var_ca	nel		var_sn	ake
	Models	Acc.	Acc.	IF	Acc. IF	Acc.	IF	Acc. IF	Acc.	IF	Acc. IF	Acc.	IF	Acc. IF	Acc.	IF	Acc. IF	Acc.	IF	Acc. IF
	Claude-3.5-Sonnet	77.4	75.6	90.2	70.1	76.8	97.0	76.2	78.7	98.8	78.7	77.4	98.8	77.4	74.4	90.9	69.5	75.0	97.0	74.4
	GPT-40	80.5	81.7	89.6	73.8	80.5	97.6	79.9	84.1	98.8	83.5	79.9	99.4	79.9	81.7	90.2	75.0	79.3	98.2	79.3
	GPT-4o-mini	78.0	73.2	42.7	34.1	76.2	98.8	76.2	75.6	86.6	65.2	78.0	99.4	78.0	79.9	95.1	77.4	75.6	99.4	75.0
	GPT-3.5-turbo-0125	62.8	65.2	51.8	37.8	61.0	97.6	59.8	63.4	87.2	56.1	62.8	98.8	62.8	64.0	41.5	30.5	62.8	98.2	61.6
	CL-7B-Py	29.3	28.0	18.3	4.9	29.9	98.8	29.9	29.9	22.0	7.3	31.7	98.2	31.7	28.7	78.7	23.2	28.7	98.2	28.7
	CL-7B-Ins	32.3	31.7	2.4	0.0	29.3	98.2	29.3	31.1	5.5	0.0	31.7	100.0	31.7	31.7	43.9	12.2	33.5	99.4	33.5
	CL-13B-Py	40.2	34.8	6.7	2.4	35.4	97.0	35.4	36.0	9.1	3.7	33.5	97.6	33.5	34.8	73.2	28.0	35.4	98.2	35.4
	CL-13B-Ins	36.0	37.2	4.3	3.0	37.2	99.4	37.2	40.9	9.8	5.5	34.8	99.4	34.8	40.2	48.8	20.7	36.0	99.4	36.0
	CL-34B-Py	31.7	27.4	19.5	5.5	28.0	95.7	27.4	29.9	21.3	6.1	26.2	97.6	26.2	26.2	80.5	22.6	25.6	97.0	25.0
	CL-34B-Ins	36.0	37.2	4.3	2.4	34.8	92.1	34.8	36.6	5.5	2.4	36.6	97.0	36.6	37.8	47.6	19.5	36.0	94.5	36.0
	WC-15B	38.4	39.6	4.3	1.2	40.9	98.2	40.9	38.4	5.5	1.2	38.4	97.6	38.4	39.0	62.8	26.2	36.0	97.6	36.0
	WC-33B	58.5	57.9	25.0	14.6	59.1	95.1	57.3	57.3	34.1	20.7	57.9	97.6	57.9	59.8	60.4	35.4	61.0	95.7	59.8
	WC-Py-7B	34.8	34.8	4.9	1.8	34.1	95.7	34.1	34.8	5.5	1.2	34.1	97.6	34.1	37.8	62.8	26.8	39.0	94.5	39.0
	WC-Py-13B	36.0	38.4	4.3	1.8	36.6	97.0	36.6	36.6	6.1	1.2	38.4	97.6	38.4	37.8	59.8	23.8	41.5	96.3	41.5
	DS-Coder-6.7B-Ins	65.2	65.2	26.2	15.9	65.9	97.6	64.6	67.7	47.0	29.9	67.7	100.0	67.7	62.8	48.2	33.5	64.0	98.2	63.4
	DS-Coder-7B-Ins	61.0	61.6	9.1	6.1	59.1	99.4	58.5	62.2	11.6	7.3	61.6	100.0	61.6	62.8	43.9	26.2	61.6	98.8	61.0
	DS-Coder-33B-Ins	65.9	64.6	73.2	51.2	65.2	97.0	64.6	62.2	99.4	61.6	64.0	100.0	64.0	68.3	73.2	50.0	63.4	97.6	62.8
	DS-Coder-V2-16B-Ins	72.0	72.0	9.8	7.9	69.5	95.1	67.7	72.6	14.0	10.4	71.3	99.4	71.3	73.2	33.5	26.2	68.9	95.1	67.1
	DS-V2.5-236B	72.0	75.0	89.0	67.7	76.2	98.8	75.6	74.4	98.8	74.4	75.0	99.4	75.0	72.0	90.9	67.1	74.4	97.6	73.2
	CodeQwen1.5-7B-Chat	76.2	75.6	12.2	9.1	76.2	97.6	75.0	76.2	15.9	11.0	79.3	99.4	78.7	76.8	57.9	43.3	76.8	96.3	75.6
	Qwen2.5-Coder-7B-Ins	78.0	81.1	17.1	14.6	81.1	97.6	78.7	82.9	36.0	31.1	81.7	100.0	81.7	80.5	40.9	32.9	81.1	97.6	78.7
	Qwen2-72B-Ins	73.2	75.6	90.2	68.3	78.7	98.2	78.0	75.6	93.9	69.5	79.9	100.0	79.9	74.4	95.7	70.7	76.8	97.6	75.0
	Qwen2.5-72B-Ins	79.3	78.7	94.5	74.4	81.1	97.0	78.0	79.9	99.4	79.3	80.5	100.0	80.5	78.0	93.3	73.8	79.3	97.6	77.4
	Mixtral-8x22B	61.0	65.2	62.2	43.3	64.6	99.4	64.0	65.9	96.3	62.2	64.0	100.0	64.0	62.8	65.9	40.2	64.0	98.2	63.4
	Llama3-8B-Ins	49.4	49.4	88.4	47.0	31.7	57.9	30.5	51.8	97.6	51.8	44.5	89.6	44.5	50.6	90.2	48.8	45.1	89.6	43.3
	Llama3-70B-Ins	65.2	70.1	93.9	65.9	65.9	97.6	64.6	68.3	99.4	68.3	66.5	100.0	66.5	68.9	91.5	64.0	67.1	97.6	66.5
	StarCoder2-15B	36.0	41.5	11.6	4.3	38.4	96.3	38.4	38.4	14.0	5.5	38.4	96.3	38.4	42.1	72.0	30.5	38.4	95.7	37.8

#### 1026 Table 6: Detailed experimental results for the Name Convention factor in the readability dimension 1027 on the RACE benchmark.

1045 Table 7: Detailed experimental results for the Length and Comment factor in the readability dimen-1046 sion on the RACE benchmark.

1048		Readability (Length)												Readability (Comment)							Maintainability (Loop Structure)					
1040		С		L_60_	20		L_70_	30		L_79_	40		by_func	tion		by_lii	ne	for				while	•			
1049	Models	Acc.	Acc.	IF	Acc. IF	Acc.	IF	Acc. IF	Acc.	IF	Acc. IF	Acc.	IF	Acc. IF	Acc.	IF	Acc. IF	Acc.	IF	Acc. IF	Acc.	IF	Acc. IF			
1050	Claude-3.5-Sonnet	77.4	50.0	49.4	37.2	67.7	78.0	57.9	68.9	83.5	61.0	75.0	95.7	75.0	73.2	74.4	56.1	70.1	98.2	68.3	66.5	98.8	65.9			
1050	GPT-40	80.5	80.5	74.4	61.6	76.2	75.0	58.5	79.9	87.2	69.5	77.4	98.2	77.4	82.3	76.8	63.4	75.0	93.3	71.3	70.1	97.0	68.9			
	GPT-4o-mini	78.0	70.7	68.3	51.8	73.2	72.0	55.5	67.1	84.1	59.8	73.8	98.8	73.8	74.4	95.1	72.0	71.3	94.5	67.1	65.2	97.0	64.6			
1051	GPT-3.5-turbo-0125	62.8	58.5	64.6	39.0	61.0	78.0	45.7	61.6	87.8	53.7	66.5	95.1	64.0	65.2	45.1	31.1	56.7	97.0	54.9	52.4	90.2	48.8			
1001	CL-7B-Py	29.3	29.9	68.3	23.8	31.1	78.0	25.6	29.3	83.5	28.0	28.7	72.6	22.0	20.7	43.3	3.0	28.7	93.9	26.2	26.8	58.5	13.4			
1050	CL-7B-Ins	32.3	29.9	50.6	20.1	31.7	57.9	23.2	33.5	70.7	26.8	29.9	100.0	29.9	30.5	52.4	14.6	31.7	95.7	31.1	29.9	42.7	11.0			
1052	CL-13B-Py	40.2	34.1	78.7	27.4	34.8	83.5	31.1	35.4	88.4	34.1	34.8	92.1	34.8	25.6	62.8	14.0	33.5	91.5	31.1	34.8	60.4	18.9			
	CL-13B-Ins	36.0	34.8	53.0	20.1	35.4	62.8	25.6	34.8	64.0	25.0	36.6	92.7	34.8	34.8	57.3	23.2	31.1	95.1	30.5	34.1	45.7	14.0			
1053	CL-34B-Py	31.7	32.3	64.6	22.0	33.5	72.6	28.7	31.7	82.3	29.3	23.2	67.7	12.8	32.3	29.9	4.3	25.6	94.5	25.0	25.6	70.7	12.2			
	CL-34B-Ins	36.0	33.5	36.0	15.9	36.6	39.0	15.9	37.2	50.0	20.7	35.4	38.4	12.8	37.2	34.1	8.5	36.0	94.5	36.0	35.4	46.3	15.2			
1054	WC-15B	38.4	42.7	50.0	20.7	40.2	67.1	28.0	42.7	77.4	34.8	41.5	99.4	41.5	38.4	30.5	14.6	42.7	97.6	41.5	40.2	59.1	21.3			
1034	WC-33B	58.5	62.2	67.7	42.7	62.8	76.2	48.2	61.6	84.1	51.8	59.8	98.2	58.5	57.9	49.4	31.1	59.8	90.2	54.9	59.8	62.8	36.0			
1055	WC-Py-7B	34.8	35.4	72.6	25.6	34.1	81.1	28.0	33.5	85.4	30.5	33.5	40.9	13.4	37.2	22.6	6.7	36.0	93.9	34.1	35.4	40.9	11.6			
1000	WC-Py-13B	36.0	40.2	75.6	32.9	37.8	84.8	31.7	37.2	89.0	34.8	43.3	98.2	43.3	43.9	37.2	17.7	43.3	91.5	38.4	39.0	43.3	16.5			
	DS-Coder-6.7B-Ins	65.2	62.2	61.0	40.9	61.0	76.8	47.6	60.4	82.9	51.2	64.0	100.0	64.0	58.5	56.7	36.0	64.0	90.9	59.1	62.8	65.9	39.6			
1056	DS-Coder-7B-Ins	61.0	61.6	57.3	37.2	62.2	71.3	47.0	64.0	84.1	53.7	62.2	99.4	62.2	63.4	66.5	45.1	61.6	95.1	58.5	57.9	64.0	39.6			
	DS-Coder-33B-Ins	65.9	62.8	73.8	47.6	65.2	84.1	53.7	67.1	90.2	59.1	68.9	100.0	68.9	64.0	61.6	39.0	66.5	91.5	60.4	68.3	70.1	48.2			
1057	DS-Coder-V2-16B-Ins	72.0	66.5	77.4	53.7	65.9	84.8	57.9	67.1	89.0	61.6	67.7	98.2	67.7	66.5	43.9	27.4	70.7	88.4	62.2	63.4	64.6	38.4			
1057	DS-V2.5-236B	72.0	72.0	86.0	63.4	72.0	93.3	67.1	74.4	90.2	67.7	73.2	96.3	72.6	75.0	78.7	59.1	70.1	92.7	65.9	67.7	95.7	66.5			
	CodeQwen1.5-7B-Chat	76.2	71.3	47.0	36.6	75.6	61.0	48.2	73.2	74.4	56.1	76.2	98.8	75.0	73.2	62.8	49.4	72.0	93.3	68.3	65.2	70.1	44.5			
1058	Qwen2.5-Coder-7B-Ins	78.0	77.4	61.6	50.6	78.0	70.7	54.3	76.8	64.0	50.6	79.3	99.4	79.3	71.3	61.0	43.3	78.7	92.7	73.8	71.3	78.7	55.5			
	Qwen2-72B-Ins	73.2	72.6	62.8	46.3	76.2	64.6	48.2	75.6	65.9	48.2	73.2	98.8	72.0	68.9	67.1	48.2	73.2	90.2	67.1	70.7	81.1	61.0			
1050	Qwen2.5-72B-Ins	79.3	75.6	92.1	71.3	76.2	96.3	73.8	80.5	87.8	71.3	82.3	98.8	81.7	78.7	79.9	64.0	78.0	92.7	72.0	72.6	97.0	70.7			
1055	Mixtral-8x22B	61.0	61.6	68.3	44.5	63.4	75.0	50.6	62.2	76.2	48.2	65.2	100.0	65.2	64.6	69.5	47.0	56.7	91.5	51.2	57.3	86.6	51.8			
1060	Llama3-8B-Ins	49.4	30.5	40.9	22.0	31.1	51.8	27.4	24.4	45.1	21.3	50.0	99.4	50.0	46.3	60.4	29.9	49.4	93.9	46.3	48.8	73.8	35.4			
1000	Llama3-70B-Ins	65.2	57.3	68.9	43.3	56.7	78.7	49.4	54.3	79.9	50.6	65.2	100.0	65.2	64.0	69.5	43.3	65.2	92.1	61.6	59.8	86.6	52.4			
1001	StarCoder2-15B	36.0	39.6	61.0	23.2	42.1	68.9	29.3	39.0	68.3	31.1	34.8	89.6	34.1	36.0	29.3	9.8	36.0	94.5	34.1	35.4	61.0	15.9			

1061

1047

1062 1063

1064 1065

1066

1067

1069

#### CASE STUDIES D

The tasks provided by our benchmark present challenges that even the most advanced LLMs find difficult to handle. Here, we show some examples with Claude-3.5-Sonnet.

```
1068
      D.1 CORRECTNESS
```

Good Case: The model successfully generates functionally correct code. 1070

1071 # Input Please generate the Python code to solve the following problem. 1072 1073 Problem: 1074 Input to this function is a string containing multiple groups of nested parentheses. Your goal 1075 is to separate those group into separate strings and return the list of those. 1076 Separate groups are balanced (each open brace is properly closed) and not nested within each 1077 other 1078 Ignore any spaces in the input string. 1079

# Output

```
1080
        def separate_parentheses_groups(s):
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
```

```
# Remove all spaces from the input string
s = s.replace(" ", "")
result = []
stack = []
current_group = ""
balance = 0
for char in s:
    if char == '(':
       balance += 1
        stack.append(char)
        current_group += char
    elif char == ')':
       balance -= 1
        stack.pop()
        current_group += char
        # If balance is 0, we've completed a group
        if balance == 0:
            result.append(current_group)
            current_group =
return result
```

1099

#### **Bad Case**: The model fails to generate functionally correct code.

```
1100
         # Input
        Please generate the Python code to solve the following problem.
1101
```

```
1102
        Problem:
1103
        You're given a list of deposit and withdrawal operations on a bank account that starts with
1104
        zero balance. Your task is to detect if at any point the balance of account fallls below zero,
             and
1105
        at that point function should return True. Otherwise it should return False.
```

```
1106
```

```
1107
        # Output
1108
        def check_negative_balance(operations):
            balance = 0
1109
            for operation in operations:
1110
                if operation.startswith("d"): # deposit
                    amount = int(operation[1:])
1111
                    balance += amount
1112
                elif operation.startswith("w"): # withdrawal
                    amount = int (operation[1:])
1113
                    balance -= amount
1114
                    if balance < 0:
                        return True
1115
            return False
```

1116 1117

#### 1118 D.2 READABILITY

# Input

#### **Good Case**: The model successfully generates code with function names employing camel case. 1120

```
1121
```

```
Please generate the Python code to solve the following problem, and use camel case for
1122
             function names.
1123
        Problem:
1124
1125
        Check if in given list of numbers, are any two numbers closer to each other than
        given threshold.
1126
1127
```

```
1128
        # Output
        def areNumbersCloserThanThreshold(numbers, threshold):
1129
1130
            Check if any two numbers in the list are closer to each other than the given threshold.
1131
            :param numbers: List of numbers to check
1132
            :param threshold: The minimum distance allowed between any two numbers
            :return: True if any two numbers are closer than the threshold, False otherwise
1133
            . . . .
            # Sort the list to make comparisons easier
```

```
1135
1136
1137
1138
1139
```

1140 1141

1142

```
sortedNumbers = sorted(numbers)
# Compare adjacent numbers in the sorted list
for i in range(len(sortedNumbers) - 1):
    if abs(sortedNumbers[i] - sortedNumbers[i+1]) < threshold:</pre>
       return True
# If we've gone through all pairs without finding any closer than the threshold
return False
```

#### **Bad Case**: The model fails to generate code with function names employing camel case.

```
1143
         # Input
        Please generate the Python code to solve the following problem, and use camel case for
1144
             function names.
1145
        Problem:
1146
1147
        Find the shortest palindrome that begins with a supplied string.
        Algorithm idea is simple:
1148
         - Find the longest postfix of supplied string that is a palindrome.
1149
        - Append to the end of the string reverse of a string prefix that comes before the palindromic
              suffix.
1150
1151
        Given the following code:
1152
        ```python
1153
        def is_palindrome(string: str) -> bool:
            """ Test if given string is a palindrome """
1154
            return string == string[::-1]
1155
        • • •
1156
```

```
1157
```

1160

1161

1162 1163

1164 1165

1166

1167

1168

1169 1170

1171

1187

return True

```
# Output
        def is_palindrome(string: str) -> bool:
1158
1159
```

```
""" Test if given string is a palindrome """
    return string == string[::-1]
def findShortestPalindrome(s: str) -> str:
    Find the shortest palindrome that begins with the supplied string.
    if not s or is_palindrome(s):
       return s
    for i in range(len(s) - 1, -1, -1):
       if is_palindrome(s[:i+1]):
           return s[i+1:][::-1] + s
    # If no palindromic suffix is found, reverse the entire string and append
    return s[::-1][1:] + s
```

**Good Case**: The model successfully generates code accompanied by line-by-line comments.

```
# Input
1172
        Please generate the Python code to solve the following problem, and add comments for each line
1173
              in each function.
1174
        Problem:
1175
        Check if in given list of numbers, are any two numbers closer to each other than
1176
        given threshold.
1177
1178
         # Output
1179
        def are_numbers_close(numbers, threshold):
            # Sort the list of numbers in ascending order
1180
            sorted_numbers = sorted(numbers)
1181
            # Iterate through the sorted list, comparing adjacent pairs
1182
            for i in range(len(sorted_numbers) - 1):
1183
                # Calculate the difference between adjacent numbers
                difference = sorted_numbers[i+1] - sorted_numbers[i]
1184
1185
                 # Check if the difference is less than the threshold
                if difference < threshold:</pre>
1186
```

# If a pair is found closer than the threshold, return True

# If no pair is found closer than the threshold, return False 1189 1190 **Bad Case**: The model fails to generate code accompanied by line-by-line comments. 1191 # Input 1192 Please generate the Python code to solve the following problem, and add comments for each line 1193 in each function. 1194 Problem: 1195 Input are two strings a and b consisting only of 1s and 0s. 1196 Perform binary XOR on these inputs and return result also as a string. 1197 1198 # Output def binary\_xor(a: str, b: str) -> str: 1199 # Function to perform binary XOR on two input strings 1200 # Ensure both input strings have the same length by padding with zeros 1201 max\_length = max(len(a), len(b)) 1202 a = a.zfill(max\_length) b = b.zfill(max\_length) 1203 1204 # Initialize an empty result string result = " 1205

```
# Initialize an empty result string
result = ""
# Iterate through each position in the input strings
for i in range(max_length):
    # Perform XOR on corresponding bits and add to result
    if a[i] == b[i]:
        result += "0"
    else:
        result += "1"
# Remove leading zeros from the result
result = result.lstrip("0")
# Return "0" if result is empty (all zeros), otherwise return result
```

```
return result if result else "0
```

# Good Case: The model successfully generates code where each line is less than 70 characters long and each function is less than 30 lines long.

```
1218
         # Input
1219
        Please generate the Python code to solve the following problem, where each line is less than
             70 characters long and each function is less than 30 lines long.
1220
1221
        Problem:
1222
        Check if in given list of numbers, are any two numbers closer to each other than
1223
        given threshold.
1224
         # Output
1225
        def are_numbers_closer_than_threshold(numbers, threshold):
1226
            if len(numbers) < 2:</pre>
                return False
1227
1228
            sorted_numbers = sorted(numbers)
1229
            for i in range(len(sorted_numbers) - 1):
1230
                if abs(sorted_numbers[i] - sorted_numbers[i+1]) < threshold:</pre>
                     return True
1231
```

```
return False
```

#### 1232 1233

1188

1206

1207

1208

1209

1210

1211

1212

1213

1214

1215

# Bad Case: The model fails to generate code where each line is less than 70 characters long and each function is less than 30 lines long.

1257 1258

1259 1260

1261

1262

```
# Output
        def sort_even_indices(1):
1244
            # Extract values at even indices
1245
            even_values = [1[i] for i in range(0, len(1), 2)]
1246
            # Sort the even values
1247
            sorted even = sorted(even values)
1248
            # Create a new list with sorted even values and original odd values
1249
            result = []
            even_index = 0
1250
1251
            for i in range(len(l)):
                if i % 2 == 0:
1252
                    result.append(sorted_even[even_index])
1253
                    even index += 1
                else:
1254
                    result.append(l[i])
1255
            return result
1256
```

#### D.3 MAINTAINABILITY

# **Good Case**: The model successfully generates code with only the given function and two addition sub-functions.

# Input 1263 You are given a string word containing distinct lowercase English letters. Telephone keypads have keys mapped with distinct collections of lowercase English letters, 1264 which can be used to form words by pushing them. For example, the key 2 is mapped with [" 1265 a","b","c"], we need to push the key one time to type "a", two times to type "b", and three times to type "c" 1266 It is allowed to remap the keys numbered 2 to 9 to distinct collections of letters. The keys 1267 can be remapped to any amount of letters, but each letter must be mapped to exactly one key. You need to find the minimum number of times the keys will be pushed to type the 1268 string word. 1269 Return the minimum number of pushes needed to type word after remapping the keys. An example mapping of letters to keys on a telephone keypad is given below. Note that 1,  $\star$ , #, 1270 and 0 do not map to any letters. 1271 1272 Example 1: 1273 1274 Input: word = "abcde" 1275 Output: 5 Explanation: The remapped keypad given in the image provides the minimum cost. 1276 "a" -> one push on key 2 1277 "b" -> one push on key 3 "c" -> one push on key 4 1278 "d" -> one push on key 5 1279 "e" -> one push on key 6 Total cost is 1 + 1 + 1 + 1 + 1 = 5. 1280 It can be shown that no other mapping can provide a lower cost. 1281 1282 Example 2: 1283 Input: word = "xycdefghij" 1284 Output: 12 1285 Explanation: The remapped keypad given in the image provides the minimum cost. "x" -> one push on key 2 1286 "y" -> two pushes on key 2 1287 "c" -> one push on key 3 "d" -> two pushes on key 3 1288 "e" -> one push on key 4 1289 "f" -> one push on key 5 "g" -> one push on key 6 1290 "h" -> one push on key 7 1291 "i" -> one push on key 8 "j" -> one push on key 9 1292 Total cost is 1 + 2 + 1 + 2 + 1 + 1 + 1 + 1 + 1 + 1 = 12. 1293 It can be shown that no other mapping can provide a lower cost. 1294 1295 Constraints:

word consists of lowercase English letters.

def minimumPushes(self, word: str) -> int:

All letters in word are distinct.

addition sub-functions:

1 <= word.length <= 26

'''python class Solution:

...

```
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
```

1296

```
# Output
        from collections import Counter
1307
1308
        class Solution:
            def minimumPushes(self, word: str) -> int:
1309
                def count frequencies (word):
1310
                    return Counter (word)
1311
                def calculate_pushes(frequencies):
1312
                    total_pushes = 0
                    for i, freq in enumerate(frequencies):
1313
                        pushes = (i / / 8) + 1
                        total_pushes += pushes * freq
1314
                    return total_pushes
1315
1316
                freq_count = count_frequencies(word)
                sorted_frequencies = sorted(freq_count.values(), reverse=True)
1317
                return calculate_pushes(sorted_frequencies)
```

## **Bad Case**: The model fails to generate code with only the given function and two addition subfunctions.

Please complete the code below to solve above problem, and use only the given function and two

1321

```
# Input
1322
        You are given three positive integers n, x, and y.
1323
        In a city, there exist houses numbered 1 to n connected by n streets. There is a street
             connecting the house numbered i with the house numbered i + 1 for all 1 <= i <= n - 1 .
1324
            An additional street connects the house numbered x with the house numbered y.
1325
        For each k, such that 1 <= k <= n, you need to find the number of pairs of houses (housel,
            house2) such that the minimum number of streets that need to be traveled to reach house2
1326
             from housel is k.
1327
        Return a 1-indexed array result of length n where result[k] represents the total number of
            pairs of houses such that the minimum streets required to reach one house from the other
1328
             is k.
1329
        Note that x and y can be equal.
1330
        Example 1:
1331
1332
        Input: n = 3, x = 1, y = 3
1333
        Output: [6,0,0]
        Explanation: Let's look at each pair of houses:
1334
        - For the pair (1, 2), we can go from house 1 to house 2 directly.
1335
        - For the pair (2, 1), we can go from house 2 to house 1 directly.
1336
        - For the pair (1, 3), we can go from house 1 to house 3 directly.
        - For the pair (3, 1), we can go from house 3 to house 1 directly.
1337
        - For the pair (2, 3), we can go from house 2 to house 3 directly.
        - For the pair (3, 2), we can go from house 3 to house 2 directly.
1338
1339
        Example 2:
1340
1341
        Input: n = 5, x = 2, y = 4
1342
        Output: [10,8,2,0,0]
        Explanation: For each distance k the pairs are:
1343
        - For k == 1, the pairs are (1, 2), (2, 1), (2, 3), (3, 2), (2, 4), (4, 2), (3, 4), (4, 3),
            (4, 5), and (5, 4).
1344
        - For k == 2, the pairs are (1, 3), (3, 1), (1, 4), (4, 1), (2, 5), (5, 2), (3, 5), and (5, 3)
1345
        - For k == 3, the pairs are (1, 5), and (5, 1).
1346
        - For k == 4 and k == 5, there are no pairs.
1347
1348
        Example 3:
1349
        Input: n = 4, x = 1, y = 1
```

```
1350
        Output: [6,4,2,0]
1351
        Explanation: For each distance k the pairs are:
        - For k == 1, the pairs are (1, 2), (2, 1), (2, 3), (3, 2), (3, 4), and (4, 3).
- For k == 2, the pairs are (1, 3), (3, 1), (2, 4), and (4, 2).
1352
1353
         - For k == 3, the pairs are (1, 4), and (4, 1).
1354
        - For k == 4, there are no pairs.
1355
        Constraints:
1356
1357
         2 <= n <= 100
        1 \le x, y \le n
1358
1359
        Please complete the code below to solve above problem, and use only the given function and two
1360
              addition sub-functions:
1361
         ''python
         class Solution:
1362
            def countOfPairs(self, n: int, x: int, y: int) -> List[int]:
1363
         • • •
1364
1365
         # Output
1366
         class Solution:
1367
             def countOfPairs(self, n: int, x: int, y: int) -> List[int]:
1368
                 # Initialize the distance matrix
                 dist = self.initializeDistanceMatrix(n, x, y)
1369
                 # Apply Floyd-Warshall algorithm
1370
                 self.floydWarshall(dist, n)
1371
                 # Count pairs for each distance
1372
                 return self.countPairsForEachDistance(dist, n)
1373
             def initializeDistanceMatrix(self, n: int, x: int, y: int) -> List[List[int]]:
1374
                 # Initialize distance matrix with infinity
1375
                 dist = [[float('inf')] * n for _ in range(n)]
1376
                 # Set distance to self as 0
1377
                 for i in range(n):
                     dist[i][i] = 0
1378
1379
                 # Set distance for adjacent houses
                 for i in range(n - 1):
1380
                     dist[i][i+1] = dist[i+1][i] = 1
1381
                 # Set distance for the additional street
1382
                 if x != y:
1383
                     dist[x-1][y-1] = dist[y-1][x-1] = 1
1384
                 return dist
1385
             def floydWarshall(self, dist: List[List[int]], n: int) -> None:
1386
                 for k in range(n):
1387
                     for i in range(n):
                         for j in range(n):
1388
                             dist[i][j] = min(dist[i][j], dist[i][k] + dist[k][j])
1389
             def countPairsForEachDistance(self, dist: List[List[int]], n: int) -> List[int]:
1390
                 result = [0] * n
1391
                 for i in range(n):
                     for j in range(n):
1392
                         if i != j:
1393
                             result[dist[i][j] - 1] += 1
                 return result
1394
1395
```

#### D.4 EFFICIENCY

1396

1397

1400

**Good Case**: The model successfully generates code with a time complexity of O(n) and a space complexity of O(n). (NI<sub>T</sub> = 100.00, NI<sub>S</sub> = 100.00)

```
# Input
You are given a 0-indexed string s, a string a, a string b, and an integer k.
An index i is beautiful if:
1403 0 <= i <= s.length - a.length
s[i..(i + a.length - 1)] == a
```

```
1404
        There exists an index j such that:
1405
         0 <= j <= s.length - b.length
1406
        s[j..(j + b.length - 1)] == b
1407
        |j - i| <= k
1408
1409
        Return the array that contains beautiful indices in sorted order from smallest to largest.
1410
1411
        Example 1:
1412
         Input: s = "isawsquirrelnearmysquirrelhouseohmy", a = "my", b = "squirrel", k = 15
1413
        Output: [16,33]
        Explanation: There are 2 beautiful indices: [16,33].
1414
         - The index 16 is beautiful as s[16..17] == "my" and there exists an index 4 with s[4..11] ==
1415
              "squirrel" and |16 - 4| <= 15.
         - The index 33 is beautiful as s[33..34] == "my" and there exists an index 18 with s[18..25]
1416
        == "squirrel" and |33 - 18| \le 15.
Thus we return [16,33] as the result.
1417
1418
        Example 2:
1419
         Input: s = "abcd", a = "a", b = "a", k = 4
1420
         Output: [0]
1421
        Explanation: There is 1 beautiful index: [0].
         - The index 0 is beautiful as s[0..0] == "a" and there exists an index 0 with s[0..0] == "a"
1422
             and |0 - 0| \le 4.
1423
        Thus we return [0] as the result.
1424
1425
        Constraints:
1426
        1 <= k <= s.length <= 5 \star 105
1427
        1 <= a.length, b.length <= 5 \star 105
        s, a, and b contain only lowercase English letters.
1428
1429
        Please complete the code below to solve above problem, and make sure that the time complexity
1430
             is 0(n) and the space complexity is 0(n):
        ```python
1431
        class Solution:
1432
            def beautifulIndices(self, s: str, a: str, b: str, k: int) -> List[int]:
1433
         • • •
1434
1435
         # Output
1436
         class Solution:
1437
             def beautifulIndices(self, s: str, a: str, b: str, k: int) -> List[int]:
                 def kmp_search(text, pattern):
1438
                    m = len(pattern)
1439
                     n = len(text)
1440
                     # Compute LPS array
1441
                     lps = [0] * m
                     length = 0
1442
                     i = 1
1443
                     while i < m:</pre>
                         if pattern[i] == pattern[length]:
1444
                             length += 1
1445
                              lps[i] = length
1446
                             i += 1
                         elif length != 0:
1447
                             length = lps[length - 1]
                         else:
1448
                             lps[i] = 0
1449
                             i += 1
1450
                     # KMP search
1451
                     result = []
                     i = j = 0
1452
                     while i < n:
1453
                         if pattern[j] == text[i]:
1454
                             i += 1
                             j += 1
1455
                         if j == m:
1456
                             result.append(i - j)
                         j = lps[j - 1]
elif i < n and pattern[j] != text[i]:</pre>
1457
                             if j != 0:
```

```
1460
1461
1462
1463
1464
```

```
1465
1466
1467
1468
```

1470

1471

1458

1459

```
else
                 i += 1
    return result
# Find all occurrences of 'a' and 'b' in 's'
a_indices = kmp_search(s, a)
b_indices = kmp_search(s, b)
# Find beautiful indices
beautiful_indices = []
j = 0
for i in a_indices:
    while j < len(b_indices) and b_indices[j] < i - k:</pre>
        j += 1
    if j < len(b_indices) and abs(b_indices[j] - i) <= k:</pre>
        beautiful_indices.append(i)
```

j = lps[j - 1]

return beautiful indices

#### **Bad Case**: The model fails to generate code with a time complexity of O(n). (NI<sub>T</sub> = 0)

```
1472
1473
1474
```

```
# Input
        You are given a 0-indexed array of positive integers nums.
1475
        In one operation, you can swap any two adjacent elements if they have the same number of set
1476
             bits. You are allowed to do this operation any number of times (including zero).
        Return true if you can sort the array, else return false.
1477
        Example 1:
1478
1479
        Input: nums = [8, 4, 2, 30, 15]
1480
        Output: true
        Explanation: Let's look at the binary representation of every element. The numbers 2, 4, and 8
1481
              have one set bit each with binary representation "10", "100", and "1000" respectively.
1482
             The numbers 15 and 30 have four set bits each with binary representation "1111" and
             "11110".
1483
        We can sort the array using 4 operations:
        - Swap nums[0] with nums[1]. This operation is valid because 8 and 4 have one set bit each.
1484
             The array becomes [4,8,2,30,15].
1485
        - Swap nums[1] with nums[2]. This operation is valid because 8 and 2 have one set bit each.
1486
             The array becomes [4,2,8,30,15].
        - Swap nums[0] with nums[1]. This operation is valid because 4 and 2 have one set bit each.
1487
             The array becomes [2,4,8,30,15].
1488
        - Swap nums[3] with nums[4]. This operation is valid because 30 and 15 have four set bits each
             . The array becomes [2,4,8,15,30].
1489
        The array has become sorted, hence we return true.
        Note that there may be other sequences of operations which also sort the array.
1490
1491
        Example 2:
1492
        Input: nums = [1, 2, 3, 4, 5]
1493
        Output: true
1494
        Explanation: The array is already sorted, hence we return true.
1495
        Example 3:
1496
        Input: nums = [3, 16, 8, 4, 2]
1497
        Output: false
1498
        Explanation: It can be shown that it is not possible to sort the input array using any number
             of operations.
1499
1500
        Constraints:
1501
1502
        1 <= nums.length <= 100
        1 <= nums[i] <= 28
1503
1504
        Please complete the code below to solve above problem, and make sure that the time complexity
1505
            is $0(n)$:
        ```python
1506
        class Solution:
1507
            def canSortArray(self, nums: List[int]) -> bool:
1508
        • • •
1509
1510
```

```
# Output
1511
         class Solution:
```

```
def canSortArray(self, nums: List[int]) -> bool:
```

```
1512
                  def count_set_bits(n):
1513
                       return bin(n).count('1')
1514
                  n = len(nums)
1515
                  set_bits = [count_set_bits(num) for num in nums]
1516
                   # Group elements with the same number of set bits
1517
                  groups = []
                  current_group = [nums[0]]
current_bits = set_bits[0]
1518
1519
                  for i in range(1, n):
    if set_bits[i] == current_bits:
1520
1521
                           current_group.append(nums[i])
                       else:
1522
                           groups.append(current_group)
1523
                           current_group = [nums[i]]
current_bits = set_bits[i]
1524
1525
                  groups.append(current_group)
1526
                  # Sort each group
1527
                  sorted_groups = [sorted(group) for group in groups]
1528
                  \# Check if the sorted groups are in ascending order <code>prev_max = 0</code>
1529
                  for group in sorted_groups:
1530
                       if group[0] < prev_max:</pre>
1531
                            return False
                       prev_max = group[-1]
1532
1533
                  return True
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
```