

# CODEINSIGHTBENCH: A BENCHMARK FOR ADVANCED CODE UNDERSTANDING AND COMPARISON IN LARGE LANGUAGE MODELS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

While Large Language Models (LLMs) have demonstrated significant progress in coding tasks, their capabilities in nuanced code analysis and deep comprehension remain insufficiently explored. To address this gap, we introduce **CodeInsightBench**, a comprehensive multilingual benchmark designed to systematically evaluate advanced code reasoning. Built upon real-world Codeforces data, it employs both multiple-choice and open-ended questions to assess three core tasks: **Semantic Code Judgment**, **Debugging Path Tracking**, and **Code Efficiency Comparison**. We conduct extensive evaluations on **22** state-of-the-art LLMs (11 closed-source, 11 open-source), revealing critical insights into their strengths and limitations. Our results reveal substantial performance gaps across tasks, with closed-source models generally outperforming open-source counterparts. Besides, models fail primarily on large-scale code transformations, indicating fundamental limitations in understanding code evolution logic. Additionally, the results indicate distinct programming language preferences in code efficiency comparison, and show that multiple sampling substantially improves semantic code judgment performance, with Pass@3 achieving 92.71% accuracy compared to 60.57% at Pass@1. By providing a comprehensive and systematic evaluation methodology, CodeInsightBench enables deeper understanding of LLM capabilities in sophisticated code comprehension tasks.

## 1 INTRODUCTION

Large Language Models (LLMs) have become increasingly influential in programming, offering significant advancements across multiple domains. In particular, they have made notable contributions in areas such as code generation (Nijkamp et al., 2023; Dong et al., 2024), debugging (Wei et al., 2023; Zhong et al., 2024), and code reasoning (Liu et al., 2024a), revolutionizing the way software development tasks are approached. Tools such as GitHub Copilot (GitHub, 2021) and Cursor (Cursor, 2023) have demonstrated impressive capabilities in automating code generation from natural language descriptions, helped developers streamline the coding process and reduce manual effort. Recent advances in LLMs have enabled them to analyze code logic, predict bugs, and offer optimization suggestions. However, these approaches mainly focused on surface-level tasks, with limited attention given to evaluating their deeper understanding of code, such as the ability to reason about code structure, track changes over time, and assess the solution efficiency (Nam et al., 2024).

Existing studies on code reasoning have largely focused on evaluating code generation accuracy (Liu et al., 2023) and efficiency (Dong et al., 2025), while underemphasizing more complex aspects such as code understanding and comparison. As shown in Table 1, popular code evaluation benchmarks suffer from key limitations: focus on isolated code snippets rather than code relationships, reliance on open-ended QA or biased MCQ formats, and use of publicly available datasets prone to data leakage. This gap highlights the need for a systematic benchmark to assess LLMs’ deeper insight into code, moving beyond code generation toward comprehensive evaluation of their comprehension and reasoning capabilities. The ability to answer questions about code represents a natural extension of LLMs’ role in programming, requiring models to exhibit profound understanding of code semantics, debugging processes, and optimization strategies.

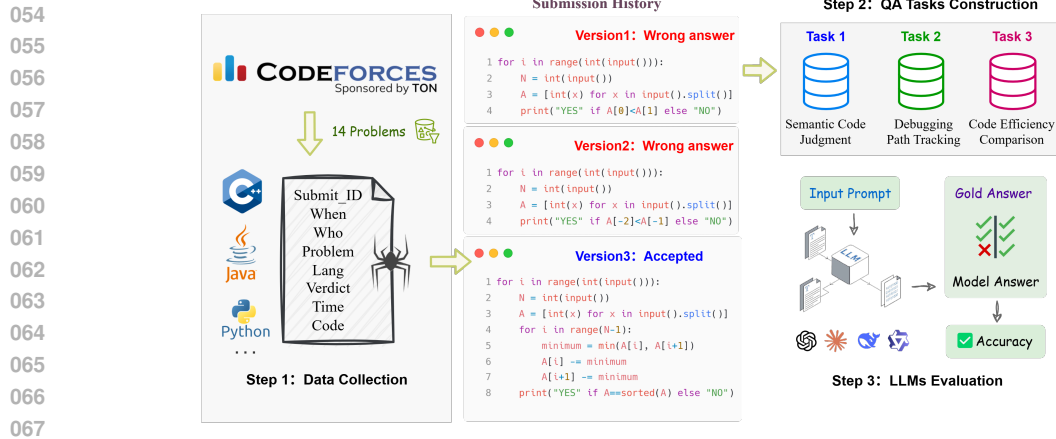


Figure 1: **CodeInsightBench Overview.** The benchmark is constructed by first collecting metadata from Codeforces, followed by the creation of Q&A tasks that focus on key aspects of code understanding, such as semantic judgment, debugging, and efficiency. The final step involves evaluating LLMs by comparing their outputs to a gold standard, enabling the assessment of their accuracy and comprehension of code.

To facilitate comprehensive evaluation of code comprehension and cross-code analysis, we present CodeInsightBench, a novel benchmark that challenges large language models with complex reasoning tasks requiring deep code understanding and comparative analysis. As shown in Figure 1, we systematically construct CodeInsightBench using data from Codeforces (Codeforces, 2025), leveraging actual programming contest submissions to create authentic question-answer pairs derived from verified code implementations. Through rigorous rule-based filtering and manual verification, we curate 14,035 high-quality pairs spanning both open-ended and multiple-choice formats. This methodology ensures meaningful evaluation that reflects genuine coding challenges, moving beyond synthetic annotations toward assessment frameworks that enable more intuitive, human-like programming interactions. Our benchmark evaluates deep code understanding through three strategically designed tasks: **Semantic Code Judgment** for semantic comprehension, **Debugging Path Tracking** for logical execution tracing, and **Code Efficiency Comparison** for algorithmic performance analysis. Unlike traditional benchmarks focused on syntactic correctness, these tasks demand profound understanding of code semantics and programming logic, providing comprehensive assessment of LLMs’ capacity to handle authentic challenges in debugging, optimization, and code reasoning. Our contributions are summarized as follows:

- We introduce CodeInsightBench, a comprehensive multi-language benchmark featuring 14,035 Q&A pairs derived from real-world code submissions. While existing benchmarks primarily assess single-function analysis and code generation capabilities, our benchmark fills a critical gap by evaluating LLMs’ ability to perform cross-code comparative reasoning across three challenging tasks.
- Through extensive evaluation of 22 state-of-the-art models (11 closed-source and 11 open-source), we reveal significant performance gaps in cross-code understanding and uncover intriguing patterns: models excel at surface-level syntax comparison but struggle with deeper code reasoning. Notably, some reasoning models like OpenAI o-series exhibit a confidence miscalibration phenomenon, where high accuracy comes at the cost of trustworthy uncertainty estimation.

## 2 RELATED WORK

**Code Benchmark.** Previous benchmarks, such as HumanEval (Chen et al., 2021), MBPP (Austin et al., 2021), and DS-1000 (Lai et al., 2023), have primarily focused on evaluating code generation capabilities but lacked in-depth assessments of models’ broader code-related abilities. As LLM coding capabilities have expanded, the scope of benchmarks has diversified to include additional tasks, such as code efficiency evaluation (e.g., EFFIBENCH (Huang et al., 2024), Mercury (Du et al., 2024)) and code repair (e.g., MdEval (Liu et al., 2024b), DebugBench (Tian et al., 2024), Feed-

Table 1: The comparison between popular code evaluation benchmarks and CodeInsightBench.

Benchmark	Task Type				Data Properties		
	Code Understanding	Code Comparison	Code Q&A Open-end	Code Q&A MCQ	Multilingual	Real-world	Test Size
EFFIBENCH	✗	✗	✗	✗	✗	✓	1,000
MdEval	✓	✓	✗	✓	✓	✓	3,897
DebugBench	✓	✓	✓	✗	✓	✗	4,253
FeedbackEval	✓	✗	✗	✗	✗	✓	3,736 *
Infibench	✓	✓	✓	✗	✓	✓	234
CRQBench	✓	✓	✓	✗	✗	✓	100
CoCo-Bench	✓	✗	✗	✗	✓	✓	705
CodeXGLUE	✓	✓	✗	✗	✓	✓	2.75M *
CodeCriticBench	✗	✗	✓	✗	✓	✗	4,300
LiveCodeBench	✓	✗	✓	✗	✗	✓	1,432
CodeJudgeBench	✓	✓	✗	✓	✗	✗	4,260
CodeMMLU	✓	✓	✗	✓	✓	✓	19,912 *
<b>CodeInsightBench</b>	✓	✓	✓	✓	✓	✓	<b>14,035</b>

\* Indicates that the data is sourced from publicly available existing datasets.

backEval (Dai et al., 2025)). Simultaneously, the evaluation landscape has evolved to incorporate novel methodological approaches, with several benchmarks adopting multi-task approaches to comprehensively evaluate code comprehension from multiple perspectives. For instance, CoCo-Bench (Yin et al., 2025) and CodeXGLUE (Lu et al., 2021) assess models across various tasks, including code inference, translation, repair, and summarization.

**QA-based Code Evaluation.** Benchmarks like Infibench (Li et al., 2024) and CRQBench (Dinella et al., 2024) focus on open-ended code-related question answering. CodeCriticBench (Zhang et al., 2025) and LiveCodeBench (Jain et al., 2024) utilize open-ended code QA to evaluate models on diverse code-related tasks. These approaches typically source question-answer pairs from open platforms or manual annotations, but the quality can vary significantly, and high annotation costs often result in smaller datasets with limited coverage. Alternatively, multiple-choice benchmarks like CodeMMLU (Manh et al., 2024) and CodeJudgeBench (Jiang et al., 2025) provide more stable evaluation with reduced annotation overhead. CodeMMLU specifically derives its data from existing public datasets and programming resources. However, CodeMMLU’s reliance on existing public datasets and programming resources raises concerns about data leakage, as these sources may have been included in LLM training data. Moreover, MCQ approaches face inherent limitations: they require carefully crafted distractors and may introduce bias when generating candidate answers through automated processes or LLM assistance, potentially limiting their effectiveness in assessing natural reasoning patterns and authentic problem-solving approaches.

### 3 CODEINSIGHTBENCH

In this section, we introduce CodeInsightBench, a novel benchmark designed to assess LLMs’ ability to understand programming contest challenges in depth. The benchmark uses a code question-answering approach with three evaluation tasks, employing both multiple-choice questions and open-ended QA pairs derived from competitive programming solutions. We detail the data collection, task formulation, and evaluation metrics used to assess model performance.

#### 3.1 DATA CURATION

CodeInsightBench utilizes data from Codeforces (Codeforces, 2025), a widely-recognized competitive programming platform, including coding challenges and authentic user submissions. The raw dataset consists of 73,895 submission records for 14 medium-difficulty problems, with ratings ranging from 1000 to 1600. These submissions come from 6,703 unique users and span the period from December 5, 2024, to April 28, 2025. The recency of the dataset ensures that no data leakage into the training data of the evaluated models has occurred. Specific details about the 14 problems can be found in Appendix C.1.

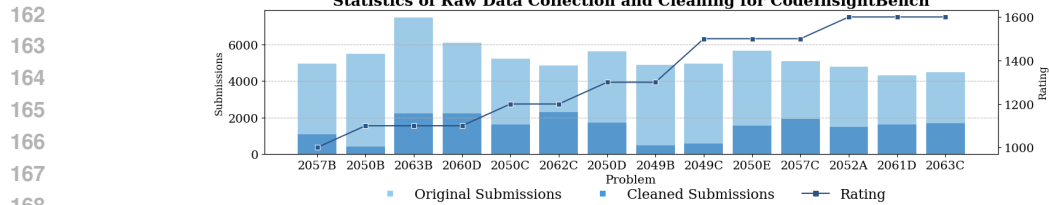


Figure 2: Dataset statistics for CodeInsightBench. The figure shows the problem rating distribution and submission counts in both the original raw dataset and the curated dataset used for question-answer pair construction.

Given the inherent noise and inconsistencies in raw user submission data, comprehensive data cleaning procedures were essential to ensure dataset quality and reliability. We implement systematic filtering criteria to eliminate problematic submissions and maintain data integrity. Specifically, we excluded submissions from users who consecutively submitted solutions in multiple programming languages for the same problem to reduce redundancy and maintain consistency. We also removed submission records that did not result in a final verdict of accepted(AC). To ensure effective code debugging trajectories, we excluded submissions with more than five attempts to focus on meaningful debugging processes while avoiding overly complex cases. Following these rigorous cleaning procedures and manual validation, we retained 20,856 high-quality submissions that form the foundation of our benchmark dataset. Figure 2 presents the statistical overview of the curated dataset utilized for question-answer pair construction. Detailed information about the constructed Q&A datasets for each task is provided in the Appendix C.2.

### 3.2 TASK CONSTRUCTION

As illustrated in Figure 3, CodeInsightBench offers multi-dimensional insights into LLMs’ code understanding and comparison capabilities through three core tasks—Semantic Code Judgment, Debugging Path Tracking, and Code Efficiency Comparison. The prompt template used in these three tasks are detailed in Appendix D.4.1.

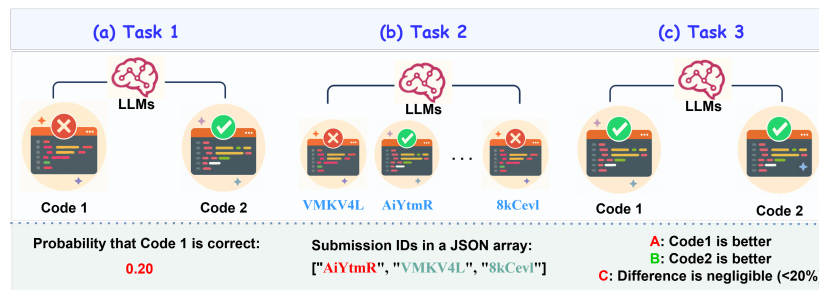


Figure 3: Overview of CodeInsightBench’s three tasks: (a) Semantic Code Judgment; (b) Debugging Path Tracking; (c) Code Efficiency Comparison. Each task presents one or more contest solutions as input and requires the model to output respectively a correctness judgment score, a step-by-step execution trace, or a relative performance decision.

**Task 1: Semantic Code Judgment.** Unlike traditional code correctness evaluation that targets obvious syntactic or logical errors, our Semantic Code Judgment task evaluates models’ ability to distinguish between semantically similar code submissions from the same user, requiring reasoning about how subtle modifications can fundamentally alter program behavior. In this task, the model is given pairs of code submissions where one is correct, and the other contains errors resulting in verdicts such as *Wrong answer*, *Compilation error*, *Runtime error*, *Time limit exceeded*, or *Memory limit exceeded*. The model’s objective is to identify the correct submission. To avoid static binary outputs, we introduce an open-ended probabilistic evaluation approach that requires the model to output the probability that the first code submission receives an *Accept* verdict.

**Task 2: Debugging Path Tracking.** Real-world debugging involves gradual refinement through iterative code modifications, making the understanding of these evolutionary patterns crucial for evaluating models’ code comprehension capabilities. The Debugging Path Tracking task evaluates the model’s ability to understand code evolution patterns by presenting sequences of submissions from the same user for a specific problem, each reflecting iterative attempts to reach an accepted solution (AC). Models must identify the underlying debugging trajectory by recognizing syntactic and semantic variations that indicate logical progression patterns. To prevent position-based bias, all submissions are randomly shuffled and anonymized with modified identifiers. Models are requested to predict submission ID sequences, which are evaluated against ground truth chronological orderings to measure proficiency in code evolution reasoning and comparative analysis.

**Task 3: Code Efficiency Comparison.** While correctness verification remains fundamental in code evaluation, assessing algorithmic efficiency represents a more sophisticated dimension of code understanding that demands deeper analytical capabilities from large language models. The Code Efficiency Comparison task evaluates a model’s ability to compare the relative efficiency of two verified correct code submissions using the same programming language. The objective is to determine which implementation demonstrates superior efficiency based on both theoretical time complexity analysis and empirical execution performance. This is an MCQ task where models choose between different scenarios. To mitigate the impact of minor performance variations inherent in different execution environments and input configurations, we establish an equivalence threshold whereby submissions are considered equally efficient when their execution time difference is less than 20% of their average execution time. This threshold effectively filters measurement noise while preserving meaningful efficiency distinctions.

### 3.3 DATASET ANALYSIS

**Task 1: Semantic Code Judgment.** Task 1 contains 1,400 question-answer pairs across 14 problems, with 100 pairs per problem. To avoid positional bias, correct and incorrect submissions are randomly assigned to code positions, resulting in a nearly balanced distribution. The problems span diverse algorithmic categories including dynamic programming, graph algorithms, and mathematical computations, ensuring comprehensive coverage of programming concepts. Detailed analysis of the character-level edit distance distribution is provided in Appendix C.3, revealing that our dataset captures both subtle semantic differences and more substantial code variations.

**Task 2: Debugging Path Tracking.** Task 2 consists of 3,314 debugging paths from individual users addressing a single problem. The distribution includes 1,869 three-submission paths, 964 four-submission paths, and 481 five-submission paths, reflecting typical debugging cycles in real-world development, where shorter iterations are more common. Each path represents an authentic debugging trace, capturing the iterative refinement process programmers undergo when solving complex problems. The paths encompass various error types including logical errors, edge case handling, and algorithm optimization, providing comprehensive coverage of debugging scenarios.

**Task 3: Code Efficiency Comparison.** Task 3 features 9,321 code pairs across three programming languages: 1,564 in C++17, 1,743 in Java, and 6,014 in Python. The label distribution is balanced, with approximately 35.2% of pairs marked as equally efficient, 32.9% favoring the Code 1, and 31.9% favoring the Code 2. The efficiency comparisons are determined through actual runtime measurements on identical test cases, which serves as a practical proxy for algorithmic efficiency since both codes solve the same problem under identical constraints.

### 3.4 EVALUATION CRITERIA

To evaluate the performance of Large Language Models across the three tasks in CodeInsightBench, we employ four main evaluation metrics: (1) **Accuracy (Acc.)** measures the proportion of correctly predicted instances among all evaluated examples. (2) **Cross-Entropy (CE)** evaluates how well the predicted probability distributions align with the ground-truth labels, with lower values indicating better performance. (3) **Pairwise Accuracy (Pair Acc.)** assesses the model’s ability to maintain correct relative rankings between pairs of items by comparing the predicted order with the ground-truth ranking. (4) **Pass@ $k$**  determines whether the correct answer appears within the top- $k$  predictions, which is particularly valuable for tasks allowing multiple valid outputs. Detailed mathematical definitions and formulations for all evaluation metrics are provided in Appendix C.4.

Table 2: Results of different models. The evaluation covers 11 closed-source and 11 open-source models across 3 tasks in CodeInsightBench.

Type	Model Name	Size	Task 1		Task 2		Task 3
			Acc.↑	CE↓	Acc.↑	Pair Acc.↑	Acc.↑
Close-sourced	GPT-5	-	<b>93.43</b>	<b>0.2125</b>	<b>61.16</b>	<b>83.01</b>	<b>45.50</b>
	O4-mini	-	86.93	1.5695	45.62	73.61	43.73
	Claude-Opus-4-1	-	75.93	0.5873	54.83	79.63	38.65
	O1-mini	-	68.14	3.5528	28.09	62.74	41.81
	Gemini-2.0-Flash-Exp	-	54.14	2.4422	33.43	66.80	35.09
	GPT-4o	-	63.00	0.7632	37.72	68.74	37.61
	GPT-4.1-Mini	-	71.07	1.6420	40.13	70.33	38.86
	GPT-3.5-Turbo	-	52.64	2.9725	15.99	52.06	33.04
	GLM-4-Plus	-	59.86	0.8374	23.78	59.18	36.19
	Claude-3-7-Sonnet	-	67.43	2.5177	44.99	73.60	36.78
	Claude-3-5-Haiku	-	52.57	0.8202	21.54	57.10	35.08
Open-sourced	Qwen3-8B	8B	62.21	3.1812	34.13	66.91	34.33
	Qwen3-14B	14B	<b>62.86</b>	3.1017	33.22	67.00	34.29
	Qwen2.5-Coder-7B-Instruct	7B	51.93	3.8621	16.02	54.49	34.42
	Qwen2.5-7B-Instruct-1M	7B	50.00	<b>1.0060</b>	14.03	51.99	33.27
	Qwen2.5-7B-Instruct	7B	52.64	1.4775	14.21	52.84	34.33
	Qwen2.5-14B-Instruct-1M	14B	60.21	1.0278	24.56	61.51	37.08
	Phi-4	14B	56.86	3.4042	21.12	58.94	33.44
	Llama-3.1-8B-Instruct	8B	52.36	3.5515	11.29	54.03	32.51
	Gemma-2-9B-It	9B	55.36	4.1816	14.24	52.73	33.27
	Deepseek-V3	671B	59.57	1.2582	<b>43.39</b>	<b>72.30</b>	<b>39.20</b>
	DeepSeek-V2	236B	60.36	1.2935	42.97	71.86	38.91

## 4 EXPERIMENTS

### 4.1 EXPERIMENTAL SETUP

We conduct a comprehensive evaluation of 22 popular LLMs from 9 different organizations, including 11 closed-source and 11 open-source models. Detailed parameter settings and model configurations in our experiments are provided in Appendix D.1.

### 4.2 MAIN RESULTS

The results in Table 2 and Figure 4 demonstrate substantial performance differences among models, highlighting the diagnostic value of CodeInsightBench. Closed-source models consistently outperform their open-source counterparts across all evaluation dimensions, with GPT-5 leading by substantial margins, achieving 93.43% accuracy on Task 1 compared to 62.86% for the top open-source model. Beyond that, we find that o1-mini and o4-mini exhibit a confidence miscalibration phenomenon in Task 1, where high accuracy comes at the cost of trustworthy uncertainty estimation. Particularly noteworthy is that o4-mini surpasses Claude-Opus-4-1 in accuracy yet exhibits nearly 3× worse calibration, suggesting that current reasoning model development may prioritize correctness over reliability. These performance patterns are further visualized in Appendix D.2.

**Task Complexity Hierarchy Reveals Fundamental Gaps in Code Understanding.** Task performance analysis reveals a clear difficulty hierarchy across the benchmark. Task 1 proves most accessible to current models, with top performers like GPT-5 achieving 93.43% accuracy, indicating that static code comparison aligns well with LLM capabilities. Task 2 presents moderate difficulty, where even the best model reaches 61.16% accuracy, reflecting the challenges of reasoning through execution traces. The model performance on Task 2 also varies across different debug path lengths, as detailed in Appendix E.1. Most notably, Task 3 emerges as the most challenging task, with all models struggling significantly—the highest accuracy peaks at only 45.50%, demonstrating that algorithmic complexity analysis and runtime behavior understanding remain substantial hurdles for current language models.

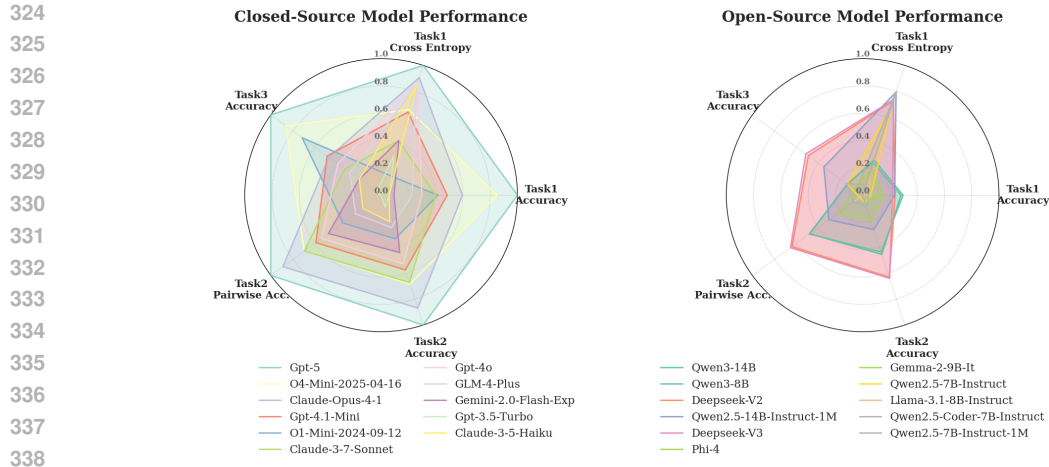


Figure 4: Radar chart comparing the performance of closed-source and open-source models across three tasks in CodeInsightBench. Task 1 evaluates Accuracy and Cross-Entropy, Task 2 assesses Accuracy and Pairwise Accuracy, and Task 3 focuses on Accuracy. The chart highlights that closed-source models generally perform better across tasks, while open-source models demonstrate specialized strengths in certain areas.

**Complex Code Transformations Expose Fundamental Limitations in Semantic Reasoning.** Our analysis reveals that current LLMs systematically fail when reasoning about substantial code evolution, indicating fundamental limitations in semantic understanding beyond syntactic pattern recognition. Table 3 presents collective model failures in the Debugging Path Tracking task, comparing cases where all models fail against those where at least 80% of models answer correctly (similarity calculation detailed in Appendix E.3). These universally challenging scenarios typically involve substantially more complex code structures. By task category, models perform relatively well on Greedy (all wrong 16.47%, 80%+ correct 20.37%) and Sortings (all wrong 6.60%, 80%+ correct 9.35%), but struggle on Bitmasks (all wrong 5.38%, 80%+ correct 3.18%) and Graph matching (all wrong 1.28%, 80%+ correct 0.19%). These findings highlight a critical weakness: models excel at detecting incremental changes but cannot comprehend the broader logical intent underlying longer, more complex, or logic-specific code transformations, indicating a systematic architectural limitation rather than model-specific weaknesses. Developing models capable of learning semantic difference representations between code versions thus represents a key research direction for deeper understanding of how code modifications affect program behavior and logic.

Table 3: Comparative Analysis of Universally Failed vs. High-Agreement Correct Cases in Task 2

Metric	Universally Incorrect Samples	High-Agreement Correct Samples ( $\geq 80\%$ )
Case Count	560	167
Submissions Count	4.10	3.11
Avg. Similarity	0.760	0.823
Min. Similarity	0.619	0.749
Avg. Diff Lines	40.08	25.71
Avg. Token	561.65	362.97

**Poor Algorithmic Reasoning Performance Across All Programming Languages.** To examine language-specific effects on algorithmic reasoning, we evaluate LLM performance across different programming languages in the Code Efficiency Comparison task, with results presented in Figure 5. The evaluation reveals a fundamental challenge: no model achieves satisfactory performance in time complexity assessment, with most accuracy scores falling below 50% regardless of the programming language used. Notably, performance on Python is generally slightly lower than on Java and C++17, and different models exhibit distinct performance patterns across languages. This indicates that LLM performance on algorithmic reasoning and time complexity evaluation depends not only on fundamental reasoning capabilities but also on language syntax, conventions, and training

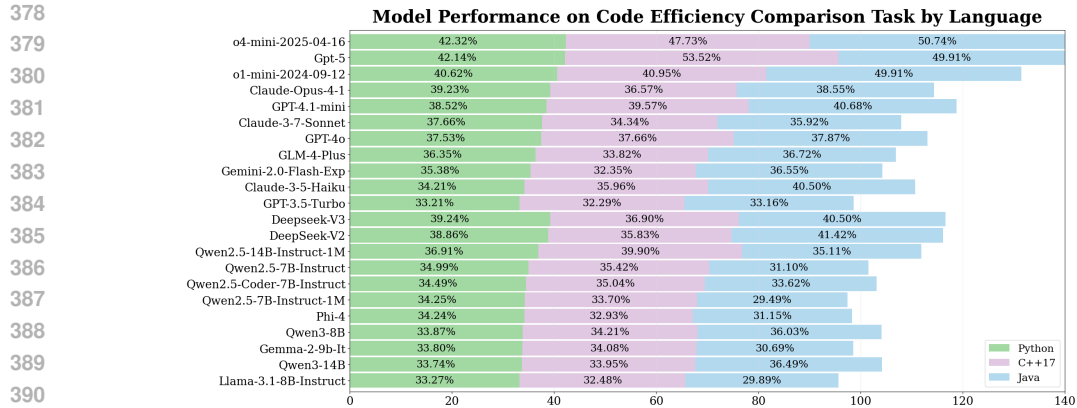


Figure 5: Performance of Models on Time Complexity Assessment Across Different Programming Languages. The stacked bar chart displays the performance of various models on the time complexity assessment task, segmented by programming language (Python, C++17, and Java). Models are compared based on their accuracy across different programming languages, with varying performance trends observed between languages.

data distribution. To further understand model behavior, we conducted comprehensive bias analysis examining overall prediction tendencies, with detailed results presented in Appendix E.2.

### 4.3 MODEL ROBUSTNESS ASSESSMENT

To evaluate whether models can achieve improved performance through multiple sampling attempts, we conduct a Pass@ $k$  analysis using GPT-4-Turbo on Task 1 and Task 3. Task 2 is excluded from this analysis due to its complex multi-step path output format, which requires specialized evaluation metrics beyond simple accuracy measures. We evaluate performance with  $k=1, 2,$  and  $3$  attempts, generating independent samples for each test instance. Table 4 presents the Pass@ $k$  results across different programming languages, revealing how sampling strategies affect model performance in code understanding and code comparing tasks.

For Task 1, GPT-4-Turbo demonstrates substantial improvement through multiple sampling: overall accuracy increases from 0.6057 at Pass@1 to 0.7493 at Pass@2 and reaches 0.9271 at Pass@3. Java shows the strongest performance at Pass@1 (0.6765) and Pass@2 (0.8235), while C++17 achieves the highest Pass@3 accuracy (0.9442). All languages demonstrate significant improvements with multiple sampling attempts, with Pass@3 accuracy exceeding 0.88 across all tested languages.

For Task 3, the baseline performance is considerably lower, with overall Pass@1 accuracy at 0.3756, improving to 0.5134 at Pass@2 and 0.5893 at Pass@3. C++17 consistently outperforms other languages across all  $k$  values, achieving the highest scores of 0.3862, 0.5294, and 0.6279 for Pass@1, Pass@2, and Pass@3 respectively. While multiple sampling provides consistent improvements, the absolute performance remains limited compared to Task 1, highlighting the greater difficulty of time complexity assessment.

Table 4: Pass@ $k$  Evaluation Results of GPT-4-Turbo. The table shows Pass@ $k$  accuracy across different programming languages on Task 1 and Task 3.

Language	Task 1			Task 3		
	Pass@1	Pass@2	Pass@3	Pass@1	Pass@2	Pass@3
Java	<b>0.6765</b>	<b>0.8235</b>	0.9118	0.3534	0.5141	0.5904
Python	0.6000	0.7286	0.8857	0.3793	0.5090	0.5790
C++17	0.6265	0.7679	<b>0.9442</b>	<b>0.3862</b>	<b>0.5294</b>	<b>0.6279</b>
C++20	0.5886	0.7275	0.9346	-	-	-
C++23	0.5862	0.7385	0.8994	-	-	-
<b>Overall</b>	<b>0.6057</b>	<b>0.7493</b>	<b>0.9271</b>	<b>0.3756</b>	<b>0.5134</b>	<b>0.5893</b>

Table 5: Task 1 Performance Under Different Prompting Strategies for Claude-3-5-Haiku and Qwen2.5-7B-Instruct-1M

multirow2*Prompt Type	Claude-3-5-Haiku		Qwen2.5-7B-Instruct-1M	
	Accuracy ↑	Cross-Entropy ↓	Accuracy ↑	Cross-Entropy ↓
Default	52.57	0.8202	50.00	1.0060
Tree-of-Thought	51.64	1.0035	51.29	1.7278
Chain-of-Draft	54.36	0.6848	51.57	4.2554
Few-shot	55.36	1.2382	51.79	4.7582
Chain-of-Thought	<b>55.50</b>	<b>1.5145</b>	<b>52.00</b>	<b>4.7800</b>

#### 4.4 PROMPTING STRATEGY IMPACT

Based on previous results, Claude-3-5-Haiku exhibited the lowest accuracy among closed-source models, while Qwen2.5-7B-Instruct-1M showed the poorest performance among open-source models. To investigate potential improvements for Task 1, we systematically evaluated four advanced prompting strategies, including Few-shot, Chain-of-Thought (Wei et al., 2022), Tree-of-Thought (Yao et al., 2023), and Chain-of-Draft (Xu et al., 2025). The detailed prompts are provided in Appendices D.4.2, D.4.3, D.4.4, and D.4.5.

As shown in Table 5, the impact of advanced prompting strategies varied between models and across different approaches. For Qwen2.5-7B-Instruct-1M, all advanced prompting strategies consistently improved accuracy over the default baseline, with Chain-of-Thought achieving the highest improvement to 52.00%. However, Claude-3-5-Haiku exhibited mixed results, with accuracy of Tree-of-Thought decreasing to 51.64%, while the other three strategies showed improvements, with Chain-of-Thought reaching the peak accuracy of 55.50%.

Among the strategies that improved accuracy, Chain-of-Draft emerged as the most balanced approach for Claude-3-5-Haiku, achieving accuracy of 54.36% while maintaining the lowest cross-entropy of 0.6848 among advanced strategies. Although Chain-of-Thought maximizes accuracy across both models, it significantly compromises confidence calibration, producing the highest cross-entropy values despite accuracy gains. For deployment scenarios prioritizing reliable uncertainty quantification, our findings suggest adopting Chain-of-Draft for Claude-3-5-Haiku and Tree-of-Thought for Qwen2.5-7B-Instruct-1M, as these strategies achieve the most favorable accuracy-calibration trade-offs within their respective model families.

## 5 CONCLUSION

We introduce CodeInsightBench, a novel benchmark designed to evaluate deeper code understanding in Large Language Models (LLMs). CodeInsightBench assesses models across three critical tasks including Semantic Code Judgment, Debugging Path Tracking, and Code Efficiency Comparison. Our evaluation of 22 state-of-the-art models (11 open-source and 11 closed-source) reveals a significant performance gap. While the best-performing model GPT-5 achieves strong accuracy (93.43%) in distinguishing correct code versions, performance drops substantially to 45.50% on efficiency analysis tasks requiring deeper algorithmic reasoning. These results provide quantitative evidence that current LLMs excel at surface-level syntactic tasks but struggle with semantic code comprehension and execution flow analysis. The consistent failure patterns across models indicate **fundamental limitations in transitioning from pattern recognition to deep code reasoning**. Our findings establish that advancing LLM code understanding requires moving beyond syntactic accuracy toward robust semantic reasoning and algorithmic analysis capabilities.

**Statement.** All datasets used in this work are publicly available and were obtained through legitimate channels, ensuring no ethical concerns regarding data access or privacy. Large language models were used in this work primarily for grammar correction and language polishing, with details provided in Appendix A. To ensure reproducibility, all code and data are made available as described in Appendix B.

## REFERENCES

- 486  
487  
488 Marah Abdin, Jyoti Aneja, Harkirat Behl, Sébastien Bubeck, Ronen Eldan, Suriya Gunasekar,  
489 Michael Harrison, Russell J Hewett, Mojan Javaheripi, Piero Kauffmann, et al. Phi-4 technical  
490 report. *arXiv preprint arXiv:2412.08905*, 2024.
- 491 Anthropic. Claude models. <https://www.anthropic.com>, 2025.
- 492  
493 Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan,  
494 Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language  
495 models. *arXiv preprint arXiv:2108.07732*, 2021.
- 496  
497 BigModel. Glm-4-plus. <https://bigmodel.cn/>, 2025.
- 498  
499 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared  
500 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large  
501 language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- 502  
503 Codeforces. Codeforces, 2025. URL <https://codeforces.com>. Accessed: 2025-07-18.
- 504  
505 Cursor. Cursor: Ai-powered code completion. <https://www.cursor.so/>, 2023.
- 506  
507 Dekun Dai, MingWei Liu, Anji Li, Jialun Cao, Yanlin Wang, Chong Wang, Xin Peng, and Zibin  
508 Zheng. FeedbackEval: A benchmark for evaluating large language models in feedback-driven  
509 code repair tasks. *arXiv preprint arXiv:2504.06939*, 2025.
- 510  
511 Google DeepMind. Gemini models. <https://deepmind.com>, 2025.
- 512  
513 DeepSeek-AI. Deepseek-v3 technical report, 2024a. URL <https://arxiv.org/abs/2412.19437>.
- 514  
515 DeepSeek-AI. Deepseek-v2: A strong, economical, and efficient mixture-of-experts language  
516 model, 2024b.
- 517  
518 Elizabeth Dinella, Satish Chandra, and Petros Maniatis. Crqbench: A benchmark of code reasoning  
519 questions. *arXiv preprint arXiv:2408.08453*, 2024.
- 520  
521 Yihong Dong, Xue Jiang, Zhi Jin, and Ge Li. Self-collaboration code generation via chatgpt. *ACM  
522 Transactions on Software Engineering and Methodology*, 33(7):1–38, 2024.
- 523  
524 Yihong Dong, Jiazheng Ding, Xue Jiang, Ge Li, Zhuo Li, and Zhi Jin. Codescore: Evaluating  
525 code generation by learning code execution. *ACM Transactions on Software Engineering and  
526 Methodology*, 34(3):1–22, 2025.
- 527  
528 Mingzhe Du, Anh Tuan Luu, Bin Ji, Qian Liu, and See-Kiong Ng. Mercury: A code efficiency  
529 benchmark for code large language models. *arXiv preprint arXiv:2402.07844*, 2024.
- 530  
531 GitHub. Github copilot. <https://copilot.github.com/>, 2021.
- 532  
533 Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad  
534 Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd  
535 of models. *arXiv preprint arXiv:2407.21783*, 2024.
- 536  
537 Dong Huang, Yuhao Qing, Weiyi Shang, Heming Cui, and Jie Zhang. Effibench: Benchmarking the  
538 efficiency of automatically generated code. *Advances in Neural Information Processing Systems*,  
539 37:11506–11544, 2024.
- Binyuan Hui, Jian Yang, Zeyu Cui, Jiayi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Kai Dang, An Yang, Rui Men, Fei Huang, Shanghaoran Quan, Xingzhang Ren, Xuancheng Ren, Jingren Zhou, and Junyang Lin. Qwen2.5-coder technical report. *ArXiv*, abs/2409.12186, 2024. URL <https://api.semanticscholar.org/CorpusID:272707390>.
- Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free evaluation of large language models for code. *arXiv preprint arXiv:2403.07974*, 2024.

- 540 Hongchao Jiang, Yiming Chen, Yushi Cao, Hung-yi Lee, and Robby T Tan. Codejudgebench:  
541 Benchmarking llm-as-a-judge for coding tasks. *arXiv preprint arXiv:2507.10535*, 2025.
- 542
- 543 Yuhang Lai, Chengxi Li, Yiming Wang, Tianyi Zhang, Ruiqi Zhong, Luke Zettlemoyer, Wen-tau  
544 Yih, Daniel Fried, Sida Wang, and Tao Yu. Ds-1000: A natural and reliable benchmark for data  
545 science code generation. In *International Conference on Machine Learning*, pp. 18319–18345.  
546 PMLR, 2023.
- 547 Linyi Li, Shijie Geng, Zhenwen Li, Yibo He, Hao Yu, Ziyue Hua, Guanghan Ning, Siwei Wang,  
548 Tao Xie, and Hongxia Yang. Infibench: Evaluating the question-answering capabilities of code  
549 large language models. *Advances in Neural Information Processing Systems*, 37:128668–128698,  
550 2024.
- 551
- 552 Changshu Liu, Shizhuo Dylan Zhang, Ali Reza Ibrahimzada, and Reyhaneh Jabbarvand. Code-  
553 mind: A framework to challenge large language models for code reasoning. *arXiv preprint*  
554 *arXiv:2402.09664*, 2024a.
- 555 Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. Is your code generated by chat-  
556 gpt really correct? rigorous evaluation of large language models for code generation. *Advances*  
557 *in Neural Information Processing Systems*, 36:21558–21572, 2023.
- 558
- 559 Shukai Liu, Linzheng Chai, Jian Yang, Jiajun Shi, He Zhu, Liran Wang, Ke Jin, Wei Zhang,  
560 Hualei Zhu, Shuyue Guo, et al. Mdeval: Massively multilingual code debugging. *arXiv preprint*  
561 *arXiv:2411.02310*, 2024b.
- 562 Shuai Lu, Daya Guo, Shuo Ren, Junjie Huang, Alexey Svyatkovskiy, Ambrosio Blanco, Colin  
563 Clement, Dawn Drain, Daxin Jiang, Duyu Tang, et al. Codexglue: A machine learning benchmark  
564 dataset for code understanding and generation. *arXiv preprint arXiv:2102.04664*, 2021.
- 565
- 566 Dung Nguyen Manh, Thang Phan Chau, Nam Le Hai, Thong T Doan, Nam V Nguyen, Quang  
567 Pham, and Nghi DQ Bui. Codemmlu: A multi-task benchmark for assessing code understanding  
568 capabilities of codellms. *arXiv preprint arXiv:2410.01999*, 2024.
- 569
- 570 Daye Nam, Andrew Macvean, Vincent Hellendoorn, Bogdan Vasilescu, and Brad Myers. Using  
571 an llm to help with code understanding. In *Proceedings of the IEEE/ACM 46th International*  
572 *Conference on Software Engineering*, pp. 1–13, 2024.
- 573 Erik Nijkamp, Hiroaki Hayashi, Caiming Xiong, Silvio Savarese, and Yingbo Zhou. Code-  
574 gen2: Lessons for training llms on programming and natural languages. *arXiv preprint*  
575 *arXiv:2305.02309*, 2023.
- 576
- 577 OpenAI. Openai’s gpt series and variants. <https://openai.com>, 2025.
- 578
- 579 Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bh-  
580 patiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, et al. Gemma  
2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*, 2024.
- 581
- 582 Runchu Tian, Yining Ye, Yujia Qin, Xin Cong, Yankai Lin, Yinxu Pan, Yesai Wu, Haotian Hui, We-  
583 ichuan Liu, Zhiyuan Liu, et al. Debugbench: Evaluating debugging capability of large language  
584 models. *arXiv preprint arXiv:2401.04621*, 2024.
- 585
- 586 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny  
587 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in*  
*neural information processing systems*, 35:24824–24837, 2022.
- 588
- 589 Yuxiang Wei, Chunqiu Steven Xia, and Lingming Zhang. Copiloting the copilots: Fusing large  
590 language models with completion engines for automated program repair. In *Proceedings of the*  
591 *31st ACM Joint European Software Engineering Conference and Symposium on the Foundations*  
592 *of Software Engineering*, pp. 172–184, 2023.
- 593
- Silei Xu, Wenhao Xie, Lingxiao Zhao, and Pengcheng He. Chain of draft: Thinking faster by writing  
less. *arXiv preprint arXiv:2502.18600*, 2025.

594 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu,  
595 Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint*  
596 *arXiv:2505.09388*, 2025a.  
597  
598 An Yang, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoyan Huang, Jiandong Jiang,  
599 Jianhong Tu, Jianwei Zhang, Jingren Zhou, et al. Qwen2. 5-1m technical report. *arXiv preprint*  
600 *arXiv:2501.15383*, 2025b.  
601  
602 Qwen An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan  
603 Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu,  
604 Jianwei Zhang, Jianxin Yang, Jiabin Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu,  
605 Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji  
606 Lin, Tianhao Li, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yi-Chao  
607 Zhang, Yunyang Wan, Yuqi Liu, Zeyu Cui, Zhenru Zhang, Zihan Qiu, Shanghaoran Quan, and  
608 Zekun Wang. Qwen2.5 technical report. *ArXiv*, abs/2412.15115, 2024. URL <https://api.semanticscholar.org/CorpusID:274859421>.  
609  
610 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik  
611 Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *Ad-*  
612 *vances in neural information processing systems*, 36:11809–11822, 2023.  
613  
614 Wenjing Yin, Tianze Sun, Yijiong Yu, Jiawei Fang, Guangyao Su, Jiancheng Wang, Zekun Wang,  
615 Wei Wang, Ran Chen, Ziyun Dai, et al. Coco-bench: A comprehensive code benchmark for  
616 multi-task large language model evaluation. *arXiv preprint arXiv:2504.20673*, 2025.  
617  
618 Alexander Zhang, Marcus Dong, Jiaheng Liu, Wei Zhang, Yejie Wang, Jian Yang, Ge Zhang, Tianyu  
619 Liu, Zhongyuan Peng, Yingshui Tan, et al. Codecriticbench: A holistic code critique benchmark  
620 for large language models. *arXiv preprint arXiv:2502.16614*, 2025.  
621  
622 Li Zhong, Zilong Wang, and Jingbo Shang. Debug like a human: A large language model debugger  
623 via verifying runtime execution step-by-step. *arXiv preprint arXiv:2402.16906*, 2024.  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647

648	APPENDIX CONTENTS	
649		
650	<b>A LLM Usage</b>	<b>14</b>
651		
652	<b>B Reproducibility Statement</b>	<b>14</b>
653		
654		
655	<b>C CodeInsightBench Analysis</b>	<b>14</b>
656	C.1 Descriptions of Selected Codeforces Problems . . . . .	14
657	C.2 Overview of the CodeInsightBench Dataset . . . . .	20
658	C.3 Character-Level Edit Distance Analysis in Task 1 . . . . .	21
659	C.4 Evaluation Metrics . . . . .	21
660		
661		
662		
663	<b>D Experiment Setup</b>	<b>22</b>
664	D.1 Evaluated Model Configuration . . . . .	22
665	D.2 Visualization of Model Performance in Closed-Source vs. Open-Source . . . . .	22
666	D.3 Visualization of Accuracy Comparison Across Languages for Different Tasks . . . . .	22
667	D.4 Prompt Templates . . . . .	24
668		
669		
670		
671	<b>E Result Analysis</b>	<b>27</b>
672	E.1 Impact of Debug Path Length on Debugging Path Tracking Performance . . . . .	27
673	E.2 Answer Distribution by Gold Answer Label in Task 3 . . . . .	28
674	E.3 Code Similarity Calculation Methodology . . . . .	28
675		
676		
677		
678		
679		
680		
681		
682		
683		
684		
685		
686		
687		
688		
689		
690		
691		
692		
693		
694		
695		
696		
697		
698		
699		
700		
701		

## A LLM USAGE

Large language models (LLMs) were used in this work solely for grammar correction and language polishing to improve the clarity and readability of the manuscript. Specifically, LLMs were employed to (1) correct grammatical errors and improve sentence structure, (2) enhance word choice and phrasing for better clarity, and (3) ensure consistent writing style throughout the paper.

LLMs were not used for research ideation, experimental design, data analysis, result interpretation, or generation of scientific content. All research contributions, including benchmark construction, methodological innovations, experimental results, and scientific conclusions, are entirely the work of the human authors. No LLM-generated content was directly incorporated into the core research findings or claims.

The authors take full responsibility for all content in this paper and have verified the accuracy of all claims and statements. All LLM-assisted text has been thoroughly reviewed and validated by the human authors to ensure factual correctness and appropriateness.

## B REPRODUCIBILITY STATEMENT

We have released our data and code on an anonymous website <https://anonymous.4open.science/r/CodeInsightBench-C2C5/>. The repository includes the complete CodeInsightBench dataset, baseline model implementations, evaluation scripts, and experimental results. All materials are provided with detailed documentation to ensure reproducibility and facilitate future research in code comprehension evaluation.

## C CODEINSIGHTBENCH ANALYSIS

### C.1 DESCRIPTIONS OF SELECTED CODEFORCES PROBLEMS

In the context of competitive programming, various problem categories require distinct approaches and strategies. The data presented in Table 6 provides information about 14 selected problems from Codeforces. These problems represent a broad spectrum of algorithmic challenges, encompassing various categories such as Brute Force, Constructive Algorithms, Graph Matching, Implementation, Greedy Algorithms, Mathematics, and Dynamic Programming (DP). Each problem in the table is classified into one or more of these categories based on the type of approach or algorithm required to solve it. This diverse selection highlights the comprehensiveness of the dataset, capturing a wide array of problem types and algorithmic techniques that are central to competitive programming.

Table 6: Details of the 14 Selected Codeforces Problems

Question ID	Question Description	Problem Tags
2049B - pspspsp	<p>Cats are attracted to ‘pspsps’, but Evirir, being a dignified dragon, is only attracted to ‘pspsps’ with oddly specific requirements. Given a string <math>s = s_1s_2\dots s_n</math> of length <math>n</math> consisting of characters ‘p’, ‘s’, and ‘.’ (dot), determine whether a permutation* <math>p</math> of length <math>n</math> exists, such that for all integers <math>i</math> (<math>1 \leq i \leq n</math>):</p> <ul style="list-style-type: none"> <li>- If <math>s_i</math> is ‘p’, then <math>[p_1, p_2, \dots, p_i]</math> forms a permutation (of length <math>i</math>);</li> <li>- If <math>s_i</math> is ‘s’, then <math>[p_i, p_{i+1}, \dots, p_n]</math> forms a permutation (of length <math>n - i + 1</math>);</li> <li>- If <math>s_i</math> is ‘.’, then there is no additional restriction.</li> </ul>	Brute force Constructive algorithms Graph matching Implementation

Continued on next page

Table 6 – continued from previous page

Question ID	Question Description	Problem Tags
2049C - MEX Cycle	<p>Evirir the dragon has many friends. They have 3 friends! That is one more than the average dragon.</p> <p>You are given integers <math>n</math>, <math>x</math>, and <math>y</math>. There are <math>n</math> dragons sitting in a circle. The dragons are numbered <math>1, 2, \dots, n</math>. For each <math>i</math> (<math>1 \leq i \leq n</math>), dragon <math>i</math> is friends with dragon <math>i - 1</math> and <math>i + 1</math>, where dragon <math>0</math> is defined to be dragon <math>n</math> and dragon <math>n + 1</math> is defined to be dragon <math>1</math>. Additionally, dragons <math>x</math> and <math>y</math> are friends with each other (if they are already friends, this changes nothing). Note that all friendships are mutual.</p> <p>Output <math>n</math> non-negative integers <math>a_1, a_2, \dots, a_n</math> such that for each dragon <math>i</math> (<math>1 \leq i \leq n</math>), the following holds: Let <math>f_1, f_2, \dots, f_k</math> be the friends of dragon <math>i</math>. Then:</p> $a_i = \text{mex}(a_{f_1}, a_{f_2}, \dots, a_{f_k})$	Brute force Constructive algorithms Greedy Implementation
2050B - Transfusion	<p>You are given an array <math>a</math> of length <math>n</math>. In one operation, you can pick an index <math>i</math> from 2 to <math>n - 1</math> inclusive, and do one of the following actions:</p> <ol style="list-style-type: none"> <li>Decrease <math>a_{i-1}</math> by 1, then increase <math>a_{i+1}</math> by 1.</li> <li>Decrease <math>a_{i+1}</math> by 1, then increase <math>a_{i-1}</math> by 1.</li> </ol> <p>After each operation, all the values must be non-negative. Can you make all the elements equal after any number of operations?</p>	Brute force Greedy Math
2050C - Uninteresting Number	<p>You are given a number <math>n</math> with a length of no more than <math>10^5</math>.</p> <p>You can perform the following operation any number of times: choose one of its digits, square it, and replace the original digit with the result. The result must be a digit (that is, if you choose the digit <math>x</math>, then the value of <math>x^2</math> must be less than 10).</p> <p>Is it possible to obtain a number that is divisible by 9 through these operations?</p>	Brute force DP Math

Continued on next page

756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809

Table 6 – continued from previous page

Question ID	Question Description	Problem Tags
2050D - Digital string maximization	<p>You are given a string <math>s</math>, consisting of digits from 0 to 9. In one operation, you can pick any digit in this string, except for 0 or the leftmost digit, decrease it by 1, and then swap it with the digit to the left of the picked digit.</p> <p>For example, in one operation from the string '1023', you can get '1103' or '1022'. Find the lexicographically maximum string you can obtain after any number of operations.</p>	Brute force Greedy Math Strings
2050E - Three Strings	<p>You are given three strings: <math>a</math>, <math>b</math>, and <math>c</math>, consisting of lowercase Latin letters. The string <math>c</math> was obtained in the following way:</p> <ul style="list-style-type: none"> <li>- At each step, either string <math>a</math> or string <math>b</math> was randomly chosen, and the first character of the chosen string was removed from it and appended to the end of string <math>c</math>, until one of the strings ran out. After that, the remaining characters of the non-empty string were added to the end of <math>c</math>.</li> <li>- Then, a certain number of characters in string <math>c</math> were randomly changed.</li> </ul> <p>For example, from the strings <math>a = abra</math> and <math>b = cada</math>, without character replacements, the strings <math>caabdraa</math>, <math>abracada</math>, <math>acadabra</math> could be obtained.</p> <p>Find the minimum number of characters that could have been changed in string <math>c</math>.</p>	DP Implementation Strings

Continued on next page

810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863

Table 6 – continued from previous page

Question ID	Question Description	Problem Tags
2052A - Adrenaline Rush	<p>Alice’s friend is a big fan of the Adrenaline Rush racing competition and always strives to attend every race. However, this time, Alice is the one watching the race. To ensure her friend does not miss any important details, Alice decides to take notes on everything that happens on the track.</p> <p>The first thing Alice notices before the race begins is the numbering of the cars. All the cars line up in front of the starting line in a specific order. The car closest to the line is numbered 1, the second car is numbered 2, and so on, up to the last car, which is numbered <math>n</math>. How convenient! — Alice thought.</p> <p>The race begins with the countdown: "Three! Two! One! Go!". Alice observes that the cars start in their original order. However, as the race progresses, their order changes. She records whenever one car overtakes another, essentially swapping places with it on the track.</p> <p>During the race, Alice notices something curious: no car overtakes another more than once. In other words, for any two cars <math>x</math> and <math>y</math>, there are at most two overtakes between them during the race: "<math>x</math> overtakes <math>y</math>" and/or "<math>y</math> overtakes <math>x</math>".</p> <p>At the end of the race, Alice carefully writes down the final order of the cars <math>c_1, c_2, \dots, c_n</math>, where <math>c_1</math> represents the winner of the race.</p> <p>Alice’s friend, however, is only interested in the final ranking and discards all of Alice’s notes except for the final ordering. As Alice is quite curious, she wonders: What is the longest possible sequence of overtakes she could have observed during the race? Your task is to help Alice answer this question.</p>	Constructive algorithms

Continued on next page

864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917

Table 6 – continued from previous page

Question ID	Question Description	Problem Tags
2057B - Gorilla and the Exam	<p>Due to a shortage of teachers in the senior class of the "T-generation", it was decided to have a huge male gorilla conduct exams for the students. However, it is not that simple; to prove his competence, he needs to solve the following problem.</p> <p>For an array <math>b</math>, we define the function <math>f(b)</math> as the smallest number of the following operations required to make the array <math>b</math> empty:</p> <ol style="list-style-type: none"> <li>1. Take two integers <math>l</math> and <math>r</math>, such that <math>l \leq r</math>, and let <math>x</math> be the <math>\min(b_l, b_{l+1}, \dots, b_r)</math>; 2. Remove all such <math>b_i</math> that <math>l \leq i \leq r</math> and <math>b_i = x</math> from the array. The deleted elements are removed, and the indices are renumbered.</li> </ol> <p>You are given an array <math>a</math> of length <math>n</math> and an integer <math>k</math>. No more than <math>k</math> times, you can choose any index <math>i</math> (<math>1 \leq i \leq n</math>) and any integer <math>p</math>, and replace <math>a_i</math> with <math>p</math>.</p> <p>Help the gorilla to determine the smallest value of <math>f(a)</math> that can be achieved after such replacements.</p>	Greedy Sortings
2057C - Trip to the Olympiad	<p>In the upcoming year, there will be many team olympiads, so the teachers of "T-generation" need to assemble a team of three pupils to participate in them. Any three pupils will show a worthy result in any team olympiad. But winning the olympiad is only half the battle; first, you need to get there...</p> <p>Each pupil has an independence level, expressed as an integer. In "T-generation", there is exactly one student with each independence level from <math>l</math> to <math>r</math>, inclusive. For a team of three pupils with independence levels <math>a</math>, <math>b</math>, and <math>c</math>, the value of their team independence is equal to:</p> $(a \oplus b) + (b \oplus c) + (a \oplus c)$ <p>where <math>\oplus</math> denotes the bitwise XOR operation.</p> <p>Your task is to choose any trio of students with the maximum possible team independence.</p>	Bitmasks Constructive algorithms Greedy Math
2060D - Subtract Min Sort	<p>You are given a sequence <math>a</math> consisting of <math>n</math> positive integers.</p> <p>You can perform the following operation any number of times:</p> <ul style="list-style-type: none"> <li>- Select an index <math>i</math> (<math>1 \leq i &lt; n</math>), and subtract <math>\min(a_i, a_{i+1})</math> from both <math>a_i</math> and <math>a_{i+1}</math>.</li> </ul> <p>Determine if it is possible to make the sequence non-decreasing by using the operation any number of times.</p>	Greedy

Continued on next page

Table 6 – continued from previous page

Question ID	Question Description	Problem Tags
2061D - Kevin and Numbers	<p>Kevin wrote an integer sequence <math>a</math> of length <math>n</math> on the blackboard.</p> <p>Kevin can perform the following operation any number of times:</p> <ul style="list-style-type: none"> <li>- Select two integers <math>x, y</math> on the blackboard such that <math> x - y  \leq 1</math>, erase them, and then write down an integer <math>x + y</math> instead.</li> </ul> <p>Kevin wants to know if it is possible to transform these integers into an integer sequence <math>b</math> of length <math>m</math> through some sequence of operations.</p> <p>Two sequences <math>a</math> and <math>b</math> are considered the same if and only if their multisets are identical. In other words, for any number <math>x</math>, the number of times it appears in <math>a</math> must be equal to the number of times it appears in <math>b</math>.</p>	<p>Bitmasks Data structures</p>
2062C - Cirno and Operations	<p>Cirno has a sequence <math>a</math> of length <math>n</math>. She can perform either of the following two operations for any (possibly, zero) times unless the current length of <math>a</math> is 1:</p> <ol style="list-style-type: none"> <li>1. <b>**Reverse the sequence.**</b> Formally, <math>[a_1, a_2, \dots, a_n]</math> becomes <math>[a_n, a_{n-1}, \dots, a_1]</math> after the operation.</li> <li>2. <b>**Replace the sequence with its difference sequence.**</b> Formally, <math>[a_1, a_2, \dots, a_n]</math> becomes <math>[a_2 - a_1, a_3 - a_2, \dots, a_n - a_{n-1}]</math> after the operation.</li> </ol> <p>Find the maximum possible sum of elements of <math>a</math> after all operations.</p>	<p>Brute force Math</p>

Continued on next page

972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

Table 6 – continued from previous page

Question ID	Question Description	Problem Tags
2063B - Subsequence Update	<p>After Little John borrowed expansion screws from auntie a few hundred times, eventually she decided to come and take back the unused ones. But as they are a crucial part of home design, Little John decides to hide some in the most unreachable places — under the eco-friendly wood veneers.</p> <p>You are given an integer sequence <math>a_1, a_2, \dots, a_n</math>, and a segment <math>[l, r]</math> (<math>1 \leq l \leq r \leq n</math>).</p> <p>You must perform the following operation on the sequence exactly once:</p> <ul style="list-style-type: none"> <li>- Choose any <b>subsequence</b> of the sequence <math>a</math>, and reverse it. Note that the subsequence does not have to be contiguous. Formally, choose any number of indices <math>i_1, i_2, \dots, i_k</math> such that <math>1 \leq i_1 &lt; i_2 &lt; \dots &lt; i_k \leq n</math>. Then, change the <math>i_x</math>-th element to the original value of the <math>i_{k-x+1}</math>-th element simultaneously for all <math>1 \leq x \leq k</math>.</li> </ul> <p>Find the minimum value of <math>a_l + a_{l+1} + \dots + a_{r-1} + a_r</math> after performing the operation.</p> <p>* A sequence <math>b</math> is a subsequence of a sequence <math>a</math> if <math>b</math> can be obtained from <math>a</math> by the deletion of several (possibly, zero or all) elements from arbitrary positions.</p>	Constructive algorithms Data structures Greedy Sortings
2063C - Remove Exactly Two	<p>Recently, Little John got a tree from his aunt to decorate his house. But as it seems, just one tree is not enough to decorate the entire house. Little John has an idea. Maybe he can remove a few vertices from the tree. That will turn it into more trees! Right?</p> <p>You are given a tree of <math>n</math> vertices. You must perform the following operation exactly twice:</p> <ol style="list-style-type: none"> <li>1. Select a vertex <math>v</math>;</li> <li>2. Remove all edges incident to <math>v</math>, and also the vertex <math>v</math>.</li> </ol> <p>Please find the maximum number of connected components after performing the operation exactly twice.</p> <p>Two vertices <math>x</math> and <math>y</math> are in the same connected component if and only if there exists a path from <math>x</math> to <math>y</math>. For clarity, note that the graph with 0 vertices has 0 connected components by definition.</p>	Brute force Data structures DFS and similar DP Graphs Greedy Sortings Trees

## C.2 OVERVIEW OF THE CODEINSIGHTBENCH DATASET

Table 7 presents a comprehensive overview of the CodeInsightBench dataset, which comprises three distinct tasks with varying complexity and scope. The dataset contains 14,035 question-answer pairs in total, distributed across Task 1 (Semantic Code Judgment, 1,400 pairs), Task 2 (Debugging Path Tracking, 3,314 pairs), and Task 3 (Code Efficiency Comparison, 9,321 pairs). Each sample is

constructed from real-world code submissions and designed to support long-context reasoning capabilities, with input lengths spanning from 700 to over 44,000 tokens. This substantial variation in context length, combined with the diverse task formulations, enables comprehensive evaluation of language models beyond mere syntactic understanding, encompassing deep semantic comprehension, logical inference capabilities, and performance-aware code analysis.

Table 7: Statistical Overview of Tasks in CodeInsightBench Dataset

Task	Task Name	Q&A Pairs	Users	Min Tokens	Max Tokens	Languages
Task 1	Semantic Code Judgment	1400	1250	729	29728	C++17, Python, C++20, Java, C++23, Other
Task 2	Debugging Path Tracking	3314	3093	718	44089	C++17, Python, C++20, Java, C++23, other
Task 3	Code Efficiency Comparison	9321	1928	739	16360	C++17, Java, Python

### C.3 CHARACTER-LEVEL EDIT DISTANCE ANALYSIS IN TASK 1

Figure 6 visualizes the character-level edit distance distribution, revealing that over 60% of submission pairs differ by fewer than 200 characters. This distribution underscores the task’s emphasis on fine-grained semantic reasoning rather than surface-level textual differences. The minimal character variations between correct and incorrect submissions create a challenging evaluation scenario where success depends on deep code comprehension rather than pattern matching or superficial heuristics.

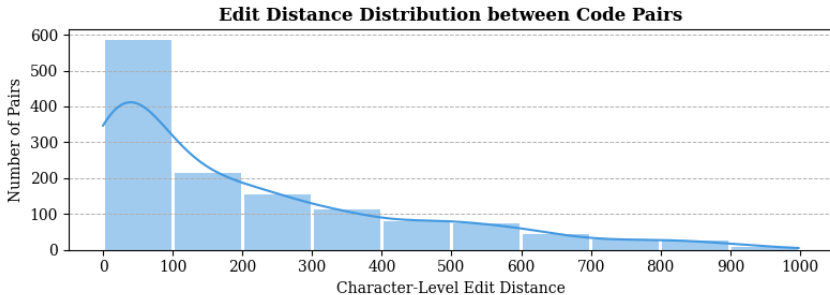


Figure 6: Edit Distance Distribution between Code Pairs

### C.4 EVALUATION METRICS

To evaluate the performance of Large Language Models across the three tasks in CodeInsightBench, we employ four main metrics: Accuracy, Cross-Entropy, Pairwise Accuracy, and Pass@k.

**Accuracy(Acc.)** measures the proportion of correctly predicted instances among all evaluated examples. If the model prediction matches the ground-truth label, it is considered correct. Higher accuracy indicates better overall predictive performance. It is defined as:

$$\text{Accuracy} = \frac{1}{N} \sum_{i=1}^N \mathbb{I}(\hat{y}_i = y_i) \tag{1}$$

where  $\hat{y}_i$  is the prediction,  $y_i$  is the ground truth, and  $\mathbb{I}(\cdot)$  is the indicator function that returns 1 when the prediction is correct.

**Cross-Entropy(CE)** measures the difference between the predicted probability distribution and the true distribution. It is commonly used in classification tasks to evaluate how well the predicted

probabilities align with the actual labels. Lower values indicate better alignment. It is defined as:

$$\text{Cross-Entropy} = -\frac{1}{N} \sum_{i=1}^N [y_i \log(\hat{p}_i) + (1 - y_i) \log(1 - \hat{p}_i)] \quad (2)$$

where  $\hat{p}_i$  is the predicted probability for class 1, and  $y_i \in \{0, 1\}$  is the ground-truth label.

**Pairwise-Accuracy (Pair Acc.)** evaluates the model’s ability to maintain the correct relative ranking between pairs of items. It considers all item pairs in the predicted order and checks whether their order matches the ground-truth ranking. The score is calculated as:

$$\text{Pairwise-Accuracy} = \frac{1}{T} \sum_{i < j} \mathbb{I}(\text{GoldRank}[\text{pred}_i] < \text{GoldRank}[\text{pred}_j]) \quad (3)$$

where  $T = \binom{N}{2}$  is the total number of distinct item pairs,  $\text{pred}_i$  is the  $i$ -th item in the model-predicted ranking, and  $\text{GoldRank}[\cdot]$  denotes the position of an item in the ground-truth order.

$\text{Pass}@k$  evaluates whether the correct answer appears within the top- $k$  predictions of the model. This metric is particularly useful in code generation or retrieval tasks where multiple outputs can be considered. Higher  $\text{Pass}@k$  indicates better performance under multiple-choice or top- $k$  generation settings.

$$\text{Pass}@k = \frac{1}{N} \sum_{i=1}^N \mathbb{I}[y_i \in \text{TopK}(x_i, k)] \quad (4)$$

where  $N$  is the total number of samples,  $x_i$  is the input for the  $i$ -th sample,  $y_i$  is the ground-truth answer,  $\text{TopK}(x_i, k)$  denotes the top- $k$  predictions, and  $\mathbb{I}[\cdot]$  is the indicator function that returns 1 if the condition is true, and 0 otherwise.

## D EXPERIMENT SETUP

### D.1 EVALUATED MODEL CONFIGURATION

Table 8 details the configuration of the 22 models. All model evaluations on CodeInsightBench were conducted using a consistent set of inference parameters. The temperature was uniformly set to 0.7, and the maximum number of output tokens was set to 4096. Inference for closed-source models and DeepSeek models (DeepSeek-V2, DeepSeek-V3) was performed via API calls, while other open-source models were evaluated using local deployments. All inference processes were conducted on NVIDIA A800-SXM4-80GB GPUs. For the detailed  $\text{Pass}@k$  performance analysis in the Semantic Code Judgment task and the Code Efficiency Comparison task, OpenAI’s GPT-4-Turbo was used to generate 3 independent samples per instance to evaluate  $\text{Pass}@k$  metrics.

### D.2 VISUALIZATION OF MODEL PERFORMANCE IN CLOSED-SOURCE VS. OPEN-SOURCE

Figure 7 presents a comparative bar chart illustrating the accuracy of various Large Language Models (LLMs) on the three distinct tasks within the CodeInsightBench: Semantic Code Judgment, Debugging Path Tracking, and Code Efficiency Comparison. The left panel displays the performance of closed-source models, while the right panel shows the performance of open-source models.

### D.3 VISUALIZATION OF ACCURACY COMPARISON ACROSS LANGUAGES FOR DIFFERENT TASKS

Figure 8 presents a bar chart comparing the accuracy across different programming languages for three distinct tasks. Relative performance among languages varies by task; for instance, Python and C++17 show strong results in Semantic Code Judgment, while Java demonstrates competitive accuracy in Debugging Path Tracking and Code Efficiency Comparison, highlighting task-specific strengths and weaknesses.

Table 8: Model List with Citations and Versions

Model	Citation	Version
<b>Close-Sourced Models</b>		
Gemini-2.0-Flash-Exp	DeepMind (2025)	Gemini-2.0-Flash-Exp
GPT-3.5-Turbo	OpenAI (2025)	gpt-3.5-turbo-1106
GPT-4-Turbo	OpenAI (2025)	gpt-4-turbo
GPT-4o	OpenAI (2025)	gpt-4o-2024-08-06
o1-mini	OpenAI (2025)	o1-mini-2024-09-12
o4-mini	OpenAI (2025)	o4-mini-2025-04-16
GPT-4.1-Mini	OpenAI (2025)	gpt-4.1-mini-2025-04-14
GPT-5	OpenAI (2025)	gpt-5-2025-08-07
Claude-3-5-Haiku	Anthropic (2025)	Claude-3-5-Haiku-20241022
Claude-3-7-Sonnet	Anthropic (2025)	Claude-3-7-Sonnet-20250219
Claude Opus 4.1	Anthropic (2025)	Claude-Opus-4-1-20250805
GLM-4-Plus	BigModel (2025)	glm-4-plus
<b>Open-Sourced Models</b>		
Phi-4	Abdin et al. (2024)	<a href="#">microsoft/phi-4</a>
Llama-3.1-8B-Instruct	Grattafiori et al. (2024)	<a href="#">meta-llama/Llama-3.1-8B-Instruct</a>
Gemma-2-9B-It	Team et al. (2024)	<a href="#">google/gemma-2-9b-it</a>
Deepseek-V3	DeepSeek-AI (2024a)	<a href="#">deepseek-ai/DeepSeek-V3</a>
DeepSeek-V2	DeepSeek-AI (2024b)	<a href="#">deepseek-ai/DeepSeek-V2-Chat</a>
Qwen2.5-Coder-7B-Instruct	Hui et al. (2024)	<a href="#">Qwen/Qwen2.5-Coder-7B-Instruct</a>
Qwen2.5-7B-Instruct	Yang et al. (2024)	<a href="#">Qwen/Qwen2.5-7B-Instruct</a>
Qwen2.5-7B-Instruct-1M	Yang et al. (2025b)	<a href="#">Qwen/Qwen2.5-7B-Instruct-1M</a>
Qwen2.5-14B-Instruct-1M	Yang et al. (2025b)	<a href="#">Qwen/Qwen2.5-14B-Instruct-1M</a>
Qwen3-8B	Yang et al. (2025a)	<a href="#">Qwen/Qwen3-8B</a>
Qwen3-14B	Yang et al. (2025a)	<a href="#">Qwen/Qwen3-14B</a>

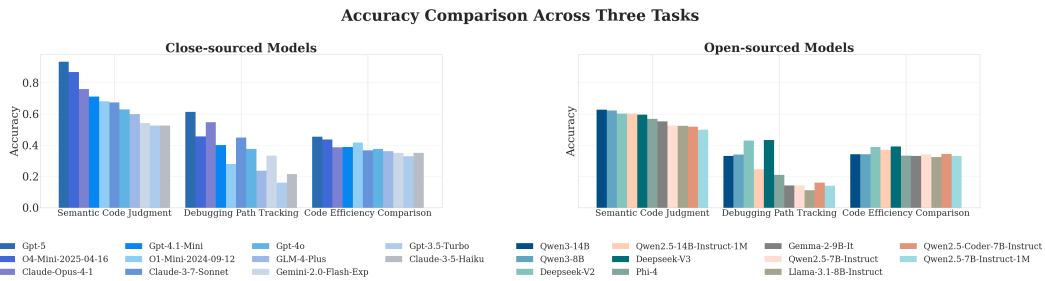


Figure 7: Bar Chart of Model Performance in Closed-Source vs. Open-Source

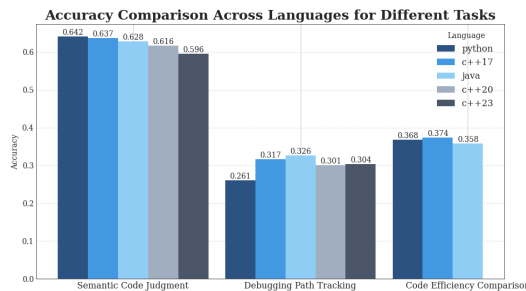


Figure 8: Bar Chart of Accuracy Comparison Across Languages for Different Tasks

1242 D.4 PROMPT TEMPLATES  
1243

1244 Each task in CodeInsightBench has a default prompt template designed to guide the model’s re-  
1245 sponse, with varying levels of detail and instructions, along with specific output format require-  
1246 ments. To further analyze the impact of different prompt strategies on model performance, we  
1247 additionally employed Few-shot, Chain-of-Thought (CoT), Tree-of-Thought (ToT), and Chain-of-  
1248 Draft (CoD) prompting strategies for evaluation on Task 1. The following subsections provide the  
1249 detailed prompt information used in the evaluation.

1250  
1251 D.4.1 DEFAULT PROMPTS FOR THREE TASK TYPES  
12521253 Prompt used for Task 1  
1254

1255 Below is a programming problem, followed by two code submissions. One of them is  
1256 correct. Please determine the probability that Code1 is correct. The probability should be  
1257 a value between 0 and 1, rounded to two decimal places. Do not provide any additional  
1258 explanation or details.

1259 Problem: {question}  
1260 Example: {example}

1261  
1262 Code1:{code1}  
1263 Code2:{code2}

1264  
1265 Probability that Code1 is correct:  
1266

1267 Prompt used for Task 2  
1268

1269 You are a code review expert tasked with analyzing a series of code submissions to  
1270 determine their logical order. You are given a problem and out-of-order code submission  
1271 history. Each submission is marked with a unique ID in the format [ID: xxxx].  
1272

1273 Your task is carefully analyze the code under each ID to determine the most likely order  
1274 they were written in, from the earliest (typically more incorrect) to the latest (typically more  
1275 correct or refined). Consider factors such as syntax errors, logic issues, completeness, and  
1276 quality improvements. Provide a JSON array listing the IDs in the inferred logical order.  
1277 Only output the ID strings in order, like this: ["abc123", "xyz789"]

1278 Problem: question  
1279 Submissions:{submission\_list}

1280  
1281 Output:  
1282

1283 Prompt used for Task 3  
1284

1285 Compare two correct solutions to the same problem, written in the same language. Focus  
1286 on their runtime performance based on both the provided code and execution times.  
1287

1288 Problem: {question}  
1289 Example: {example}

1290  
1291 Code1:{code1}  
1292 Code2:{code2}

1293  
1294 Which code is more efficient?  
1295

1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349

The answer is: Only return one number: 1 if Code1 is faster, 2 if Code2 is faster, or 0 if the difference is negligible ( $\leq 20\%$ ). Do not explain or output anything else.

#### D.4.2 FEW-SHOT PROMPT FOR TASK1

##### Few-shot prompt used for Task 1

Below is a programming problem, followed by two code submissions. One of them is correct, and the other may have a verdict of 'Wrong answer', 'Compilation error', 'Runtime error', 'Time limit exceeded', or 'Memory limit exceeded'. Please think carefully and predict the probability that Code1 is correct. The probability should be a value between 0 and 1, rounded to two decimal places. A value of 0.5 indicates uncertainty, so try to avoid outputting 0.5. If the probability is less than 0.5, it means Code2 is more likely to be correct. Do not provide any additional explanation or details.

Problem: Check if a number is prime.

Example 1:

Code1:

```
def is_prime(n):
    if n <= 1:
        return False
    for i in range(2, int(n ** 0.5) + 1):
        if n % i == 0:
            return False
    return True
```

Code2:

```
def is_prime(n):
    if n <= 1:
        return True # Incorrect
    for i in range(2, n):
        if n % i == 0:
            return False
    return True
```

Probability that Code1 is correct: 1.00

Problem: Compute the sum of elements in a list.

Example 2:

Code1:

```
def sum_list(arr):
    total = 0
    for x in arr:
        total += x
    return total
```

Code2:

```
def sum_list(arr):
    total = 0
    i = 0
    while i < len(arr):
        total += arr[i]
        # missing i += 1 causes an infinite loop
    return total
```

Probability that Code1 is correct: 1.00

Problem: {question}

Example: {example}

1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1370  
1371  
1372  
1373  
1374  
1375  
1376  
1377  
1378  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
1402  
1403

Code1: {code1}  
Code2: {code2}

Probability that Code1 is correct:

#### D.4.3 CHAIN-OF-THOUGHT PROMPT FOR TASK 1

Chain-of-Thought prompt used for Task 1

Below is a programming problem, followed by two code submissions. One of them is correct, and the other may have a verdict of 'Wrong answer', 'Compilation error', 'Runtime error', 'Time limit exceeded', or 'Memory limit exceeded'. Please think carefully and predict the probability that Code1 is correct. The probability should be a value between 0 and 1, rounded to two decimal places. A value of 0.5 indicates uncertainty, so try to avoid outputting 0.5. If the probability is less than 0.5, it means Code2 is more likely to be correct. Do not provide any additional explanation or details. Break down your reasoning into clear steps to make your decision, and evaluate the logic of the code in both submissions.

Problem: {question}  
Example: {example}

Code1:{code1}  
Code2:{code2}

Thinking step 1: Evaluate the core logic of Code1 and Code2. What does each code do to solve the problem? Is the logic correct?

Thinking step 2: Identify if there are any errors in Code1 and Code2, such as logical or boundary condition errors. Does either code return incorrect results for specific edge cases?

Thinking step 3: Based on the previous steps, decide which code is more likely to be correct. Take into account the logic and correctness of both submissions.

Probability that Code1 is correct:

#### D.4.4 TREE-OF-THOUGHT PROMPT FOR TASK 1

Tree-of-Thought prompt used for Task 1

Below is a programming problem, followed by two code submissions. Imagine three experts analyzing the codes through structured reasoning paths:

- Path 1 (Algorithm Logic): Does the code implement the correct approach for the problem? Vote Code1 or Code2.

- Path 2 (Edge Case Handling): Which code better handles boundary conditions (empty inputs, extreme values, etc.)? Vote accordingly.

- Path 3 (Syntax and Runtime Safety): Does the code have compilation errors, runtime crashes, or infinite loops? Vote the safer one.

- Path 4 (Efficiency): Does the code avoid time/memory limit issues (for applicable problems)? Vote the more efficient one.

Each valid vote for Code1 adds 0.25 to its score; each for Code2 adds 0.25. Total score is normalized to a probability (0.00-1.00), rounded to two decimals. Avoid 0.5 by leaning into reasoning biases. Do not provide any additional explanation or details.

Problem: {question}  
Example: {example}

```
Code1:{code1}
Code2:{code2}
```

Probability that Code1 is correct:

#### D.4.5 CHAIN-OF-DRAFT PROMPT FOR TASK 1

Chain-of-Draft prompt used for Task 1

You are given a programming problem and two code submissions. One is correct, the other has a bug. Think in three micro-drafts. Each draft must be less than 5 words. At the end, output the probability that Code1 is correct (rounded to two decimals). Output only the three drafts and the final probability. No explanations.

```
Problem: {question}
Example: {example}
```

```
Code1:{code1}
Code2:{code2}
```

Probability that Code1 is correct:

## E RESULT ANALYSIS

### E.1 IMPACT OF DEBUG PATH LENGTH ON DEBUGGING PATH TRACKING PERFORMANCE

Figure 9 displays model accuracy on the Debugging Path Tracking task, segmented by submission sequence length (Version Count: 3, 4, or 5). A clear trend is the superior performance of several closed-source models, notably o4-mini-2025-04-16 and Claude-3-7-Sonnet, which consistently outperform most open-source counterparts, especially on shorter sequences (Version Count 3, green bars). Large open-source models like Deepseek-V3 and DeepSeek-V2 also demonstrate robust capabilities across all sequence lengths. Generally, accuracy tends to decrease as the submission sequence lengthens from 3 to 5 versions, indicating that tracing longer, more complex debugging paths is more challenging for all models. Many smaller open-source models show markedly lower performance, particularly struggling with longer sequences, highlighting the difficulty of multi-step logical reasoning for these LLMs.

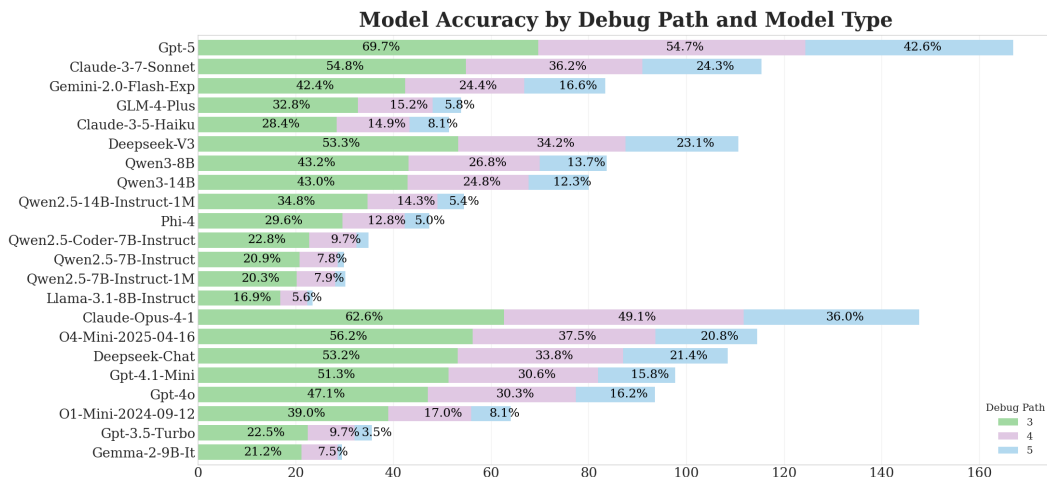


Figure 9: Model Accuracy by Debug Path and Model Type

## E.2 ANSWER DISTRIBUTION BY GOLD ANSWER LABEL IN TASK 3

The Table 9 presents the prediction distribution of various models across different gold answer labels (Gold=0, 1, 2) in Task 3, along with their overall accuracy and standard deviation (STD). It provides a comparative view of both closed-source and open-source models, revealing notable label bias in many models. Specifically, a strong tendency to predict Gold=1 or Gold=2 is observed. The standard deviation reflects the imbalance in predictions across labels, indicating the degree of prediction skewness or bias.

Although the ground-truth labels in the Code Efficiency Comparison task are evenly distributed among the three classes, we observe clear prediction imbalances across many models. For instance, Qwen3-14B predicts label “2” in over 95% of cases, with minimal correct predictions for labels “0” and “1”, resulting in a high standard deviation of 0.4231. Similar skewed patterns are observed in models such as Qwen3-8B (std = 0.4297) and Qwen2.5-7B-Instruct-1M (std = 0.3971), suggesting a strong positional bias. In contrast, more capable models like GPT-4o (std = 0.0976) and Qwen2.5-14B (std = 0.1270) exhibit more balanced output distributions. These discrepancies are not caused by the data distribution itself, but rather reflect the models’ tendencies to favor specific answer positions—possibly due to format sensitivity or insufficient reasoning ability. This analysis highlights that even when task design is balanced, certain models can still exhibit systematic biases, which in turn limits their effectiveness in nuanced comparative reasoning tasks.

Table 9: Answer Distribution by Gold Answer Label in Task 3

Type	Model Name	Accuracy	Gold=0	Gold=1	Gold=2	STD
Closed-source	GPT-5	0.4550	0.1541	<b>0.6283</b>	0.6087	0.2191
	GPT-4o	0.3761	0.2716	<b>0.5075</b>	0.3563	0.0976
	Claude-Opus-4-1	0.3865	0.0271	<b>0.5101</b>	0.6561	0.2688
	Claude-3-5-Haiku	0.3508	0.1347	<b>0.5093</b>	0.4291	0.1611
	o1-mini-2024-09-12	0.4181	0.1256	0.5387	<b>0.6475</b>	0.2248
	GPT-3.5-Turbo	0.3304	0.0341	0.4053	<b>0.5805</b>	0.2278
	o4-mini-2025-04-16	0.4373	0.1094	<b>0.6624</b>	0.5726	0.2423
	GPT-4.1-mini	0.3886	0.0442	0.5363	<b>0.6263</b>	0.2559
	Claude-3-7-Sonnet	0.3678	0.0037	0.4814	<b>0.6528</b>	0.2747
	Gemini-2.0-Flash-Exp	0.3509	0.0198	0.3664	<b>0.7005</b>	0.2779
GLM-4-Plus	0.3619	0.0042	0.4447	<b>0.6849</b>	0.2819	
Open-source	Qwen2.5-14B-Instruct-1M	0.3691	0.3042	<b>0.5439</b>	0.2522	0.1270
	Qwen2.5-Coder-7B-Instruct	0.3424	<b>0.4793</b>	0.4243	0.1048	0.1651
	Qwen2.5-7B-Instruct	0.3470	0.1832	0.2371	<b>0.6440</b>	0.2057
	Phi-4	0.3344	0.1185	<b>0.6270</b>	0.2716	0.2130
	Gemma-2-9B-It	0.3327	0.0149	<b>0.6897</b>	0.3160	0.2760
	DeepSeek-V2	0.3891	0.0385	<b>0.7457</b>	0.4090	0.2888
	DeepSeek-V3	0.3920	0.0388	<b>0.7576</b>	0.4052	0.2935
	Llama-3.1-8B-Instruct	0.3251	0.0000	<b>0.7897</b>	0.2057	0.3345
	Qwen2.5-7B-Instruct-1M	0.3425	0.0014	<b>0.8966</b>	0.1195	0.3971
	Qwen3-14B	0.3429	0.0030	0.1159	<b>0.9516</b>	0.4231
Qwen3-8B	0.3387	0.0039	0.1009	<b>0.9600</b>	0.4297	

## E.3 CODE SIMILARITY CALCULATION METHODOLOGY

To quantify the degree of code modification between different submissions in the debugging path tracking analysis, we employ standardized similarity metrics using Python’s `difflib` library.

Code similarity is computed using `difflib.SequenceMatcher`, which implements the Ratcliff-Obershelp algorithm. The similarity score is calculated as:

$$\text{Similarity} = \frac{2 \times \text{Number of Matching Characters}}{\text{Length of Code}_1 + \text{Length of Code}_2}$$

This metric ranges from 0 to 1, where 1 indicates identical code and 0 represents completely different code. The algorithm identifies the longest common subsequences between two code snippets and

1512 computes the ratio based on the total matching characters relative to the combined length of both  
1513 code segments.

1514 To measure the extent of code modifications, we calculate the number of altered lines using  
1515 `difflib.unified_diff`. This function generates a unified diff format that identifies added,  
1516 deleted, and modified lines between two code versions. The difference line count represents the  
1517 total number of lines that differ between the original and modified code submissions.

1518 These metrics provide complementary perspectives on code evolution that similarity scores capture  
1519 overall structural preservation, while difference line counts quantify the granular extent of modifi-  
1520 cations. Together, they enable systematic analysis of the relationship between code transformation  
1521 complexity and model performance degradation.

1522  
1523  
1524  
1525  
1526  
1527  
1528  
1529  
1530  
1531  
1532  
1533  
1534  
1535  
1536  
1537  
1538  
1539  
1540  
1541  
1542  
1543  
1544  
1545  
1546  
1547  
1548  
1549  
1550  
1551  
1552  
1553  
1554  
1555  
1556  
1557  
1558  
1559  
1560  
1561  
1562  
1563  
1564  
1565