

000 001 002 003 004 005 006 007 008 009 010 011 012 013 014 015 016 017 018 019 020 021 022 023 024 025 026 027 028 029 030 031 032 033 034 035 036 037 038 039 040 041 042 043 044 045 046 047 048 049 050 051 052 053 AUTOCODEBENCH: LARGE LANGUAGE MODELS ARE AUTOMATIC CODE BENCHMARK GENERATORS

Anonymous authors

Paper under double-blind review

ABSTRACT

Large Language Models (LLMs) have shown impressive performance across diverse domains, with code generation emerging as a particularly prominent application. However, existing benchmarks designed to evaluate code generation exhibit several critical limitations. First, most rely on manual annotations, which are time-consuming and difficult to scale across programming languages and problem complexities. Second, the majority focus primarily on Python, while the few multilingual benchmarks suffer from limited difficulty and imbalanced language coverage. To overcome these challenges, we present **AutoCodeGen**, an automated framework for constructing high-difficulty, multilingual code generation datasets without manual annotations. Our approach guarantees correctness and completeness of test cases by generating test inputs with LLMs, obtaining test outputs within a multilingual sandbox, and further enhancing quality through reverse problem generation and multi-stage filtering. Based on this novel method, we introduce **AutoCodeBench**, a large-scale benchmark suite spanning 20 programming languages with balanced coverage. AutoCodeBench is designed to rigorously evaluate LLMs on diverse, challenging, and realistic multilingual programming tasks. Extensive experiments reveal that even state-of-the-art models struggle on these tasks, particularly in low-resource languages. Besides, we release complementary training and evaluation resources, including a large-scale, verifiable multilingual training set generated via the same pipeline, as well as a multilingual sandbox with high-concurrency support. We hope these contributions will provide a solid foundation for future research and inspire the community to explore more automatic and scalable approaches to multilingual code generation.¹

1 INTRODUCTION

Recently, Large Language Models (LLMs) have advanced rapidly, achieving strong performance across a wide range of tasks (OpenAI, 2024; Gemini, 2025; DeepSeek-AI, 2025b; Anthropic, 2025b). Among these, code generation has emerged as a central indicator of both intelligence and practical utility, drawing increasing attention from academia and industry alike (Chen et al., 2021; Jimenez et al., 2024; Jiang et al., 2024). Many powerful LLMs, such as Claude Opus 4.1 (Anthropic, 2025a), are already widely adopted in AI-assisted coding tools (Cursor, 2025; Anthropic, 2025). Through the ability to generate executable code, LLMs can substantially accelerate programming automation and reduce manual development effort.

To evaluate and advance these capabilities, a series of benchmarks have been developed (Wang et al., 2025a). Early efforts such as HumanEval (Chen et al., 2021) and MBPP (Austin et al., 2021) focused on short, algorithmic Python tasks. More recent benchmarks (Peng et al., 2024; Jimenez et al., 2024; Jain et al., 2025; Zheng et al., 2025; Zhu et al., 2025; Chai et al., 2025; Bytedance, 2025; Zan et al., 2025) target more challenging and realistic programming tasks, including competition-level problems and multilingual scenarios. However, constructing such benchmarks remains costly and labor-intensive, which hinders scalability and makes it difficult to guarantee both high difficulty and broad coverage. As shown in Table 1, widely used multilingual benchmarks

¹The code and benchmark are available at <https://anonymous.4open.science/r/autocodebench-iclr-6E45>. The remaining resources will be released after the double-blind review process due to their large size.

054
 055
 056
 057
 058
 Table 1: Comparison of Code Generation Benchmarks. **MLing**: MultiLingual; **MLogi**: MultiLogical,
 refers to programming problems that require the model to simultaneously implement multiple core
 functionalities. **HFree**: Human-Free; **BDist**: Balanced Distribution of multiple languages. **DSize**:
 Data Size. **PLen**: Problem Length. Further details are provided in the Appendix B.

| Benchmark | MLing | MLogi | HFree | BDist | Difficulty | Category | DSIZE | PLen |
|----------------------|-------|-------|-------|-------|------------|----------|-------------|--------------|
| HumanEval | ✗ | ✗ | ✗ | / | ★ | 5 | 164 | 134.1 |
| MBPP | ✗ | ✗ | ✗ | / | ★ | 6 | 378 | 50.5 |
| LiveCodeBench | ✗ | ✗ | ✗ | / | ★★★★ | 4 | 1100 | 469.6 |
| FullStackBench | ✓ | ✗ | ✗ | ✗ | ★★ | 12 | 1687 | 184.3 |
| McEval | ✓ | ✗ | ✗ | ✓ | ★★ | 9 | 2007 | 146.7 |
| AutoCodeBench | ✓ | ✓ | ✓ | ✓ | ★★★★ | 14 | 3920 | 498.2 |

060
 061
 062
 063
 064
 065
 066
 067
 such as FullStackBench (Bytedance, 2025) and McEval (Chai et al., 2025) suffer from imbalanced
 069 language distributions and limited diversity. These issues stem from the inherent bias of manual
 070 annotation. For example, most annotators are proficient in Python algorithmic tasks but lack expertise
 071 in domains such as Elixir-based communication development. Moreover, the rapid progress of LLMs
 072 has rendered many of the overly simple problems in these benchmarks obsolete. These limitations
 073 naturally raise a critical question: *Can we automatically construct high-quality code generation*
 074 *benchmarks that scale while ensuring both comprehensiveness and diversity?*

075
 076
 077 In this paper, we propose **AutoCodeGen**, an automated workflow centered on LLM–sandbox interac-
 078 tion, to synthesize challenging multilingual code generation datasets without manual annotations.
 079 Unlike previous data synthesis approaches (Luo et al., 2024; Wei et al., 2024b; Xu et al., 2025;
 080 Ahmad et al., 2025) that directly rely on LLMs to generate test functions, we adopt a more intuitive
 081 and reliable strategy. The LLMs first produce test inputs, and a multilingual sandbox executes these
 082 inputs to obtain the corresponding outputs. This design effectively mitigates a common limitation
 083 of LLMs, namely the tendency to generate incorrect outputs when confronted with high-difficulty
 084 problems. Furthermore, we introduce a reverse-generation paradigm, where code solutions and test
 085 functions are synthesized first, followed by the construction of the programming problem itself. This
 086 ensures that the resulting tasks are not only sufficiently challenging but also verifiably correct.

087
 088 Based on the automation workflow, we introduce **AutoCodeBench**, a large-scale, fully automated
 089 code generation benchmark, as shown in Table 1. Compared with previous multilingual bench-
 090 marks (Cassano et al., 2022; Bytedance, 2025; Chai et al., 2025), ours simultaneously offers high
 091 difficulty, diversity, and practicality, with a balanced distribution of problems across 20 programming
 092 languages. We intentionally include some multi-logical problems to test the LLMs’ capacity for
 093 multitasking in a single problem. The key contributions of this paper are as follows:

- 094 1. **AutoCodeGen**. We propose an automated workflow based on LLM-sandbox interaction,
 095 where LLMs generate test inputs and a sandbox executes them to obtain outputs, to create
 096 high-quality multilingual code generation tasks.
- 097 2. **AutoCodeBench**. We introduce AutoCodeBench, a large-scale, fully automatic code
 098 generation benchmark with 3,920 problems, evenly distributed across 20 programming
 099 languages, featuring high difficulty, practicality, and diversity. We also construct Lite and
 100 Complete versions to enable efficient and high-quality evaluation. The evaluation results
 101 show that current LLMs still struggle with complex and diverse multilingual programming
 102 tasks, especially in multi-logical and low-resource language scenarios.
- 103 3. **Training and Evaluation Resources**. We propose **AutoCodeInstruct**, a multilingual code
 104 generation training set constructed using the same pipeline as AutoCodeBench, ensuring
 105 comparable quality. We design a two-stage GRPO training to demonstrate the potential of
 106 this dataset in enhancing code generation capabilities. Besides, we release a **Multilingual**
 107 **Sandbox** with high-concurrency support, which can be employed for both model evaluation
 108 and RL training.

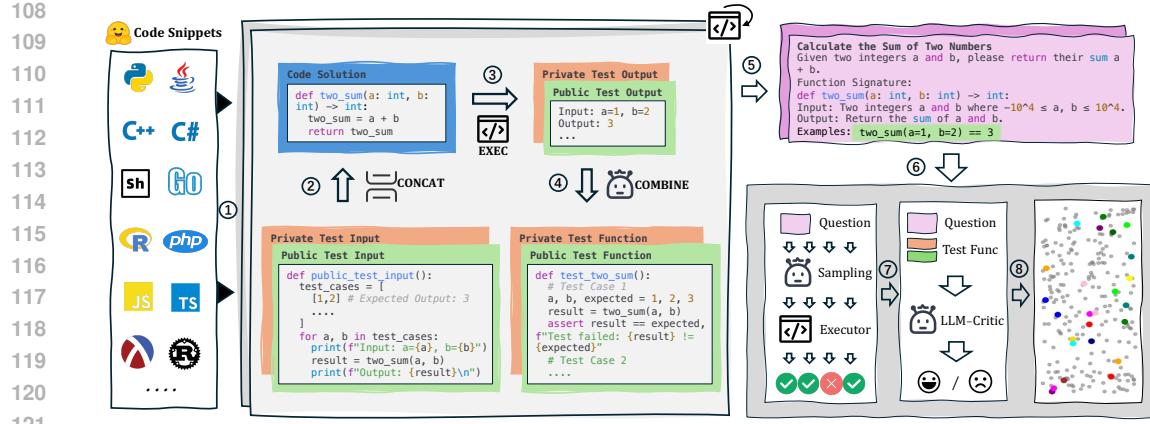


Figure 1: **The overview of AutoCodeGen.** It first generates code solution and the corresponding public/private test input functions based on multilingual code snippets (①). They are concatenated and executed in a sandbox to obtain test outputs, which are then combined by the LLM into complete test functions (②,③,④). Based on the code solution and test function, the LLM is prompted to generate accurate programming problems (⑤). Finally, a three-stage data filtering is applied: multiple sampling to remove too easy problems (⑥), LLM-as-Critic to discard low-quality ones (⑦), and diversity-based tagging to ensure distributional variety (⑧).

2 AUTOCodeGen & AUTOCodeBench

In this section, we first present how AutoCodeGen constructs the AutoCodeBench family of benchmarks, and then provide an overview of AutoCodeBench.

2.1 AUTOCodeGen

Our AutoCodeGen is a fully automated workflow based on LLM-sandbox interaction for constructing verifiable code generation datasets. It first generates large-scale multilingual data with guaranteed executability and correctness, then applies a three-stage filtering strategy to ensure the benchmark is challenging, high-quality, and diverse. As illustrated in Figure 1, the workflow includes four key stages: Code Solution Generation (①), Test Function Generation (②,③,④), Programming Problem Generation (⑤), and Data Filtering (⑥,⑦,⑧).

2.1.1 CODE SOLUTION GENERATION

We begin by extracting multilingual code snippets from Stack-Edu (Allal et al., 2025), a large-scale dataset of educational code filtered from The Stack v2 (Lozhkov et al., 2024), as seeds. These seeds span function-level, class-level, and file-level code, sourced from real GitHub repositories, ensuring diversity and practicality. Using a language-specific few-shot prompt, we guide DeepSeek-V3-0324 to refine and evolve these seeds into verifiable and self-contained code solutions. During this process, the model removes non-essential logic and adds appropriate comments for clarity. We then validate the correctness of the generated solutions by multilingual sandbox.

2.1.2 TEST FUNCTION GENERATION

We enhance efficiency and edge-case coverage by first generating test inputs via LLMs and then executing them in a sandbox to obtain the corresponding outputs. Specifically, it is divided into the following three steps:

Test Input Generation The test input functions (both public and private) are generated alongside the above code solution, ensuring alignment between the solution and its inputs. The public test input function includes no more than 3 basic cases and serves demonstration purposes; it will be embedded into the final programming problem as an illustrative usage. In contrast, the private test input function

Table 2: Programming language translation pairs.

| Origin | Target | Origin | Target | Origin | Target | Origin | Target | Origin | Target |
|--------|--------|--------|--------|--------|--------|------------|------------|--------|--------|
| Python | R | Python | Ruby | Java | Scala | Java | C# | Shell | Perl |
| Python | Elixir | Python | Julia | Java | Kotlin | JavaScript | PHP | C++ | Rust |
| Python | Swift | Python | Racket | Java | Dart | JavaScript | Typescript | | |

contains 7+ cases, including edge cases, and functions as the comprehensive test for verifying the correctness of the code solution.

Test Output Generation We concatenate the code solution with test input functions and execute them in the sandbox to obtain the corresponding test outputs.

Input-Output Integration We prompt DeepSeek-V3-0324 with both the test input functions and output results from the sandbox to generate coherent and verifiable test functions. Finally, we validate the correctness by executing the code solution together with the generated public and private test functions in the sandbox.

2.1.3 PROGRAMMING PROBLEM GENERATION

Generating high-quality programming problems is challenging, as it requires detailed and accurate problem descriptions. We find that models often omit key information when generating programming problems, such as the entry point specified in the test function. Therefore, we define a set of specifications, such as explicit input/output formats, function and class names, to ensure that the generated problems are detailed and well-structured. We prompt DeepSeek-V3-0324 to generate high-quality programming problems based on the code solution with appropriate comments and the corresponding test function, while embedding the public test function as example usage.

Through these three steps, we obtain a large-scale multilingual dataset, where each instance is represented as a tuple <programming problem, code solution, public test function, private test function>.

2.1.4 DATA FILTERING

Finally, We apply three filtering and sampling steps to ensure the high-difficulty, high-quality, and diversity of the final benchmark.

Difficulty Control Too simple programming problems are barely meaningful for evaluating the code generation capabilities of current LLMs. To address this, we employ a moderately capable code model, DeepSeek-Coder-V2-Lite, to filter out too easy problems. Specifically, we sample answers for each problem ten times using the model and validate the correctness via sandbox execution. We discard problems that are solved in all attempts. Take Python as an example, DeepSeek-Coder-V2-Lite can filter out 25.1% of the whole problems.

Quality Control During the aforementioned problem generation stage, we define six specifications to guide the generation of detailed and accurate programming problems. To further ensure high quality, we employ DeepSeek-R1-0528 to critique each <problem, test function> pair. Only the data whose test functions are completely accurate and fully aligned with the problem are retained.

Diversity Sampling We aim for our benchmark to cover as many real-world scenarios as possible. To this end, we perform diversity-based sampling on the existing data to construct the final benchmark. We use DeepSeek-V3-0324 to label each problem. We then divide the problems into different pools by category and perform cyclic sampling, ensuring a broad representation of programming scenarios.

2.1.5 APPROXIMATE LANGUAGE TRANSLATION

For Python, C++, Shell, Java, JavaScript, and Go, we directly use the workflow described above. For the other 14 languages, while the proposed workflow is still applicable, we choose to employ an approximate language translation approach due to their limited data resources and lack of diversity. We extract unused data from the dataset generated in Section 2.1.3 and translate them into the target

216 Table 3: Statistics of Dataset. **ACB**: AutoCodeBench; **Langs**: Languages; **Prob**: Problem; **Solu**:
 217 Solution; **Len**: Length; **E/M/H**: Easy/Medium/Hard. The difficulty level is determined by the number
 218 of passes in ten samplings of DeepSeek-Coder-V2-Lite. Problems with zero correct solutions
 219 are classified as hard, 1-5 correct solutions as medium, and those with more than five as easy.

| | #Problems | #Test Cases | #Langs | Prob Len | Solu Len | Difficulty (E/M/H) |
|--------------|-----------|-------------|--------|----------|----------|--------------------|
| ACB-Full | 3,920 | 37,777 | 20 | 498.2 | 487.5 | 646/846/2428 |
| ACB-Lite | 1,586 | 15,341 | 20 | 517.2 | 469.3 | 263/421/902 |
| ACB-Complete | 1,000 | 9,608 | 20 | 505.2 | 461.2 | 169/265/566 |

225
 226 low-resource language, as shown in Table 2. This ensures a sufficient and diverse dataset, which is
 227 further refined through the Data Filtering process in Section 2.1.4.
 228

229 2.2 AUTOCODEBENCH

231 2.2.1 DATA OVERVIEW

233 As shown in Table 1 and 3, AutoCodeBench(-Full) is a large-scale, high-difficulty multilingual
 234 benchmark. Over 60% of the problems are classified as hard problems, with each problem averaging
 235 498.2 characters and accompanied by 9.6 test cases, providing a challenging and comprehensive
 236 evaluation standard. The 20 languages are as follows: *Python, C++, Java, JavaScript(JS), Go, Shell,*
 237 *C#, Dart, Elixir, Julia, Kotlin, Perl, PHP, Racket, R, Ruby, Rust, Scala, Swift, TypeScript(TS)*.

238 To analyze the diversity and language coverage of AutoCodeBench, we first use Claude Sonnet
 239 4 to generate 20 language-agnostic task categories, and then employ DeepSeek-V3-0324 to
 240 classify each problem accordingly. Categories with less than 2% representation are merged into
 241 the “Other” group. AutoCodeBench covers 14 categories, demonstrating comprehensive coverage
 242 of practical programming scenarios. Besides, we analyze the distribution of problems across the
 243 20 programming languages. AutoCodeBench exhibits a relatively balanced distribution across
 244 languages, with no significant bias toward any specific one, further validating its completeness and
 245 representativeness as a multilingual benchmark. Detailed category and language distribution are
 246 provided in Appendix B.3.

247 2.2.2 AUTOCODEBENCH-LITE AND AUTOCODEBENCH-COMPLETE CONSTRUCTION

249 To facilitate quicker and more efficient model evaluations, we create **AutoCodeBench-Lite**, a
 250 simplified subset of AutoCodeBench. Specifically, we collect the problem-solving results from all
 251 models and sort the problems in ascending order based on the number of passes. After discarding
 252 problems with fewer than 2 passes, we select approximately 1,500 problems based on their pass count
 253 in ascending order. These problems, which have been solved correctly by existing models at least
 254 twice and have a certain level of difficulty, are selected to amplify the differences between the models.
 255 We use these problems as the set for the Lite version.

256 To enable evaluating base models, we further present **AutoCodeBench-Complete**, a completion-
 257 based version of ACB. Concretely, we select 1,000 data points from ACB-Lite to ensure a balanced
 258 distribution of 50 problems per programming language and use 3-shot demonstrations to evaluate the
 259 performance of base models. ACB-Complete can serve as a comprehensive benchmark for evaluating
 260 the multilingual code generation capabilities of base models.

261 3 EVALUATION

262 3.1 EVALUATION SETUP

266 We use the Pass@1 (%) (Chen et al., 2021) as the default evaluation metric. In terms of inference
 267 parameters, for proprietary models and open-source models that provide APIs, we access them
 268 through direct API calls. Other models are evaluated with official parameters when available, or
 269 with greedy decoding otherwise. All models are provided with our custom system prompt, which
 standardizes the output format: ***You are an expert programmer. Your task is to provide a code***

| | Average | Python | C++ | Java | JS | Go | Shell | Csharp | Dart | Elixir | Julia | Kotlin | Perl | PHP | Racket | R | Ruby | Rust | Scala | Swift | TS |
|--|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| Count | 3920 | 196 | 186 | 188 | 184 | 191 | 189 | 188 | 200 | 198 | 200 | 200 | 200 | 199 | 196 | 198 | 200 | 199 | 199 | 200 | 199 |
| <i>Current Upper Bound</i> | 75.3 | 65.3 | 75.8 | 80.9 | 60.9 | 71.7 | 72.9 | 88.4 | 78.0 | 97.5 | 78.5 | 90.5 | 64.5 | 53.8 | 89.8 | 75.8 | 81.0 | 62.8 | 78.4 | 78.5 | 61.3 |
| Proprietary Models and 2008+ Open-source Models | | | | | | | | | | | | | | | | | | | | | |
| Claude Opus 4.1 (20250805) | 55.4 | 42.3 | 49.5 | 56.4 | 42.9 | 44.0 | 50.0 | 78.4 | 59.5 | 86.9 | 59.5 | 74.5 | 47.5 | 31.2 | 73.0 | 55.6 | 61.0 | 39.2 | 50.8 | 57.0 | 47.7 |
| Claude Sonnet 4 (20250514) | 51.1 | 37.2 | 46.8 | 52.7 | 34.8 | 41.9 | 48.9 | 72.4 | 53.5 | 81.8 | 49.0 | 71.5 | 45.0 | 34.7 | 68.9 | 50.5 | 54.5 | 36.2 | 48.2 | 48.0 | 44.2 |
| Claude Opus 4.1 (20250805) | 52.6 | 38.3 | 48.4 | 53.7 | 41.3 | 38.7 | 46.8 | 75.4 | 55.0 | 80.3 | 57.5 | 76.0 | 45.0 | 29.6 | 64.8 | 51.0 | 55.5 | 39.7 | 51.8 | 53.0 | 47.2 |
| Claude Sonnet 4 (20250514) | 49.3 | 35.7 | 47.3 | 52.7 | 38.0 | 37.7 | 47.9 | 72.9 | 51.0 | 74.2 | 51.0 | 72.0 | 44.0 | 30.7 | 63.8 | 44.4 | 51.5 | 35.2 | 45.2 | 45.5 | 44.2 |
| GPT-5 (20250807) | 53.5 | 44.4 | 51.6 | 48.9 | 44.6 | 45.0 | 47.3 | 75.9 | 53.5 | 84.3 | 53.0 | 70.5 | 43.7 | 36.2 | 58.5 | 54.0 | 60.5 | 40.7 | 49.7 | 58.5 | 46.7 |
| o3-high (20250416) | 51.1 | 40.8 | 47.3 | 53.2 | 40.8 | 22.0 | 49.5 | 68.3 | 55.0 | 80.8 | 54.5 | 72.0 | 44.0 | 32.7 | 53.1 | 47.5 | 59.0 | 42.2 | 51.3 | 59.0 | 47.2 |
| o4-mini (2025-04-16) | 50.0 | 42.3 | 46.8 | 51.6 | 40.2 | 31.4 | 45.2 | 68.3 | 54.0 | 82.3 | 49.0 | 74.0 | 44.0 | 30.2 | 45.4 | 43.4 | 59.0 | 40.2 | 50.3 | 54.0 | 45.7 |
| GPT4.1 (2025-04-14) | 49.7 | 37.2 | 46.4 | 48.9 | 34.8 | 37.2 | 36.7 | 74.4 | 46.6 | 76.8 | 50.0 | 72.0 | 43.0 | 29.2 | 50.5 | 42.4 | 53.7 | 37.2 | 44.9 | 45.5 | 46.2 |
| Grok4.2 | 50.9 | 41.6 | 48.7 | 50.9 | 37.5 | 37.5 | 47.9 | 75.0 | 53.0 | 80.8 | 49.0 | 74.0 | 44.0 | 27.1 | 53.8 | 46.5 | 61.2 | 39.7 | 49.5 | 51.0 | 47.2 |
| Grok4.2 Pro | 48.7 | 40.3 | 47.5 | 50.9 | 37.0 | 37.2 | 45.2 | 70.9 | 54.0 | 86.7 | 54.0 | 70.0 | 41.0 | 29.5 | 50.5 | 49.5 | 64.7 | 36.6 | 49.5 | 41.7 | 41.7 |
| Grok4.2 Flash | 45.7 | 39.3 | 44.1 | 50.0 | 33.2 | 33.0 | 37.8 | 68.3 | 49.5 | 64.0 | 47.5 | 70.0 | 39.5 | 24.1 | 38.3 | 51.5 | 53.0 | 36.2 | 44.2 | 46.5 | 41.2 |
| DeepSeek-V1.3-250821 | 48.2 | 39.3 | 47.3 | 53.7 | 37.0 | 30.4 | 38.3 | 71.9 | 49.5 | 75.8 | 53.0 | 67.5 | 43.0 | 29.6 | 52.6 | 48.0 | 54.0 | 39.2 | 45.2 | 49.0 | 38.2 |
| DeepSeek-V1.3-250821 | 46.2 | 35.7 | 44.1 | 54.3 | 35.3 | 29.3 | 36.2 | 68.3 | 44.0 | 72.7 | 49.5 | 64.5 | 44.0 | 29.6 | 52.6 | 47.0 | 51.0 | 33.2 | 42.2 | 47.0 | 42.7 |
| DeepSeek-Coder-V2-Instruct | 37.7 | 29.1 | 34.9 | 34.0 | 27.7 | 29.8 | 31.4 | 63.8 | 33.5 | 60.6 | 37.5 | 58.5 | 35.5 | 25.1 | 41.8 | 35.4 | 45.0 | 22.6 | 33.2 | 38.0 | 35.7 |
| Hunyuan-Turbos-20250716 | 43.8 | 32.4 | 34.9 | 47.9 | 32.6 | 34.6 | 38.3 | 64.8 | 44.5 | 70.7 | 47.0 | 62.0 | 42.0 | 30.2 | 45.9 | 39.9 | 53.0 | 30.7 | 39.2 | 39.5 | 42.2 |
| GLM-4.5-enable | 46.6 | 41.0 | 43.2 | 47.9 | 34.8 | 37.8 | 43.9 | 70.5 | 42.0 | 72.5 | 47.5 | 66.0 | 43.5 | 28.6 | 50.0 | 45.0 | 54.5 | 31.6 | 41.0 | 46.0 | 42.2 |
| KimI-20905-preview | 46.8 | 36.2 | 38.2 | 47.3 | 37.0 | 35.1 | 41.5 | 68.3 | 50.5 | 78.8 | 48.5 | 66.5 | 41.5 | 30.7 | 55.6 | 40.4 | 49.5 | 31.7 | 45.7 | 48.5 | 42.7 |
| ERNIE-XI-Turbo-32K | 39.6 | 39.4 | 17.8 | 33.2 | 32.6 | 37.4 | 33.9 | 46.0 | 33.0 | 68.9 | 54.0 | 49.5 | 39.5 | 23.9 | 45.3 | 44.3 | 48.0 | 20.8 | 40.4 | 44.0 | 37.7 |
| Qwen-235B-AB2B-Thinking-2507 | 47.7 | 37.8 | 41.9 | 48.4 | 39.7 | 39.8 | 45.2 | 71.9 | 46.0 | 79.8 | 48.5 | 58.0 | 40.5 | 29.1 | 56.6 | 49.0 | 55.0 | 35.7 | 40.4 | 46.0 | 44.2 |
| Qwen-3-Code4-AB8B-Instruct | 44.8 | 39.4 | 41.1 | 51.1 | 27.9 | 31.4 | 41.1 | 63.0 | 36.5 | 73.7 | 49.5 | 61.3 | 41.0 | 27.2 | 56.3 | 42.7 | 51.5 | 25.4 | 42.1 | 47.5 | 41.8 |
| Qwen-3-Code4-AB8B-Instruct | 43.1 | 35.7 | 38.2 | 49.5 | 29.3 | 33.5 | 40.4 | 67.3 | 39.5 | 59.1 | 46.0 | 59.5 | 44.5 | 26.1 | 49.5 | 44.0 | 46.5 | 24.6 | 37.7 | 40.0 | 43.2 |
| Seed1.6-Thinking-250715 | 45.0 | 40.3 | 45.2 | 50.0 | 33.2 | 38.2 | 39.9 | 67.3 | 36.5 | 67.7 | 51.0 | 61.0 | 41.0 | 26.1 | 51.0 | 44.9 | 55.5 | 27.6 | 37.2 | 46.5 | 38.7 |
| Seed1.6-enabled (250615) | 45.3 | 39.8 | 44.6 | 46.3 | 28.3 | 40.8 | 44.1 | 60.3 | 39.5 | 69.7 | 51.0 | 58.0 | 41.5 | 25.6 | 52.6 | 51.0 | 52.0 | 28.6 | 41.7 | 47.5 | 41.2 |
| Seed1.6-enabled (250615) | 42.9 | 35.2 | 40.3 | 46.8 | 32.6 | 34.6 | 35.1 | 70.9 | 42.5 | 69.7 | 45.0 | 49.5 | 39.0 | 23.1 | 49.5 | 40.4 | 46.5 | 28.1 | 37.2 | 40.7 | 42.7 |
| Open-source Models below 200B | | | | | | | | | | | | | | | | | | | | | |
| GLM-4.5-Air-enable | 40.8 | 39.3 | 37.6 | 39.4 | 31.0 | 39.8 | 36.7 | 66.3 | 38.0 | 61.5 | 42.0 | 53.0 | 40.5 | 27.1 | 40.3 | 39.0 | 47.0 | 25.1 | 30.5 | 38.5 | 42.7 |
| Qwen3-Next-80B-AB3-Thinking | 40.6 | 38.3 | 39.8 | 43.6 | 38.0 | 33.0 | 37.4 | 66.3 | 24.0 | 59.1 | 43.5 | 43.5 | 42.5 | 25.1 | 34.2 | 46.4 | 50.5 | 26.1 | 35.2 | 43.7 | 41.2 |
| Qwen3-Next-80B-AB3-Instruct | 39.6 | 36.7 | 35.5 | 44.1 | 32.6 | 29.8 | 33.2 | 62.8 | 35.5 | 63.1 | 39.0 | 41.0 | 42.0 | 27.6 | 41.3 | 39.4 | 39.9 | 26.4 | 34.7 | 42.2 | 40.2 |
| Qwen3-32B | 47.7 | 37.8 | 40.4 | 46.4 | 34.4 | 36.3 | 39.4 | 67.8 | 44.5 | 64.5 | 40.0 | 40.0 | 42.0 | 29.4 | 34.8 | 37.3 | 39.0 | 27.4 | 32.0 | 38.7 | 40.2 |
| Qwen3-14B | 37.6 | 37.8 | 35.5 | 35.1 | 30.4 | 34.4 | 36.2 | 60.8 | 26.0 | 62.1 | 34.5 | 44.5 | 37.5 | 37.5 | 44.9 | 36.9 | 45.5 | 24.6 | 36.0 | 38.7 | 37.7 |
| Qwen3-8B | 28.5 | 28.1 | 22.6 | 21.8 | 28.3 | 29.3 | 27.1 | 52.8 | 21.0 | 43.9 | 29.0 | 35.5 | 35.5 | 18.6 | 13.3 | 30.8 | 37.0 | 12.6 | 21.1 | 22.0 | 37.7 |
| Qwen3-1.7B | 11.2 | 16.8 | 5.4 | 4.8 | 12.5 | 9.9 | 11.7 | 19.6 | 7.0 | 20.7 | 11.0 | 9.0 | 19.5 | 7.5 | 5.6 | 9.6 | 21.0 | 0.0 | 2.5 | 10.0 | 19.6 |
| Qwen3-32B | 31.0 | 26.5 | 21.5 | 29.5 | 28.0 | 25.5 | 24.0 | 59.3 | 27.5 | 52.0 | 28.0 | 45.0 | 34.5 | 21.6 | 28.6 | 22.7 | 36.5 | 16.1 | 26.6 | 30.0 | 35.2 |
| Qwen3-14B | 28.6 | 24.5 | 22.6 | 32.4 | 27.2 | 16.8 | 23.0 | 50.8 | 21.0 | 42.4 | 22.5 | 42.0 | 34.5 | 24.6 | 28.1 | 26.3 | 33.5 | 17.1 | 20.6 | 26.5 | 32.7 |
| Qwen3-8B | 23.3 | 22.4 | 11.3 | 25.0 | 22.8 | 18.3 | 22.3 | 42.2 | 17.0 | 41.4 | 18.5 | 36.0 | 29.5 | 18.1 | 23.5 | 16.2 | 27.0 | 7.5 | 19.1 | 17.5 | 29.1 |
| Qwen3-1.7B | 7.9 | 8.7 | 1.1 | 2.7 | 8.2 | 3.7 | 11.7 | 9.5 | 3.5 | 17.2 | 6.6 | 11.5 | 15.5 | 7.5 | 4.6 | 7.6 | 14.5 | 0.5 | 3.0 | 6.0 | 14.6 |
| Qwen2.5-Coder-32B-Instruct | 35.8 | 29.6 | 27.4 | 33.0 | 29.9 | 23.0 | 29.3 | 58.3 | 34.5 | 59.6 | 35.5 | 58.0 | 38.5 | 26.1 | 35.7 | 31.0 | 40.4 | 23.1 | 29.6 | 39.0 | 35.2 |
| Qwen2.5-Coder-7B-Instruct | 22.5 | 19.9 | 8.6 | 22.3 | 21.2 | 12.6 | 21.3 | 38.7 | 18.5 | 47.0 | 18.0 | 39.0 | 27.5 | 15.1 | 24.0 | 17.7 | 29.5 | 7.0 | 19.1 | 17.0 | 24.6 |
| Qwen2.5-Coder-1.5B-Instruct | 10.3 | 12.2 | 2.7 | 4.8 | 12.0 | 7.3 | 14.4 | 17.6 | 6.5 | 35.4 | 4.0 | 16.5 | 15.0 | 5.0 | 7.1 | 7.6 | 15.5 | 1.0 | 4.5 | 5.0 | 11.1 |
| DeepSeek-Coder-33B-Instruct | 28.5 | 25.0 | 24.2 | 29.3 | 24.5 | 29.8 | 22.3 | 54.8 | 17.5 | 67.7 | 14.5 | 52.0 | 29.5 | 19.1 | 28.1 | 18.7 | 33.0 | 8.0 | 24.1 | 18.0 | 29.1 |
| DeepSeek-Coder-6.7B-Instruct | 20.5 | 18.9 | 12.9 | 19.7 | 19.6 | 21.5 | 16.0 | 44.2 | 11.5 | 47.5 | 15.5 | 45.5 | 21.5 | 10.6 | 13.5 | 27.5 | 6.0 | 11.1 | 8.0 | 23.6 | |
| Seed-Coder-8B-Instruct | 32.3 | 23.5 | 23.7 | 33.5 | 28.8 | 22.5 | 20.7 | 54.8 | 30.5 | 57.1 | 33.0 | 52.5 | 34.0 | 25.1 | 36.7 | 28.3 | 35.5 | 15.6 | 29.6 | 28.0 | 31.2 |
| OpenCoder-8B-Instruct | 19.3 | 14.3 | 9.7 | 19.1 | 12.5 | 17.3 | 21.3 | 32.3 | 15.5 | 34.3 | 15.5 | 26.5 | 14.6 | 7.5 | 15.8 | 17.2 | 29.5 | 6.0 | 17.1 | 15.0 | 21.1 |

Table 4: Pass@1 (%) performance of different models for AutoCodeBench. **Current Upper Bound** represents the Pass@1 value calculated by taking the union of problems correctly solved by all models. **Blue** and **Green** denotes reasoning and non-reasoning modes.

solution within a single Markdown code block for the given programming problem. Do not include any direct execution commands, test cases, or usage examples within the code block.

3.2 MAIN RESULTS

We comprehensively evaluate the performance on ACB, with results across different programming languages shown in Tables 4. The results of ACB-Lite and leaderboards are shown in Table 9, Figure 7 and Figure 8.

Results show that ACB is highly challenging, as no model surpasses 55.5 average score, indicating that current LLMs still struggle with complex, practical multilingual problems. Among all models, Claude Opus 4.1 consistently achieves the best performance in both reasoning and non-reasoning modes, confirming its strength across diverse coding tasks and aligning with observations from SWE-bench (Jimenez et al., 2024). Finally, while individual models perform moderately, their combined upper bound reaches 75.3, revealing complementary strengths and substantial room for improvement, as no single model dominates across all languages.

3.3 PERFORMANCE ACROSS POPULAR AND LOW-RESOURCE PROGRAMMING LANGUAGES

We select five models with similar performance levels (ranging from 47.7 to 49.3) and evaluate their performance differences across popular and low-resource languages. As shown in Figure 2, the difference in average Pass@1 scores among the models for popular languages is small (Δ 3.1). However, the performance gap between models widens (Δ 6.3) in low-resource languages, suggesting that low-resource programming languages may have received insufficient attention in model development. Besides, since we use the moderately capable DeepSeek-Coder-V2-Lite as a filter to remove simple problems, the Pass@1 scores of top models on popular languages are relatively low. However, because the filter itself performs poorly on low-resource languages, many problems that appear trivial to top models are not filtered out, resulting in higher Pass@1 scores in these languages than in popular ones. This further highlights the pronounced disparities in low-resource language capabilities across models and underscores the need for greater community attention to this issue.

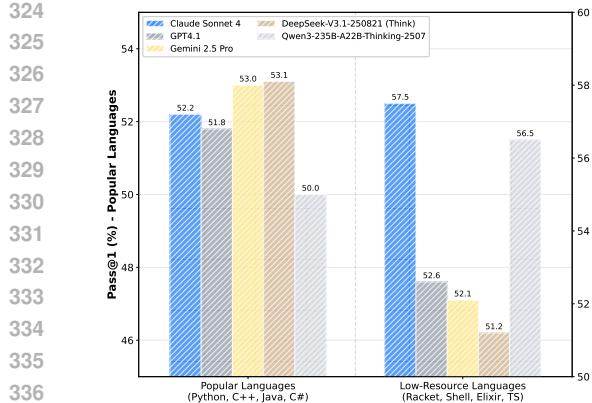


Figure 2: The performance comparison of different models across two language sets.

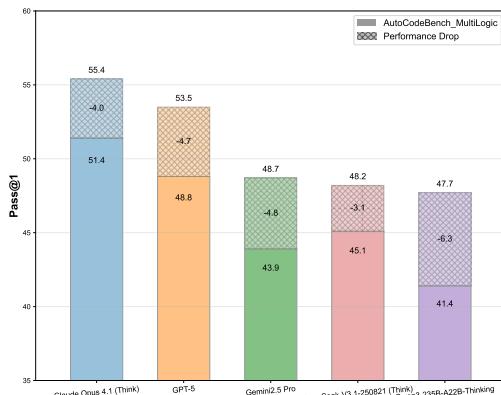


Figure 3: Performance drop of models on multi-logic problems (1,622) compared to full dataset.

3.4 PERFORMANCE ACROSS MULTI-LOGIC PROGRAMMING PROBLEMS

A key feature that distinguishes AutoCodeBench from prior benchmarks is the inclusion of multi-logical problems. These problems require models to implement multiple distinct functions or classes within a single task, challenging their ability to handle multiple core demands simultaneously. We use DeepSeek-V3-0324 to identify all multi-logical problems in AutoCodeBench and evaluate model performance on them. The results, shown in Figure 3, reveal a significant performance drop for all models when faced with multi-logical tasks. Among them, Claude Opus 4.1 and DeepSeek-V3.1 exhibit relatively smaller declines, while the other models show larger drops. These findings highlight a key limitation: current models still struggle with multi-logical problem solving, an ability that is particularly critical for real-world code agent applications.

3.5 PERFORMANCE ANALYSIS OF MULTI-TURN REFINEMENT WITH SANDBOX FEEDBACK

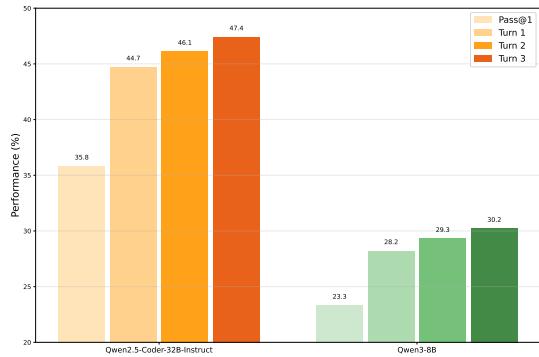
As shown in Figure 4, we evaluate how models leverage execution error messages to iteratively refine their code solutions. The results highlight the substantial value of our multilingual sandbox error feedback across all evaluated models. Qwen2.5-Coder-32B-Instruct achieves remarkable improvement from 35.8% to 47.4% after three refinement turns, while Qwen3-8B shows consistent progress from 23.3% to 30.2%. The most significant performance gains occur during the first refinement turn, with diminishing returns in subsequent iterations. This pattern suggests that models can effectively leverage execution feedback to identify and correct common coding errors, though the complexity of remaining problems increases with each iteration. The consistent improvement across different model scales indicates that multi-turn refinement with sandbox feedback is a valid strategy for enhancing code generation quality.

3.6 AUTOCodeBENCH-COMPLETE: EVALUATING BASE MODEL CAPABILITIES

Table 5 presents a performance comparison between base models on ACB-Complete and chat models on ACB-Full. Among models with 8B parameters or fewer, Seed-Coder-8B demonstrates superior performance in ACB-Complete, consistent with its strong showing on ACB-Full. This consistency suggests that the pretraining process effectively equipped Seed-Coder-8B models with strong multilingual programming capabilities, enabling them to handle diverse coding scenarios across multiple languages. Besides, an interesting observation arises when comparing Qwen2.5-Coder-7B and OpenCoder-8B. While it outperforms OpenCoder-8B on ACB-Full, the trend reverses on ACB-Complete. This suggests that Qwen2.5-Coder-7B may have undergone more effective post-training on multilingual code generation data.

378
379 Table 5: The pass@1 values of chat models
380 (ACB-Full) and base models (ACB-Complete).
381

| | ACB-Full | ACB-Complete |
|---------------------|----------|--------------|
| 30B+ Models | | |
| DeepSeek-Coder-V2 | 37.7 | 39.0 |
| Qwen2.5-72B | 34.3 | 35.9 |
| Qwen2.5-Coder-32B | 35.8 | 35.5 |
| ~8B Models | | |
| Seed-Coder-8B | 32.3 | 31.6 |
| Qwen3-8B | 23.3 | 22.6 |
| OpenCoder-8B | 19.3 | 26.1 |
| Qwen2.5-Coder-7B | 22.5 | 24.6 |
| DeepSeek-Coder-6.7B | 20.5 | 22.9 |

387
388 Figure 4: Performance improvement across multi-
389 turn refinement with sandbox feedback.
390393 Table 6: Results of two-stage GRPO and SFT with AutoCodeInstruct.
394

| Model | ACB-Full | ACB-Lite | LiveCodeBench-V6 | FullStackBench | McEval |
|----------------------------|----------------------|----------------------|----------------------|----------------------|----------------------|
| Qwen2.5-Coder-7B-Instruct | 22.5 | 21.5 | 18.3 | 41.1 | 57.2 |
| + first-stage GRPO | 25.0 ^{↑2.5} | 24.8 ^{↑3.3} | 18.3 ^{↑0.0} | 46.9 ^{↑5.8} | 58.6 ^{↑1.4} |
| + second-stage GRPO | 27.4 ^{↑4.9} | 27.6 ^{↑6.1} | 17.1 ^{↓1.2} | 47.7 ^{↑6.6} | 58.4 ^{↑1.2} |
| + SFT | 28.9 ^{↑6.4} | 29.0 ^{↑6.5} | 17.7 ^{↓0.6} | 47.7 ^{↑6.6} | 63.1 ^{↑5.9} |
| Qwen2.5-Coder-32B-Instruct | 35.8 | 37.4 | 24.0 | 57.1 | 64.5 |
| + first-stage GRPO | 38.3 ^{↑2.5} | 39.5 ^{↑2.1} | 25.1 ^{↑1.1} | 58.3 ^{↑1.2} | 65.4 ^{↑0.9} |
| + second-stage GRPO | 41.6 ^{↑5.8} | 45.3 ^{↑7.9} | 28.0 ^{↑4.0} | 59.7 ^{↑2.6} | 66.1 ^{↑1.6} |
| + SFT | 41.9 ^{↑6.1} | 46.2 ^{↑8.8} | 30.3 ^{↑6.3} | 58.7 ^{↑1.6} | 69.5 ^{↑5.0} |

404
405

4 AUTOCODEINSTRUCT

406

408 **AutoCodeInstruct** Besides the evaluation benchmarks, we further construct AutoCodeInstruct, a
409 training dataset of comparable quality to AutoCodeBench. Specifically, we collect data generated
410 during the AutoCodeGen process that does not overlap with AutoCodeBench, and repeatedly sample
411 DeepSeek-V3-0324. Problems with excessively high pass rates ($>80\%$) or low pass rates ($<40\%$)
412 are filtered out to ensure both solvability and appropriate difficulty. We further apply a two-stage
413 deduplication strategy (MinHash + LLM-as-Judge) across existing code benchmarks. The resulting
414 dataset contains 37K verifiable problems spanning 20 programming languages.

415 **Training Setup** We conduct RL experiments via a two-stage GRPO (Shao et al.,
416 2024) training strategy to unleash the potential of AutoCodeInstruct, based on the
417 Qwen2.5-Coder-7B/32B-Instruct models. Concretely, we apply a data filtering strategy
418 by sampling the responses 15 times from the Instruct models and filter out easy problems with
419 pass rates above 0.6. The remaining problems are divided into solve-partial and solve-none parts
420 depending on whether the pass rate is zero. In the first stage, only the solve-partial problems join
421 training, with a rollout size of 8. In the second stage, we incorporate both solve-partial and solve-none
422 problems, and increase the rollout size to 16 to enable better exploration for harder problems. **Besides,**
423 **for all solve-partial and solve-none problems, we additionally obtain correct code solutions from**
424 **DeepSeek-V3-0324 and perform SFT on Qwen2.5-Coder-7B/32B-Instruct.** Details of
the training configurations are presented in Appendix I.

425 **Results** As shown in Table 6, after the first-stage GRPO, both models achieve noticeable gains
426 on in-domain benchmarks (ACB-Full and ACB-Lite), which suggests that they begin to learn
427 how to stably consolidate existing knowledge. The second-stage GRPO enables the models to
428 tackle harder problems, effectively pushing their multilingual capability boundaries and leading
429 to significant performance improvements. Surprisingly, the models also show consistent gains on
430 out-of-domain multilingual benchmarks such as FullStackBench and McEval. The performance
431 of Qwen2.5-Coder-32B-Instruct on the programming contest task LiveCodeBench-V6
(20250201–20250501) also improves by 4 points. **In addition, after the SFT stage, both models**

432 achieve even larger performance improvements than those obtained from GRPO, thanks to the distilled
 433 correct code solutions from DeepSeek-V3-0324. These results indicate that AutoCodeInstruct
 434 enhances the comprehensive code generation capability of models and demonstrate the effectiveness
 435 of our approach in synthesizing high-quality training datasets for code LLMs.
 436

437 5 RELATED WORK

438
 439 **Code Generation Benchmarks** The rapid evolution of code LLMs, ranging from open-source models (Roziere et al., 2023; Zhu et al., 2024; Hui et al., 2024b) to proprietary LLMs (Anthropic, 2025a; OpenAI, 2024; 2025a; Gemini, 2025) series, has reshaped code generation, create a demand for robust and contemporary code generation benchmarks. Pioneering benchmarks like HumanEval (Chen et al., 2021) and MBPP (Austin et al., 2021) established foundational correctness on small Python tasks but suffer from contamination and limited language coverage. Later benchmarks target more complex settings, such as competition-level challenges (Hendrycks et al., 2021; Li et al., 2022; Jain et al., 2025; Wang et al., 2025b; Zheng et al., 2025) and multilingual scenarios (Cassano et al., 2022; Peng et al., 2024; Zhang et al., 2024; Jimenez et al., 2024; Chai et al., 2025; Bytedance, 2025; Zhang et al., 2025b;a;b). McEval (Chai et al., 2025) is a massively multilingual benchmark covering 40 languages for generation, explanation, and completion tasks. FullStackBench (Bytedance, 2025) assesses LLMs in realistic, multi-domain scenarios across 16 languages, employing a novel execution environment. However, due to the challenges of manual annotation, these benchmarks suffer from issues such as limited diversity and insufficient difficulty, making them difficult to scale in line with the evolving demand for high-quality evaluation. By comparison, our AutoCodeBench series adopts a fully automated and scalable approach to create realistic, diverse, and high-difficulty tasks. A recent trend, exemplified by the SWE-Bench series (Jimenez et al., 2024; Zan et al., 2025; Rashid et al., 2025; He et al., 2025), focuses on evaluating LLMs in real-world software engineering tasks such as GitHub issue solving, thereby assessing models’ comprehensive capabilities beyond atomic-level code generation. By comparison, AutoCodeBench specifically targets LLMs’ atomic-level code generation abilities, which remain a crucial foundation for overall model performance.

450
 451
 452
 453
 454
 455
 456
 457
 458
 459 **Code Data Synthesis** To reduce dependence on manually curated data, a growing body of research explores automatic data synthesis to augment the training of Code LLMs (Luo et al., 2024; Wei et al., 2024b; Zheng et al., 2024; Wu et al., 2024; Yu et al., 2024; Ahmad et al., 2025; Xu et al., 2025). For instance, Evol-Instruct (Luo et al., 2024) uses heuristic prompts to guide LLMs in evolving existing programming problems, thereby increasing their diversity and difficulty. OSS-Instruct (Wei et al., 2024b) prompts LLMs to generate new coding problems and solutions from raw, open-source code snippets. KodCode (Xu et al., 2025) synthesizes a broad spectrum of Python coding tasks—including questions, solutions, and test cases—and ensures correctness through a systematic self-verification procedure. Some other methods focus on model self-improvement (Wu et al., 2024; Wei et al., 2024a; Chen et al., 2025b; Zhou et al., 2025; Zhang et al., 2025c). For instance, Inverse-Instruct (Wu et al., 2024) is a self-improvement technique that generates new instructions by “back-translating” code from an LLM’s own training set, reducing the need to distill from more powerful proprietary models. Collectively, these data synthesis methods significantly reduce the reliance on manual curation and enable the continuous expansion of the problem space for training. Our work extends this paradigm of automation from data augmentation to the benchmark creation process. By leveraging extensive LLM-sandbox interaction, our pipeline not only automates the synthesis of verifiable test problems but can also be naturally repurposed for synthesizing high-quality training datasets.

475 6 FUTURE WORKS

476
 477
 478 Ensuring high-quality and reliable code data synthesis remains fundamentally challenging. Although
 479 the sandbox provides strong guarantees regarding the correctness of code solutions and their alignment
 480 with test functions, the intrinsic quality of the synthesized programming problems and the
 481 completeness of their test coverage cannot be fully ensured. AutoCodeGen incorporates an LLM-as-
 482 Critic stage and careful prompt engineering to mitigate such risks, yet these mechanisms inevitably
 483 operate under uncertainty. To further improve reliability, future iterations of AutoCodeGen will
 484 incorporate repeated verification rounds using increasingly capable LLMs. Double- or triple-check
 485 validation loops are expected to improve accuracy, albeit at the cost of increased computational
 overhead and reduced pipeline efficiency.

486 Moreover, repository-level evaluation and data synthesis represent a more realistic and demanding
 487 setting, while the current multi-logic analysis provides only an initial bridge between fine-grained
 488 function behavior and higher-level software engineering reasoning. Moving forward, we aim to
 489 extend AutoCodeGen toward SWE-Bench- and Terminal-Bench-style domains, enabling automated
 490 synthesis of repo-level tasks together with the stateful sandbox environments required for such
 491 scenarios. This represents a crucial step toward scalable and fully autonomous code evaluation
 492 frameworks.

493 7 CONCLUSION

494 In this paper, we explored the large-scale and automated construction of code generation benchmarks.
 495 We introduce AutoCodeGen, an automated workflow based on LLM-Sandbox interaction, designed
 496 to generate multilingual verifiable code data without any manual annotation. Through this novel
 497 approach, we have successfully built AutoCodeBench, a large-scale, human-free code generation
 498 benchmark. AutoCodeBench contains 3,920 problems, evenly distributed across 20 programming
 499 languages, and is characterized by its high difficulty, practicality, and diversity. We also provide
 500 AutoCodeBench-Lite and AutoCodeBench-Complete, for efficient and high-quality evaluation of
 501 both chat and base LLMs. Our evaluation of more than 40 open-source and proprietary LLMs
 502 reveals that even the most advanced models still face challenges when confronted with the complex
 503 and diverse multilingual tasks set by AutoCodeBench. Besides, we construct AutoCodeInstruct, a
 504 large-scale, high-quality multilingual training dataset, and validate its effectiveness through GRPO.
 505

506 REFERENCES

507 Wasi Uddin Ahmad, Aleksander Ficek, Mehrzad Samadi, Jocelyn Huang, Vahid Noroozi, Somshubra
 508 Majumdar, and Boris Ginsburg. Opencodeinstruct: A large-scale instruction tuning dataset for
 509 code llms, 2025. URL <https://arxiv.org/abs/2504.04030>.

510 Loubna Ben Allal, Anton Lozhkov, Elie Bakouch, Gabriel Martín Blázquez, Guilherme Penedo,
 511 Lewis Tunstall, Andrés Marafioti, Hynek Kydliček, Agustín Piqueres Lajarín, Vaibhav Srivastav,
 512 Joshua Lochner, Caleb Fahlgren, Xuan-Son Nguyen, Clémentine Fourrier, Ben Burtenshaw, Hugo
 513 Larcher, Haojun Zhao, Cyril Zakka, Mathieu Morlon, Colin Raffel, Leandro von Werra, and
 514 Thomas Wolf. Smollm2: When smol goes big – data-centric training of a small language model,
 515 2025. URL <https://arxiv.org/abs/2502.02737>.

516 Anthropic. Claude code, 2025. URL <https://claude.com/product/claude-code>.

517 Anthropic. Introducing claude 4, 2025a. URL <https://www.anthropic.com/news/claude-4>.

518 Anthropic. Claude opus 4.1, 2025b. URL <https://www.anthropic.com/news/claude-4-1>.

519 Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan,
 520 Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, and Charles Sutton. Program synthesis with large
 521 language models, 2021. URL <https://arxiv.org/abs/2108.07732>.

522 Baidu. Ernie-x1-turbo-32k, 2025. URL <https://cloud.baidu.com/doc/WENXINWORKSHOP/s/Wm9cvy6rl>.

523 Bytedance. Fullstack bench: Evaluating llms as full stack coders, 2025. URL <https://arxiv.org/abs/2412.00535>.

524 Federico Cassano, John Gouwar, Daniel Nguyen, Sydney Nguyen, Luna Phipps-Costin, Donald
 525 Pinckney, Ming-Ho Yee, Yangtian Zi, Carolyn Jane Anderson, Molly Q Feldman, Arjun Guha,
 526 Michael Greenberg, and Abhinav Jangda. Multipl-e: A scalable and extensible approach to
 527 benchmarking neural code generation, 2022. URL <https://arxiv.org/abs/2208.08227>.

540 Linzheng Chai, Shukai Liu, Jian Yang, Yuwei Yin, JinKe, Jiaheng Liu, Tao Sun, Ge Zhang, Changyu
 541 Ren, Hongcheng Guo, Noah Wang, Boyang Wang, Xianjie Wu, Bing Wang, Tongliang Li, Liqun
 542 Yang, Sufeng Duan, Zhaoxiang Zhang, and Zhoujun Li. Mceval: Massively multilingual code
 543 evaluation. In *The Thirteenth International Conference on Learning Representations*, 2025. URL
 544 <https://openreview.net/forum?id=UunCPtP01z>.

545 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared
 546 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri,
 547 Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan,
 548 Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian,
 549 Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios
 550 Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino,
 551 Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders,
 552 Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa,
 553 Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob
 554 McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating
 555 large language models trained on code, 2021. URL <https://arxiv.org/abs/2107.03374>.

556 Wei-Lin Chen, Zhepei Wei, Xinyu Zhu, Shi Feng, and Yu Meng. Do llm evaluators prefer themselves
 557 for a reason?, 2025a. URL <https://arxiv.org/abs/2504.03846>.

558 Xiancai Chen, Zhengwei Tao, Kechi Zhang, Changzhi Zhou, Xinyu Zhang, Wanli Gu, Yuanpeng
 559 He, Mengdi Zhang, Xunliang Cai, Haiyan Zhao, and Zhi Jin. Revisit self-debugging with self-
 560 generated tests for code generation. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and
 561 Mohammad Taher Pilehvar (eds.), *Proceedings of the 63rd Annual Meeting of the Association
 562 for Computational Linguistics (Volume 1: Long Papers)*, pp. 18003–18023, Vienna, Austria, July
 563 2025b. Association for Computational Linguistics. ISBN 979-8-89176-251-0. URL <https://aclanthology.org/2025.acl-long.881/>.

564 Cursor. The ai code editor, 2025. URL <https://cursor.com/en>.

565 DeepSeek-AI. Deepseek-v3 technical report, 2025a. URL <https://arxiv.org/abs/2412.19437>.

566 DeepSeek-AI. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning,
 567 2025b. URL <https://arxiv.org/abs/2501.12948>.

568 DeepSeek-AI, Qihao Zhu, Daya Guo, Zhihong Shao, Dejian Yang, Peiyi Wang, Runxin Xu, Y. Wu,
 569 Yukun Li, Huazuo Gao, Shirong Ma, Wangding Zeng, Xiao Bi, Zihui Gu, Hanwei Xu, Damai
 570 Dai, Kai Dong, Liyue Zhang, Yishi Piao, Zhibin Gou, Zhenda Xie, Zhewen Hao, Bingxuan Wang,
 571 Junxiao Song, Deli Chen, Xin Xie, Kang Guan, Yuxiang You, Aixin Liu, Qiushi Du, Wenjun Gao,
 572 Xuan Lu, Qinyu Chen, Yaohui Wang, Chengqi Deng, Jiashi Li, Chenggang Zhao, Chong Ruan,
 573 Fuli Luo, and Wenfeng Liang. Deepseek-coder-v2: Breaking the barrier of closed-source models
 574 in code intelligence, 2024. URL <https://arxiv.org/abs/2406.11931>.

575 Gemini. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and
 576 next generation agentic capabilities, 2025. URL <https://arxiv.org/abs/2507.06261>.

577 Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen,
 578 Xiao Bi, Y. Wu, Y. K. Li, Fuli Luo, Yingfei Xiong, and Wenfeng Liang. Deepseek-coder:
 579 When the large language model meets programming – the rise of code intelligence, 2024. URL
 580 <https://arxiv.org/abs/2401.14196>.

581 Xinyi He, Qian Liu, Mingzhe Du, Lin Yan, Zhiping Fan, Yiming Huang, Zejian Yuan, and Zejun Ma.
 582 Swe-perf: Can language models optimize code performance on real-world repositories?, 2025.
 583 URL <https://arxiv.org/abs/2507.12415>.

584 Dan Hendrycks, Steven Basart, Saurav Kadavath, Mantas Mazeika, Akul Arora, Ethan Guo, Collin
 585 Burns, Samir Puranik, Horace He, Dawn Song, et al. Measuring coding challenge competence
 586 with apps. *arXiv preprint arXiv:2105.09938*, 2021.

594 Siming Huang, Tianhao Cheng, J. K. Liu, Jiaran Hao, Liuyihan Song, Yang Xu, J. Yang, Jiaheng Liu,
 595 Chenchen Zhang, Linzheng Chai, Ruiyuan Yuan, Zhaoxiang Zhang, Jie Fu, Qian Liu, Ge Zhang,
 596 Zili Wang, Yuan Qi, Yinghui Xu, and Wei Chu. Opencoder: The open cookbook for top-tier code
 597 large language models, 2025. URL <https://arxiv.org/abs/2411.04905>.

598 Binyuan Hui, Jian Yang, Zeyu Cui, Jiaxi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun
 599 Zhang, Bowen Yu, Keming Lu, Kai Dang, Yang Fan, Yichang Zhang, An Yang, Rui Men, Fei
 600 Huang, Bo Zheng, Yibo Miao, Shanghaoran Quan, Yunlong Feng, Xingzhang Ren, Xuancheng
 601 Ren, Jingren Zhou, and Junyang Lin. Qwen2.5-coder technical report, 2024a. URL <https://arxiv.org/abs/2409.12186>.

602 Binyuan Hui, Jian Yang, Zeyu Cui, Jiaxi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang,
 603 Bowen Yu, Keming Lu, et al. Qwen2. 5-coder technical report. [arXiv preprint arXiv:2409.12186](https://arxiv.org/abs/2409.12186),
 604 2024b.

605 Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando
 606 Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free
 607 evaluation of large language models for code. In The Thirteenth International Conference on
 608 Learning Representations, 2025. URL <https://openreview.net/forum?id=chfJJYC3iL>.

609 Juyong Jiang, Fan Wang, Jiasi Shen, Sungju Kim, and Sunghun Kim. A survey on large language
 610 models for code generation, 2024. URL <https://arxiv.org/abs/2406.00515>.

611 Carlos E. Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik
 612 Narasimhan. Swe-bench: Can language models resolve real-world github issues?, 2024. URL
 613 <https://arxiv.org/abs/2310.06770>.

614 Kimi-Team. Kimi k2: Open agentic intelligence, 2025. URL <https://arxiv.org/abs/2507.20534>.

615 Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom
 616 Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, Thomas Hubert, Peter Choy, Cyprien
 617 de Masson d'Autume, Igor Babuschkin, Xinyun Chen, Po-Sen Huang, Johannes Welbl, Sven
 618 Gowal, Alexey Cherepanov, James Molloy, Daniel J. Mankowitz, Esme Sutherland Robson,
 619 Pushmeet Kohli, Nando de Freitas, Koray Kavukcuoglu, and Oriol Vinyals. Competition-level code
 620 generation with alphacode. *Science*, 378(6624):1092–1097, December 2022. ISSN 1095-9203. doi:
 621 [10.1126/science.abq1158](https://doi.org/10.1126/science.abq1158). URL <http://dx.doi.org/10.1126/science.abq1158>.

622 Anton Lozhkov, Raymond Li, Loubna Ben Allal, Federico Cassano, Joel Lamy-Poirier, Nouamane
 623 Tazi, Ao Tang, Dmytro Pykhtar, Jiawei Liu, Yuxiang Wei, Tianyang Liu, Max Tian, Denis Kocetkov,
 624 Arthur Zucker, Younes Belkada, Zijian Wang, Qian Liu, Dmitry Abulkhanov, Indraneil Paul,
 625 Zhuang Li, Wen-Ding Li, Megan Risdal, Jia Li, Jian Zhu, Terry Yue Zhuo, Evgenii Zheltonozhskii,
 626 Nii Osae Osae Dade, Wenhao Yu, Lucas Krauß, Naman Jain, Yixuan Su, Xuanli He, Manan
 627 Dey, Edoardo Abati, Yekun Chai, Niklas Muennighoff, Xiangru Tang, Muhtasham Oblokulov,
 628 Christopher Akiki, Marc Marone, Chenghao Mou, Mayank Mishra, Alex Gu, Binyuan Hui, Tri
 629 Dao, Armel Zebaze, Olivier Dehaene, Nicolas Patry, Canwen Xu, Julian McAuley, Han Hu, Torsten
 630 Scholak, Sebastien Paquet, Jennifer Robinson, Carolyn Jane Anderson, Nicolas Chapados, Mostafa
 631 Patwary, Nima Tajbakhsh, Yacine Jernite, Carlos Muñoz Ferrandis, Lingming Zhang, Sean Hughes,
 632 Thomas Wolf, Arjun Guha, Leandro von Werra, and Harm de Vries. Starcoder 2 and the stack v2:
 633 The next generation, 2024.

634 Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xiubo Geng, Wenxiang Hu, Chongyang Tao, Jing
 635 Ma, Qingwei Lin, and Dixin Jiang. Wizardcoder: Empowering code large language models with
 636 evol-instruct. In The Twelfth International Conference on Learning Representations, 2024. URL
 637 <https://openreview.net/forum?id=UnUwSIGK5W>.

638 OpenAI. Hello GPT-4o, 2024. URL <https://openai.com/index/hello-gpt-4o/>.

639 OpenAI. Introducing gpt-4.1 in the api, 2025a. URL <https://openai.com/index/gpt-4-1/>.

640 OpenAI. Gpt-5, 2025b. URL <https://openai.com/index/introducing-gpt-5/>.

648 OpenAI. Introducing openai o3 and o4-mini, 2025. URL <https://openai.com/index/introducing-o3-and-o4-mini/>.

649

650

651 Arjun Panickssery, Samuel R. Bowman, and Shi Feng. LLM evaluators recognize and favor their own

652 generations. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*,

653 2024. URL <https://openreview.net/forum?id=4NJBV6Wp0h>.

654

655 Qiwei Peng, Yekun Chai, and Xuhong Li. Humaneval-xl: A multilingual code generation benchmark

656 for cross-lingual natural language generalization, 2024. URL <https://arxiv.org/abs/2402.16694>.

657

658 Qwen. Qwen3-coder: Agentic coding in the world, 2025. URL <https://qwenlm.github.io/blog/qwen3-coder/>.

659

660 Qwen. Qwen3-next: Towards ultimate training & inference efficiency, 2025. URL <https://qwen.ai>.

661

662 Muhammad Shihab Rashid, Christian Bock, Yuan Zhuang, Alexander Buchholz, Tim Esler, Simon

663 Valentin, Luca Franceschi, Martin Wistuba, Prabhu Teja Sivaprasad, Woo Jung Kim, Anoop

664 Deoras, Giovanni Zappella, and Laurent Callot. Swe-polybench: A multi-language benchmark for

665 repository level evaluation of coding agents, 2025. URL <https://arxiv.org/abs/2504.08703>.

666

667 Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi

668 Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, et al. Code llama: Open foundation models for

669 code. *arXiv preprint arXiv:2308.12950*, 2023.

670

671 Seed. Introduction to techniques used in seed1.6, 2025. URL https://seed/bytedance.com/en/seed1_6.

672

673 ByteDance Seed, Yuyu Zhang, Jing Su, Yifan Sun, Chenguang Xi, Xia Xiao, Shen Zheng, Anxiang

674 Zhang, Kaibo Liu, Daoguang Zan, Tao Sun, Jinhua Zhu, Shulin Xin, Dong Huang, Yetao Bai,

675 Lixin Dong, Chao Li, Jianchong Chen, Hanzhi Zhou, Yifan Huang, Guanghan Ning, Xierui Song,

676 Jiaze Chen, Siyao Liu, Kai Shen, Liang Xiang, and Yonghui Wu. Seed-coder: Let the code model

677 curate data for itself, 2025. URL <https://arxiv.org/abs/2506.03524>.

678

679 Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang,

680 Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. Deepseekmath: Pushing the limits of

681 mathematical reasoning in open language models, 2024. URL <https://arxiv.org/abs/2402.03300>.

682

683 Tencent. Hunyuan-turbos: Advancing large language models through mamba-transformer synergy

684 and adaptive chain-of-thought, 2025. URL <https://arxiv.org/abs/2505.15431>.

685

686 Kaixin Wang, Tianlin Li, Xiaoyu Zhang, Chong Wang, Weisong Sun, Yang Liu, and Bin Shi. Software

687 development life cycle perspective: A survey of benchmarks for code large language models and

688 agents, 2025a. URL <https://arxiv.org/abs/2505.05283>.

689

690 Zhexu Wang, Yiping Liu, Yejie Wang, Wenyang He, Bofei Gao, Muxi Diao, Yanxu Chen, Kelin Fu,

691 Flood Sung, Zhilin Yang, et al. Objbench: A competition level code benchmark for large language

692 models. *arXiv preprint arXiv:2506.16395*, 2025b.

693

694 Yuxiang Wei, Federico Cassano, Jiawei Liu, Yifeng Ding, Naman Jain, Zachary Mueller, Harm

695 de Vries, Leandro von Werra, Arjun Guha, and Lingming Zhang. Selfcodealign: Self-alignment

696 for code generation. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak,

697 and C. Zhang (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 62787–

698 62874. Curran Associates, Inc., 2024a. URL https://proceedings.neurips.cc/paper_files/paper/2024/file/72da102da91a8042a0b2aa968429a9f9-Paper-Conference.pdf.

699

700 Yuxiang Wei, Zhe Wang, Jiawei Liu, Yifeng Ding, and Lingming Zhang. Magicoder: Empower-

701 ing code generation with OSS-instruct. In *Proceedings of the 41st International Conference on*

702 *Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 52632–52657.

703 PMLR, 21–27 Jul 2024b. URL <https://proceedings.mlr.press/v235/wei24h.html>.

702 Yutong Wu, Di Huang, Wenxuan Shi, Wei Wang, Lingzhe Gao, Shihao Liu, Ziyuan Nan, Kaizhao
 703 Yuan, Rui Zhang, Xishan Zhang, Zidong Du, Qi Guo, Yewen Pu, Dawei Yin, Xing Hu, and Yunji
 704 Chen. Inversecoder: Self-improving instruction-tuned code llms with inverse-instruct, 2024. URL
 705 <https://arxiv.org/abs/2407.05700>.

706 Zhangchen Xu, Yang Liu, Yueqin Yin, Mingyuan Zhou, and Radha Poovendran. Kodcode: A diverse,
 707 challenging, and verifiable synthetic dataset for coding, 2025. URL <https://arxiv.org/abs/2503.02951>.

708

710 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang
 711 Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. [arXiv preprint arXiv:2505.09388](https://arxiv.org/abs/2505.09388),
 712 2025.

713 Zhaojian Yu, Xin Zhang, Ning Shang, Yangyu Huang, Can Xu, Yishujie Zhao, Wenxiang Hu, and
 714 Qiu Feng Yin. Wavecoder: Widespread and versatile enhancement for code large language models
 715 by instruction tuning, 2024. URL <https://arxiv.org/abs/2312.14187>.

716

717 Daoguang Zan, Zhirong Huang, Wei Liu, Hanwu Chen, Linhao Zhang, Shulin Xin, Lu Chen, Qi Liu,
 718 Xiaojian Zhong, Aoyan Li, Siyao Liu, Yongsheng Xiao, Liangqiang Chen, Yuyu Zhang, Jing Su,
 719 Tianyu Liu, Rui Long, Kai Shen, and Liang Xiang. Multi-swe-bench: A multilingual benchmark
 720 for issue resolving, 2025. URL <https://arxiv.org/abs/2504.02605>.

721

722 Alexander Zhang, Marcus Dong, Jiaheng Liu, Wei Zhang, Yejie Wang, Jian Yang, Ge Zhang, Tianyu
 723 Liu, Zhongyuan Peng, Yingshui Tan, Yuanxing Zhang, Zhexu Wang, Weixun Wang, Yancheng
 724 He, Ken Deng, Wangchunshu Zhou, Wenhao Huang, and Zhaoxiang Zhang. Codecriticbench: A
 725 holistic code critique benchmark for large language models, 2025a. URL <https://arxiv.org/abs/2502.16614>.

726

727 Chenchen Zhang, Yuhang Li, Can Xu, Jiaheng Liu, Ao Liu, Shihui Hu, Dengpeng Wu, Guanhua
 728 Huang, Kejiao Li, Qi Yi, Ruibin Xiong, Haotian Zhu, Yuanxing Zhang, Yuhao Jiang, Yue Zhang,
 729 Zenan Xu, Bohui Zhai, Guoxiang He, Hebin Li, Jie Zhao, Le Zhang, Lingyun Tan, Pengyu Guo,
 730 Xianshu Pang, Yang Ruan, Zhifeng Zhang, Zhonghu Wang, Ziyan Xu, Zuopu Yin, Wiggin Zhou,
 731 Chayse Zhou, and Fengzong Lian. Artifactsbench: Bridging the visual-interactive gap in llm code
 732 generation evaluation, 2025b. URL <https://arxiv.org/abs/2507.04952>.

733

734 Shudan Zhang, Hanlin Zhao, Xiao Liu, Qinkai Zheng, Zehan Qi, Xiaotao Gu, Xiaohan Zhang, Yuxiao
 735 Dong, and Jie Tang. Naturalcodebench: Examining coding performance mismatch on humaneval
 736 and natural user prompts, 2024. URL <https://arxiv.org/abs/2405.04520>.

737

738 Xinyu Zhang, Changzhi Zhou, Linmei Hu, Luhao Zhang, Xiancai Chen, Haomin Fu, Yang Yang, and
 739 Mengdi Zhang. Scoder: Iterative self-distillation for bootstrapping small-scale data synthesizers to
 740 empower code llms, 2025c. URL <https://arxiv.org/abs/2509.07858>.

741

742 Tianyu Zheng, Ge Zhang, Tianhao Shen, Xueling Liu, Bill Yuchen Lin, Jie Fu, Wenhui Chen, and
 743 Xiang Yue. OpenCodeInterpreter: Integrating code generation with execution and refinement.
 744 In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), [Findings of the Association for
 745 Computational Linguistics: ACL 2024](https://aclanthology.org/2024.findings-acl.762), pp. 12834–12859, Bangkok, Thailand, August 2024.
 746 Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.762. URL
 747 <https://aclanthology.org/2024.findings-acl.762/>.

748

749 Zihan Zheng, Zerui Cheng, Zeyu Shen, Shang Zhou, Kaiyuan Liu, Hansen He, Dongruixuan Li,
 750 Stanley Wei, Hangyi Hao, Jianzhu Yao, Peiyao Sheng, Zixuan Wang, Wenhao Chai, Aleksandra
 751 Korolova, Peter Henderson, Sanjeev Arora, Pramod Viswanath, Jingbo Shang, and Saining Xie.
 752 Livecodebench pro: How do olympiad medalists judge llms in competitive programming?, 2025.
 753 URL <https://arxiv.org/abs/2506.11928>.

754

755 Zhipu. Gilm-4.5: Reasoning, coding, and agentic abilities, 2025. URL <https://z.ai/blog/gilm-4.5>.

756

757 Changzhi Zhou, Xinyu Zhang, Dandan Song, Xiancai Chen, Wanli Gu, Huipeng Ma, Yuhang Tian,
 758 Mengdi Zhang, and Linmei Hu. Refinecoder: Iterative improving of large language models via
 759 adaptive critique refinement for code generation, 2025. URL <https://arxiv.org/abs/2502.09183>.

756 Qihao Zhu, Daya Guo, Zhihong Shao, Dejian Yang, Peiyi Wang, Runxin Xu, Y Wu, Yukun Li,
757 Huazuo Gao, Shirong Ma, et al. Deepseek-coder-v2: Breaking the barrier of closed-source models
758 in code intelligence. [arXiv preprint arXiv:2406.11931](https://arxiv.org/abs/2406.11931), 2024.

759
760 Yaoming Zhu, Junxin Wang, Yiyang Li, Lin Qiu, ZongYu Wang, Jun Xu, Xuezhi Cao, Yuhuai Wei,
761 Mingshi Wang, Xunliang Cai, and Rong Ma. Oibench: Benchmarking strong reasoning models
762 with olympiad in informatics, 2025. URL <https://arxiv.org/abs/2506.10481>.

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

810
811
812 Table 7: Comparison of Accuracy, Upper Bound, and Model Performance.
813
814
815
816

| | Average | Python | C++ | Java | JS | Go | Shell |
|---------------------------|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|
| Problem Accuracy | 87.6 | 83.5 | 88.0 | 86.0 | 89.0 | 86.0 | 93.3 |
| Current Upper Bound | 66.9 _{△20.7} | 61.7 _{△21.8} | 71.5 _{△16.5} | 76.1 _{△9.9} | 58.2 _{△30.8} | 65.4 _{△20.6} | 68.6 _{△24.7} |
| Claude Opus 4 (Reasoning) | 44.6 _{△43.0} | 40.3 _{△43.2} | 44.1 _{△43.9} | 55.9 _{△30.1} | 38.6 _{△50.4} | 37.2 _{△48.8} | 51.6 _{△41.7} |

817
818 A MANUAL VERIFICATION
819

820 Although our pipeline enforces quality control through specifications and an LLM-as-Critic mechanism,
821 we further validate AutoCodeBench with human annotators. We employ a Human-LLM
822 collaboration approach for data quality validation. Specifically, we design prompts in the native lan-
823 guages of the annotators and use the DeepSeek-R1-0528 to generate detailed reasoning processes
824 and checklist-based annotation results. The prompt is shown in Figure 10. During the annotation
825 process, we assume that the programming problems are completely correct. The primary task of the
826 annotators is to assess the correctness of the test functions and their alignment with the programming
827 problem, based on the LLM’s output. We allow for test cases that may not cover all boundary condi-
828 tions, focusing primarily on the correctness of the test functions rather than their comprehensiveness.
829 The annotators pay particular attention to the following aspects:

- 830 • Whether the function names, class names, variable definitions, and return types are consistent
831 with the problem description;
- 832 • Whether the test cases exhibit randomness or non-reproducibility;
- 833 • Whether the test cases contradict the logic presented in the problem statement;
- 834 • Whether there are any precision issues with the test cases;
- 835 • Whether the test functions include test cases that are not addressed in the problem description.

836 We calculate the problem accuracy rates for different programming languages (Python, C++, Java,
837 JavaScript, Go, Shell), as shown in Table 7. The results indicate that, despite the presence of some
838 noisy data, our benchmark model still demonstrates high accuracy (87.6%). Furthermore, even after
839 removing the noise, the current SOTA model shows significant room for improvement ($\Delta 43.0$),
840 further validating the high difficulty level of our benchmark. The performance of Claude Opus 4
841 (Reasoning) in these six languages are only 44.6($\Delta 43.0$), highlighting the significant potential for
842 improvement. Besides, we find that, compared to logic errors in the problem description and errors
843 in the test functions, the most frequently occurring issue is **incomplete problem descriptions**. For
844 example, some test functions reference class or function names that are essential but not explicitly
845 mentioned in the problem description, or they require natural language outputs for edge cases that are
846 not explicitly specified in the problem statement, leading to mismatches between the generated code
847 and the test functions. Interestingly, we observe similar issues in manually annotated benchmarks,
848 highlighting the significant challenge of creating comprehensive and accurate programming problems
849 for annotators.

850
851 B SETUP OF BENCHMARK COMPARISONS
852853 B.1 MULTI-LOGIC
854

855 AutoCodeBench contains tasks that demand executing multiple functionalities, such as implement-
856 ing both an addition and a multiplication function simultaneously. In contrast, the tasks in other
857 benchmarks are almost exclusively focused on implementing a single core functionality.

858 B.2 DIFFICULTY
859

860 We rate the difficulty of each benchmark based on the performance of DeepSeek-V3-0324.
861 Specifically, benchmarks with pass@1 below 40 are assigned five stars; those between 40–50 receive
862 four stars (e.g., LiveCodeBench-v6: 46.9, AutoCodeBench: 48.1); between 50–60 receive three

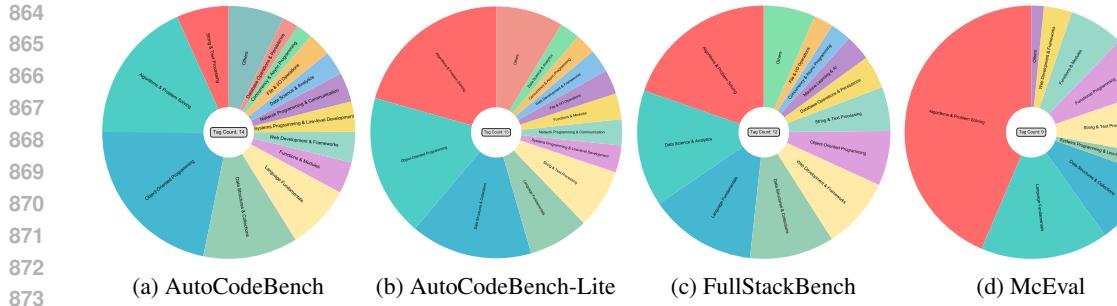


Figure 5: Category Distribution of Different Benchmarks.

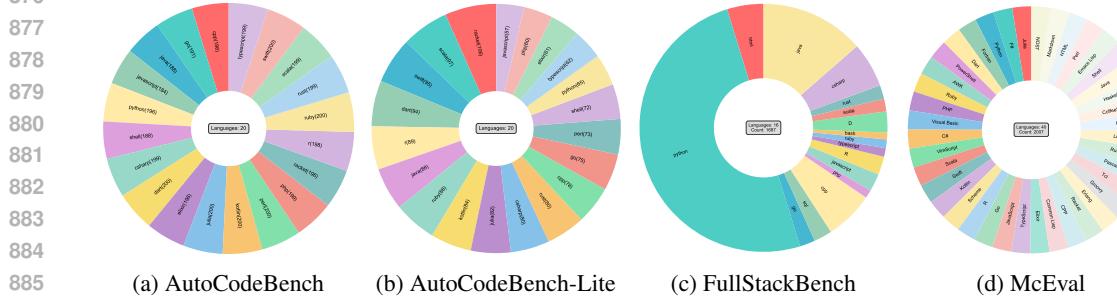


Figure 6: Language Distribution of Different Benchmarks.

stars; between 60–80 receive two stars (e.g., FullStackBench: 67.0, McEval: 72.3); and those above 80 receive one star. HumanEval and MBPP are excluded from evaluation due to extensive data leakage and its overly simple problems. For reference, DeepSeek-V3 (October 2024) already achieves a pass@1 of 91.5 in HumanEval.

B.3 CATEGORY DISTRIBUTION AND LANGUAGE DISTRIBUTION

We prompt Claude Sonnet 4 to generate 20 language-agnostic category labels for classification:

- **Core Programming Concepts:** *Language Fundamentals, Functions & Modules, Object-Oriented Programming, Functional Programming, Memory Management & Performance, Error Handling & Debugging*
- **Data and Algorithms:** *Data Structures & Collections, Algorithms & Problem Solving, String & Text Processing, File & I/O Operations, Concurrency & Async Programming*
- **Application Domains:** *Network Programming & Communication, Database Operations & Persistence, Web Development & Frameworks, Mobile & Cross-platform Development, Systems Programming & Low-level Development*
- **Advanced Topics and Tooling:** *Data Science & Analytics, Machine Learning & AI, Testing & Quality Assurance, Development Tools & Ecosystem*

In addition to AutoCodeBench, we conduct task tagging and language distribution analysis for AutoCodeBench-Lite, FullStackBench, and McEval. The results are presented in Figures 5 and 6. FullStackBench demonstrates comparable category diversity to AutoCodeBench(-Lite) but suffers from an imbalanced language distribution. In contrast, McEval exhibits a well-balanced multilingual distribution but lacks diversity and balance in its category coverage. Our AutoCodeBench(-Lite) achieves the most comprehensive category coverage while maintaining a balanced multilingual distribution, enabling thorough and accurate evaluation of LLMs’ multilingual code generation capabilities.

918 Table 8: Pass@1 (%) performance of different base models for 3-shot AutoCodeBench-Complete.
919

| Count | Average | Python 50 | Cpp 50 | Java 50 | JS 50 | Go 50 | Shell 50 | Csharp 50 | Dart 50 | Elixir 50 | Julia 50 | Kotlin 50 | Perl 50 | PHP 50 | Racket 50 | R 50 | Ruby 50 | Rust 50 | Scala 50 | Swift 50 | TS 50 |
|--------------------------|-------------|-----------|--------|---------|-------|-------|----------|-----------|---------|-----------|----------|-----------|---------|--------|-----------|------|---------|---------|----------|----------|-------|
| 30B+ Models | | | | | | | | | | | | | | | | | | | | | |
| DeepSeek-Coder-V2-Base | 39.0 | 24.0 | 32.0 | 40.0 | 44.0 | 34.0 | 26.0 | 64.0 | 38.0 | 52.0 | 46.0 | 56.0 | 38.0 | 36.0 | 32.0 | 26.0 | 40.0 | 26.0 | 36.0 | 42.0 | 48.0 |
| Qwen2.5-Coder-32B | 35.5 | 36.0 | 34.0 | 32.0 | 32.0 | 38.0 | 34.0 | 58.0 | 30.0 | 42.0 | 38.0 | 52.0 | 40.0 | 32.0 | 30.0 | 26.0 | 30.0 | 18.0 | 30.0 | 34.0 | 44.0 |
| Qwen2.5-72B | 35.9 | 32.0 | 22.0 | 38.0 | 40.0 | 22.0 | 34.0 | 62.0 | 22.0 | 42.0 | 38.0 | 46.0 | 42.0 | 28.0 | 26.0 | 38.0 | 28.0 | 28.0 | 30.0 | 30.0 | 54.0 |
| ~8B Models | | | | | | | | | | | | | | | | | | | | | |
| Seed-Coder-8B-Base | 31.6 | 26.0 | 22.0 | 40.0 | 30.0 | 32.0 | 12.0 | 54.0 | 24.0 | 48.0 | 30.0 | 48.0 | 28.0 | 36.0 | 22.0 | 26.0 | 32.0 | 18.0 | 20.0 | 36.0 | 48.0 |
| OpenCoder-8B-Base | 26.1 | 22.0 | 6.0 | 28.0 | 34.0 | 30.0 | 24.0 | 52.0 | 10.0 | 32.0 | 28.0 | 26.0 | 24.0 | 20.0 | 20.0 | 28.0 | 14.0 | 26.0 | 14.0 | 42.0 | |
| Qwen2.5-Coder-7B | 24.6 | 20.0 | 10.0 | 22.0 | 28.0 | 24.0 | 14.0 | 46.0 | 8.0 | 46.0 | 32.0 | 42.0 | 30.0 | 30.0 | 14.0 | 20.0 | 18.0 | 16.0 | 24.0 | 14.0 | 34.0 |
| DeepSeek-Coder-6.7B-Base | 22.9 | 20.0 | 14.0 | 26.0 | 34.0 | 18.0 | 18.0 | 50.0 | 8.0 | 44.0 | 28.0 | 18.0 | 12.0 | 14.0 | 34.0 | 6.0 | 10.0 | 4.0 | 42.0 | | |
| Qwen3-8B-Base | 22.6 | 20.0 | 14.0 | 18.0 | 34.0 | 20.0 | 12.0 | 50.0 | 6.0 | 34.0 | 26.0 | 24.0 | 32.0 | 30.0 | 8.0 | 20.0 | 30.0 | 8.0 | 14.0 | 16.0 | 36.0 |

926 Table 9: Pass@1 (%) performance of different models for AutoCodeBench-Lite.
927

| Count | Average | Python | Cpp | Java | JS | Go | Shell | Csharp | Dart | Elixir | Julia | Kotlin | Perl | PHP | Racket | R | Ruby | Rust | Scala | Swift | TS | | |
|---|----------------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|------|------|------|
| Count | Current Upper Bound | 65 | 78 | 88 | 57 | 75 | 72 | 80 | 94 | 61 | 82 | 84 | 73 | 60 | 106 | 89 | 88 | 80 | 97 | 95 | 62 | | |
| Count | Current Upper Bound | 100.0 | | | |
| Proprietary Models and 200B+ Open-Source Models | | | | | | | | | | | | | | | | | | | | | | | |
| Claude Opus 4.1 (20250805) | 69.9 | 61.5 | 62.8 | 65.9 | 68.4 | 68.0 | 70.8 | 90.0 | 69.1 | 78.7 | 73.2 | 73.8 | 67.1 | 55.0 | 80.2 | 68.5 | 72.7 | 67.5 | 56.7 | 67.4 | 77.4 | | |
| Claude Sonnet 4 (20250514) | 62.0 | 53.9 | 60.3 | 60.2 | 59.0 | 64.0 | 65.3 | 78.8 | 61.7 | 63.9 | 53.7 | 66.7 | 64.4 | 61.7 | 76.4 | 59.6 | 62.5 | 57.5 | 56.7 | 52.6 | 66.1 | | |
| Claude Opus 4.1 (20250805) | 63.8 | 61.5 | 61.4 | 63.2 | 61.3 | 62.5 | 83.7 | 68.3 | 68.3 | 73.8 | 61.6 | 48.3 | 68.9 | 60.7 | 61.4 | 65.0 | 56.7 | 60.0 | 72.6 | | | | |
| Claude Sonnet 4 (20250514) | 59.8 | 44.6 | 61.5 | 60.2 | 57.9 | 56.0 | 65.3 | 83.8 | 56.4 | 62.3 | 53.7 | 69.1 | 65.8 | 55.0 | 71.7 | 50.6 | 56.8 | 56.3 | 49.5 | 49.5 | 71.0 | | |
| GPT-3 (20250807) | 67.0 | 66.2 | 64.1 | 52.3 | 73.7 | 66.7 | 73.6 | 80.0 | 67.0 | 72.1 | 61.0 | 64.3 | 62.5 | 65.0 | 61.9 | 71.9 | 70.5 | 68.7 | 58.8 | 68.4 | 80.6 | | |
| o3-mlm (20250416) | 63.5 | 61.9 | 59.1 | 60.2 | 64.0 | 53.3 | 68.1 | 70.0 | 67.0 | 64.0 | 55.0 | 61.0 | 60.0 | 66.0 | 59.0 | 58.0 | 58.0 | 57.0 | 58.0 | 57.0 | 77.4 | | |
| o3-mlm (2025-04-16) | 60.5 | 63.1 | 55.1 | 54.6 | 57.9 | 53.5 | 63.2 | 66.0 | 53.8 | 56.1 | 67.9 | 64.4 | 55.3 | 46.2 | 50.6 | 70.5 | 68.8 | 58.8 | 63.2 | 72.6 | | | |
| GPT-4 (2025-04-14) | 56.9 | 57.7 | 48.9 | 52.6 | 38.7 | 47.2 | 83.8 | 54.3 | 59.0 | 63.4 | 69.1 | 63.0 | 41.7 | 72.6 | 59.6 | 73.9 | 65.0 | 60.8 | 54.7 | 66.1 | | | |
| Grok-4 | 63.0 | 60.9 | 64.1 | 55.7 | 56.1 | 61.3 | 68.1 | 82.5 | 55.3 | 59.0 | 57.4 | 62.3 | 64.6 | 71.4 | 58.9 | 51.7 | 53.8 | 60.7 | 63.6 | 37.5 | 53.6 | 57.9 | 61.3 |
| GemmIn2.5 Pro | 59.1 | 56.9 | 59.7 | 59.1 | 54.4 | 50.7 | 62.5 | 75.0 | 57.4 | 62.3 | 51.2 | 63.1 | 60.3 | 50.0 | 55.7 | 57.3 | 60.2 | 61.2 | 48.5 | 55.8 | 57.9 | 58.1 | |
| GemmIn2.5 Flash | 52.9 | 49.2 | 56.4 | 54.6 | 50.9 | 45.4 | 45.8 | 75.0 | 56.4 | 45.0 | 51.2 | 63.1 | 54.1 | 40.0 | 34.0 | 59.6 | 58.0 | 58.8 | 46.4 | 51.6 | 61.3 | | |
| DeepSeek-V3.1 | 57.5 | 53.8 | 61.5 | 60.2 | 52.6 | 37.3 | 47.2 | 80.0 | 57.4 | 67.0 | 61.0 | 63.1 | 60.3 | 50.0 | 55.7 | 57.3 | 60.2 | 61.2 | 48.5 | 55.8 | 58.1 | | |
| DeepSeek-V3.1 | 52.2 | 44.6 | 51.3 | 64.8 | 43.9 | 33.3 | 43.1 | 72.5 | 43.6 | 59.0 | 56.1 | 57.1 | 61.6 | 50.0 | 51.9 | 55.1 | 51.1 | 50.0 | 42.3 | 48.4 | 66.1 | | |
| DeepSeek-Coder-V2-Instruct | 40.5 | 29.2 | 37.2 | 37.5 | 26.5 | 42.7 | 36.1 | 66.1 | 30.9 | 36.1 | 34.5 | 42.7 | 48.0 | 40.0 | 40.6 | 46.5 | 47.0 | 30.0 | 39.0 | 46.8 | | | |
| Human-Chat (2025-04-15) | 50.3 | 47.1 | 54.6 | 54.6 | 45.6 | 47.9 | 53.1 | 75.0 | 56.9 | 58.9 | 56.1 | 56.9 | 53.4 | 49.5 | 42.5 | 50.0 | 52.3 | 42.5 | 42.2 | 43.2 | 64.5 | | |
| GLM-4.5-Enable | 55.0 | 56.1 | 54.7 | 56.6 | 49.0 | 54.2 | 59.4 | 76.3 | 39.5 | 59.3 | 52.4 | 63.1 | 61.6 | 50.0 | 49.5 | 57.8 | 58.0 | 50.7 | 42.9 | 52.1 | 61.3 | | |
| KimJ-K2-0905-preview | 53.7 | 46.2 | 44.9 | 51.1 | 52.6 | 49.3 | 51.4 | 73.7 | 58.5 | 65.6 | 51.2 | 56.0 | 56.2 | 50.0 | 48.3 | 48.9 | 45.0 | 48.5 | 51.6 | 66.1 | | | |
| ERNIE-XL-Turbo-32K | 44.4 | 50.8 | 26.3 | 37.9 | 49.1 | 48.0 | 32.4 | 46.8 | 28.7 | 53.5 | 63.0 | 42.9 | 53.4 | 38.3 | 45.3 | 49.4 | 50.0 | 31.2 | 42.7 | 45.3 | 61.3 | | |
| Qwen3-23B-A22B-Thinking-2507 | 57.3 | 55.4 | 52.6 | 54.6 | 64.9 | 53.3 | 62.5 | 80.0 | 50.0 | 68.9 | 56.1 | 47.6 | 50.7 | 46.7 | 65.1 | 62.5 | 60.7 | 57.5 | 44.3 | 51.6 | 66.1 | | |
| Qwen3-Coder-480B-A35B-Instruct | 51.5 | 52.3 | 44.9 | 59.1 | 45.6 | 42.7 | 52.8 | 68.8 | 37.2 | 59.0 | 56.1 | 56.0 | 57.5 | 46.7 | 57.6 | 43.8 | 53.4 | 38.4 | 40.2 | 53.7 | 67.7 | | |
| Qwen3-23B-A22B-Instruct-2507 | 49.8 | 43.1 | 47.4 | 55.7 | 43.9 | 44.0 | 54.2 | 73.8 | 38.3 | 45.9 | 41.5 | 52.5 | 59.0 | 60.7 | 54.8 | 43.3 | 51.7 | 47.7 | 40.4 | 37.1 | 48.4 | 66.1 | |
| Seed1.6-Thinking-20715 | 53.9 | 56.9 | 60.3 | 58.0 | 50.9 | 58.7 | 48.6 | 73.8 | 38.3 | 45.9 | 41.5 | 50.0 | 57.1 | 40.0 | 51.9 | 55.1 | 60.2 | 46.3 | 47.4 | 47.4 | 62.9 | | |
| Seed1.6-disabled (250615) | 53.2 | 52.2 | 51.1 | 51.1 | 45.6 | 60.0 | 55.6 | 62.5 | 43.6 | 52.3 | 57.3 | 44.8 | 57.5 | 41.7 | 52.8 | 60.0 | 58.0 | 43.8 | 41.5 | 46.1 | 66.1 | | |
| Seed1.6-disabled (250615) | 48.8 | 41.5 | 51.3 | 50.0 | 47.4 | 50.7 | 47.1 | 78.8 | 35.0 | 52.4 | 56.2 | 33.3 | 50.0 | 46.0 | 41.7 | 48.8 | 29.9 | 40.0 | 71.0 | | | | |
| 20B+ Open-source Models | | | | | | | | | | | | | | | | | | | | | | | |
| GLM-4.5-Airable | 46.2 | 53.9 | 44.9 | 45.5 | 45.6 | 50.7 | 46.8 | 70.0 | 34.8 | 50.8 | 42.7 | 44.1 | 46.2 | 45.0 | 41.5 | 48.3 | 47.7 | 33.8 | 29.9 | 40.0 | 67.7 | | |
| Qwen3-Next-80B-A3B-Thinking | 46.3 | 50.8 | 41.0 | 48.9 | 59.6 | 42.7 | 47.6 | 76.2 | 34.5 | 42.5 | 47.6 | 39.3 | 60.3 | 41.7 | 29.2 | 53.9 | 54.5 | 38.7 | 38.1 | 44.7 | 61.3 | | |
| Qwen3-Next-80B-A3B-Instruct | 42.6 | 43.1 | 35.9 | 46.6 | 42.1 | 34.7 | 47.2 | 66.2 | 35.1 | 44.3 | 43.9 | 28.6 | 58.9 | 45.0 | 36.8 | 43.8 | 40.9 | 38.7 | 30.9 | 42.1 | 59.7 | | |
| Qwen3-32B | 47.6 | 50.8 | 43.6 | 46.6 | 43.9 | 46.8 | 46.8 | 72.5 | 36.2 | 55.7 | 47.6 | 44.1 | 60.3 | 45.0 | 35.9 | 52.8 | 44.3 | 43.8 | 37.1 | 46.3 | 64.5 | | |
| Qwen3-14B | 40.7 | 46.2 | 39.7 | 37.5 | 39.8 | 42.1 | 37.3 | 40.3 | 62.5 | 24.5 | 52.5 | 41.5 | 33.3 | 50.7 | 33.3 | 41.5 | 41.6 | 43.2 | 42.5 | 22.7 | 33.7 | 61.3 | |
| Qwen3-8B | 28.9 | 29.2 | 24.4 | 19.3 | 35.1 | 38.7 | 29.2 | 51.3 | 18.8 | 39.3 | 23.8 | 45.2 | 25.0 | 38.3 | 25.5 | 22.5 | 35.2 | 25.0 | 23.7 | 25.3 | 51.6 | | |
| Qwen3-1.7B | 21.4 | 26.1 | 6.4 | 31.0 | 28.1 | 20.0 | 30.8 | 38.0 | 16.0 | 22.7 | 31.2 | 19.5 | 34.5 | 42.2 | 22.0 | 29.2 | 28.4 | 25.0 | 13.4 | 22.1 | 45.2 | | |
| Qwen3-1.7B | 7.3 | 7.0 | 0.0 | 1.1 | 7.0 | 5.3 | 12.5 | 8.8 | 3.2 | 8.2 | 6.1 | 8.3 | 21.9 | 13.3 | 2.8 | 9.0 | 12.5 | 0.0 | 2.1 | 4.0 | 22.6 | | |
| Qwen2.5-Coder-32B-Instruct | 37.0 | 33.9 | 26.9 | 38.6 | 40.4 | 28.0 | 30.6 | 57.5 | 29.8 | 52.5 | 31.7 | 46.4 | 30.2 | 34.8 | 40.9 | 32.8 | 22.7 | 35.8 | 48.4 | | | | |
| Qwen2.5-Coder-7B-Instruct | 21.5 | 23.1 | 6.4 | 21.6 | 24.6 | 12.0 | 22.2 | 42.5 | 16.0 | 32.8 | 18.3 | 33.3 | 34.3 | 25.0 | 17.9 | 21.4 | 25.0 | 5.0 | 16.5 | 10.5 | 33.9 | | |
| Qwen2.5-Coder-1.5B-Instruct | 10.2 | | | | | | | | | | | | | | | | | | | | | | |

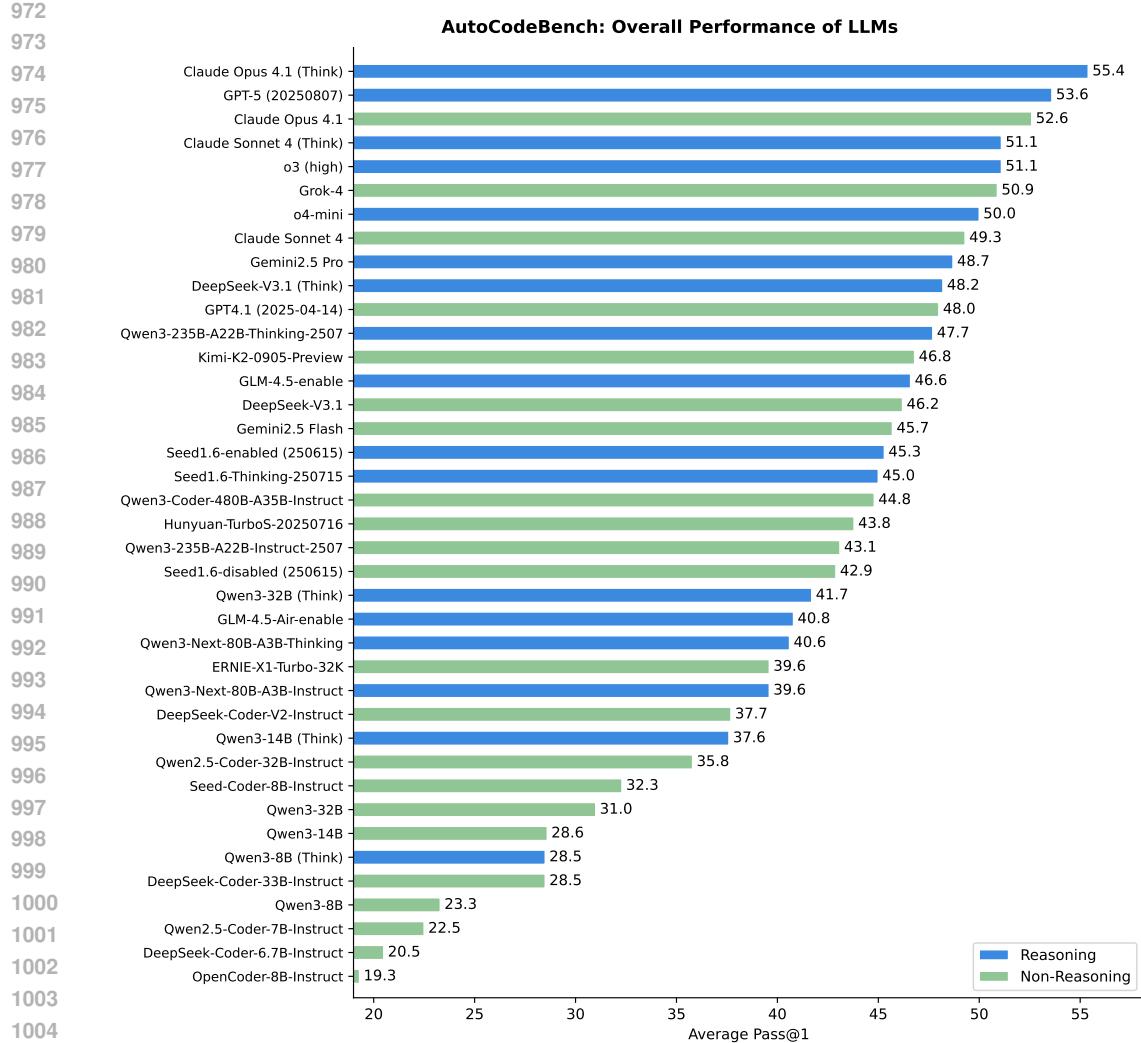


Figure 7: AutoCodeBench leaderboard showing Pass@1 performance of various LLMs.

benefits larger models. The test-time sampling scaling law (right) reveals more uniform behavior, with three models showing similar improvement rates from increased sampling during inference. These results suggest that while test-time sampling provides consistent benefits regardless of model size, reasoning capabilities scale more aggressively with model size.

G BASELINES

We evaluate a diverse set of open-source models with sizes ranging from 1.5B to 1T parameters, as well as leading proprietary models. These models are classified based on their families:

- **OpenAI**: GPT-5 (OpenAI, 2025b), o3 and o4-mini (OpenAI, 2025), and GPT4.1 (OpenAI, 2025a).
- **Claude**: Claude Opus 4.1 (Anthropic, 2025b) and Claude Sonnet 4 (Anthropic, 2025a).
- **Gemini**: Gemini 2.5 Pro and Gemini 2.5 Flash (Gemini, 2025).
- **DeepSeek**: DeepSeek-V3.1 (DeepSeek-AI, 2025a;b) and DeepSeek-Coder Series (DeepSeek-AI et al., 2024; Guo et al., 2024).

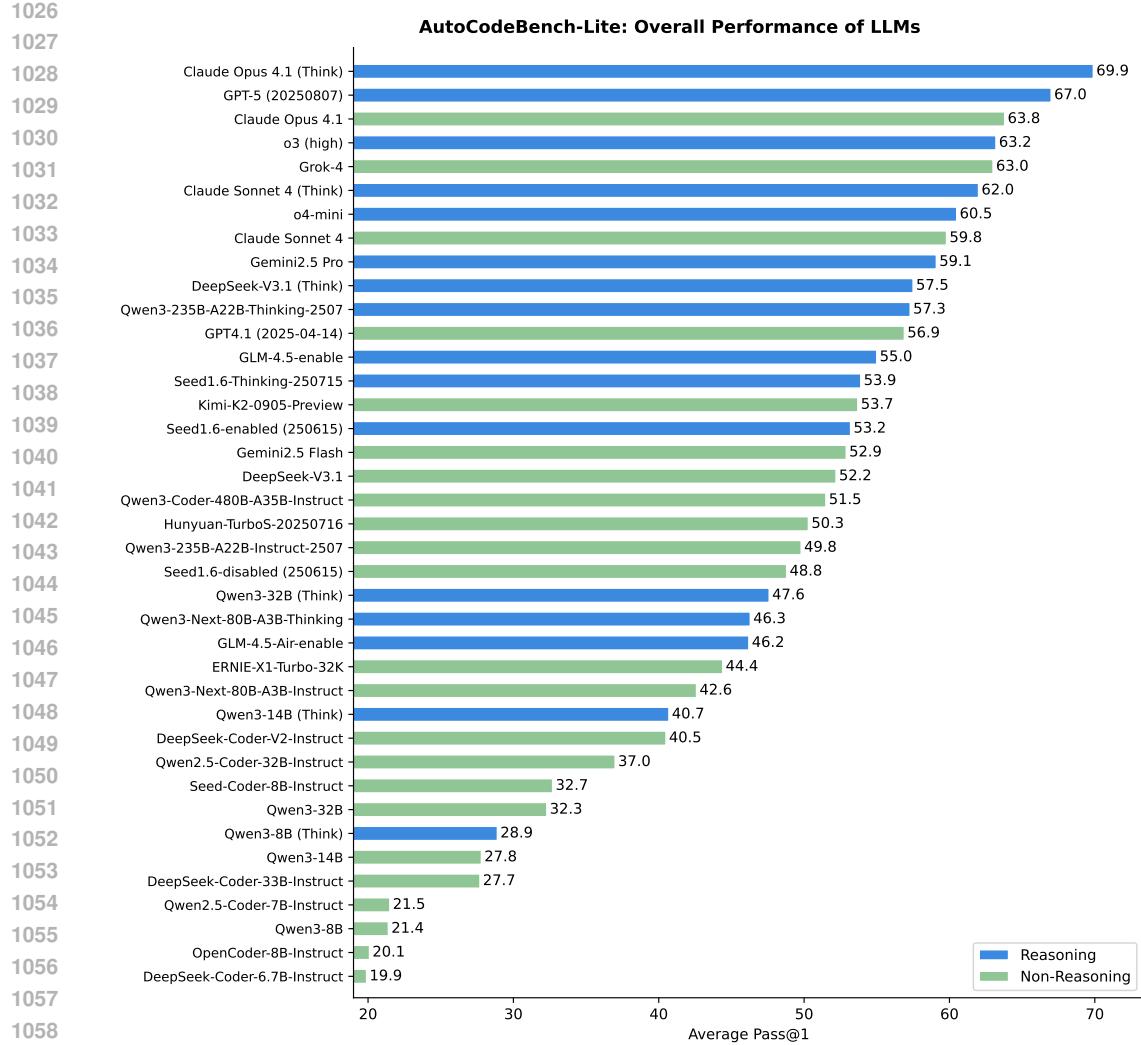


Figure 8: AutoCodeBench-Lite leaderboard showing Pass@1 performance of various LLMs.

- **Hunyuan:** Hunyuan-TurboS (Tencent, 2025).
- **Qwen:** Qwen3-Next-80B-A3B (Qwen, 2025), Qwen3-235B-A22B-Thinking-2507, Qwen3-235B-A22B-Instruct-2507, and Qwen3 Series (Yang et al., 2025), Qwen3-Coder-480B-A35B-Instruct (Qwen, 2025), Qwen2.5-Coder Series (Hui et al., 2024a).
- **Seed:** Seed1.6-Thinking (Seed, 2025), Seed1.6 (Seed, 2025) and Seed-Coder-8B (Seed et al., 2025).
- **GLM:** GLM-4.5 and GLM-4.5-Air (Zhipu, 2025).
- **Other Models:** ERNIE-X1-Turbo-32K (Baidu, 2025), Kimi-K2 (Kimi-Team, 2025), and OpenCoder-8B (Huang et al., 2025).

H HYPOTHESES ON MODEL BIAS IN THE GENERATION PROCESS

It is well-known that models exhibit inherent biases, particularly their tendency to favor their own outputs—a common phenomenon in automated data synthesis and evaluation tasks (Panickssery et al., 2024; Chen et al., 2025a). Our automated workflow is no exception to this issue. While completely eliminating such bias is challenging, we employ several mitigation strategies. Specifically,

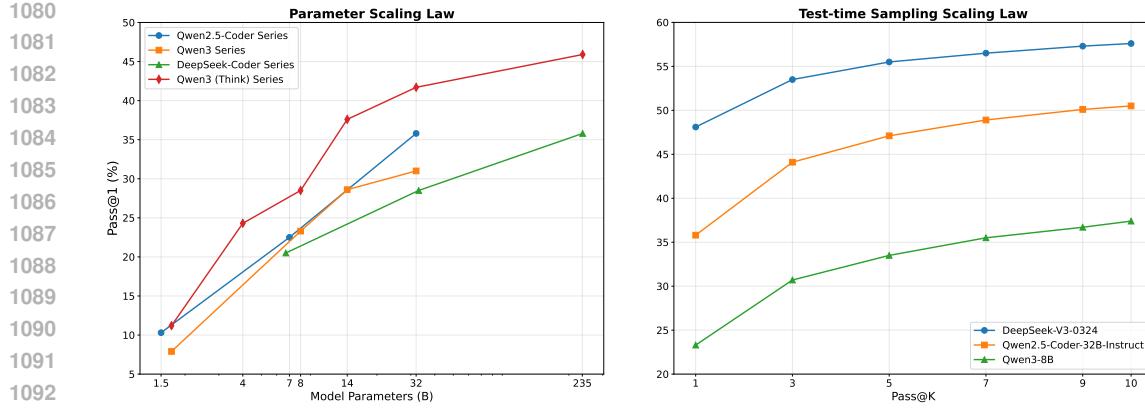


Figure 9: Scaling laws for different models.

Table 10: The average pass@1 scores and rankings of models at different stages.

| | Initial Stage (Rank) | After Simple Problem Filtering (Rank) | After Critic Filtering (Rank) |
|----------------------------|----------------------|---------------------------------------|-------------------------------|
| DeepSeek-V3-0324 | 47.1 (3) | 25.7 \downarrow 21.4 (4) | 31.6 \downarrow 5.9 (4) |
| DeepSeek-R1-0528 | 48.9 (2) | 28.7 \downarrow 20.2 (2) | 36.2 \downarrow 7.5 (2) |
| o3 | 46.4 (4) | 28.1 \downarrow 18.3 (3) | 34.9 \downarrow 6.8 (3) |
| Gemini2.5 Pro | 51.4 (1) | 31.6 \downarrow 19.8 (1) | 38.7 \downarrow 7.0 (1) |
| Qwen2.5-Coder-32B-Instruct | 39.9 (5) | 17.1 \downarrow 22.8 (5) | 22.0 \downarrow 4.9 (5) |

we intentionally only use DeepSeek series models in the workflow to prevent bias from affecting other model families. We hypothesize that using DeepSeek-V3-0324 for code generation and DeepSeek-R1-0528 for the Critic process may introduce favorable bias toward DeepSeek families. To counteract this, we employ DeepSeek-Coder-V2-Lite during the simple problem filtering phase, creating a "push-and-pull" mechanism that balances potential biases across different stages.

To quantitatively assess bias, we sampled 3,600 data points across six programming languages (Python, C++, Java, JS, Go, and Shell) and tracked performance changes at each generation stage, as shown in Table 10. The results reveal nuanced bias patterns: simple problem filtering negatively impacts smaller models (Qwen2.5-Coder-32B-Instruct) more than DeepSeek series, while the Critic process benefits DeepSeek-R1-0528 but surprisingly provides greater improvements to reasoning models (o3 and Gemini 2.5 Pro) than to DeepSeek-V3-0324. This suggests that model bias depends not only on model family but also on factors like model size and reasoning modes. Furthermore, as mutual distillation between models from different families continues, this bias becomes increasingly difficult to measure. In conclusion, we believe that our automated process may introduce a favorable bias toward the DeepSeek family of models, but the impact is minimal.

I AUTOCODEINSTRUCT EXPERIMENTAL DETAILS

We apply separate data filtering for the two Instruct models, as the performance varies across models. This results in 8684 solve-partial ($0 < pass_rate < 0.6$) prompts and 4882 solve-none ($pass_rate = 0$) prompts for Qwen2.5-Coder-7B-Instruct, 10518 solve-partial and 3294 solve-none prompts for Qwen2.5-Coder-32B-Instruct. During the two-stage GRPO training, the batch size is set to 128 and 64, respectively for 32B and 7B experiments. The learning rate is set to 1×10^{-6} and the maximum input/output lengths are 8192/8192. The 32B model is trained for 60 steps in the first stage and 80 steps in the second stage, while the 7B model is trained for 70 steps in both stages. **During SFT, the batch size is set to 64, the gradient accumulation steps are set to 2, the learning rate is 1×10^{-5} , and both the 32B and 7B models are trained for two epochs.** For evaluation, we adopt greedy decoding for all the models and the maximum output length is set to 16384.

```

1134
1135 #CONTEXT#
1136 You are a **{language}** test development expert.
1137 The system will provide you with two inputs:
1138 1. **[QUESTION]** - The problem description, outlining the functionality and constraints that the code under test
1139 should fulfill.
1140 2. **[TEST_FUNCTION]** - The test script, containing several test cases.
1141 Your task is to rigorously review the **[TEST_FUNCTION]** to ensure it truly verifies the requirements of the
1142 **[QUESTION]**, and provide structured review results for the data annotators.
1143
1144 # OBJECTIVE #
1145 Review the **[TEST_FUNCTION]** according to the following audit rules (in fixed order), providing a boolean value
1146 and a 20-40 word justification for each:
1147 1. Naming/Signature Mismatch - Function names, classes, variables, and return types in the test do not match the
1148 ones described in the question.
1149 2. Randomness/Non-Determinism - The test cases contain random factors, and the results are unstable (e.g., calling
1150 random without a seed or relying on system time).
1151 3. Incorrect Test Target - The test case verifies functionality that is not described in the problem.
1152 4. Precision Handling Issues - High precision requirements use == instead of math.isclose() or similar approximate
1153 comparisons, only using them when precision is needed.
1154 5. Exception Swallowing - The test case catches exceptions (try-except) which obscure the actual errors.
1155 6. Unexecutable Test - The default environment dependencies are okay, but the entry point (like if __name__ ==
1156 "__main__":) is missing, making the test not executable.
1157 7. Irrelevant Requirements - The test case checks for functionality not required by the problem (reasonable edge
1158 cases are exceptions).
1159 8. Other Issues - Any defects not covered by the above rules.
1160
1161 # STYLE #
1162 - Structured, concise, engineering tone
1163 - Standard JSON format; fields should use snake_case
1164 - Justification should be 20-40 words, in Chinese
1165
1166 # TONE #
1167 Professional, objective, direct
1168
1169 # AUDIENCE #
1170 Data annotators with senior development experience, who need to judge the quality of test scripts based on this.
1171
1172 #RESPONSE#
1173 Only output the following JSON structure (without Markdown code block tags):
1174
1175 ````json
1176 {
1177     "rule_results": [
1178         {
1179             "rule": "Naming/Signature Mismatch",
1180             "result": true,
1181             "reason": "Example: Function name 'add' does not match 'sum' in the problem"
1182         },
1183         // ... 8 items in total in the fixed order
1184         {
1185             "rule": "Randomness/Non-Determinism",
1186             "result": false,
1187             "reason": "Test cases rely on random factors, leading to unstable results"
1188         }
1189     ],
1190     "summary": {
1191         "overall_pass": false,
1192         "failed_rules": ["Naming/Signature Mismatch", "Randomness/Non-Determinism"],
1193         "key_points": "Random factors lead to unstable results; test target deviates from the problem's
1194         requirement"
1195     }
1196 }
1197
1198 # rule_results: List of 8 rules in fixed order (result should be true or false).
1199 # summary.overall_pass: false if any rule result is false; otherwise true.
1200 # summary.failed_rules: List of all failed rule descriptions; empty if all pass.
1201 # summary.key_points: ≤60 words summarizing the main flaws.
1202
1203 # USER INPUT #
1204 <QUESTION>
1205 {question}
1206 </QUESTION>
1207 <TEST_FUNCTION>
1208 {test_function}
1209 </TEST_FUNCTION>
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241
1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
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
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
1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457
1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
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
1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619
1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673
1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727
1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781
1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835
1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889
1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943
1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997
1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051
2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105
2106
2107
2108
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159
2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199
2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213
2214
2215
2216
2217
2218
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2250
2251
2252
2253
2254
2255
2256
2257
2258
2259
2260
2261
2262
2263
2264
2265
2266
2267
2268
2269
2270
2271
2272
2273
2274
2275
2276
2277
2278
2279
2280
2281
2282
2283
2284
2285
2286
2287
2288
2289
2290
2291
2292
2293
2294
2295
2296
2297
2298
2299
2300
2301
2302
2303
2304
2305
2306
2307
2308
2309
2310
2311
2312
2313
2314
2315
2316
2317
2318
2319
2320
2321
2322
2323
2324
2325
2326
2327
2328
2329
2330
2331
2332
2333
2334
2335
2336
2337
2338
2339
2340
2341
2342
2343
2344
2345
2346
2347
2348
2349
2350
2351
2352
2353
2354
2355
2356
2357
2358
2359
2360
2361
2362
2363
2364
2365
2366
2367
2368
2369
2370
2371
2372
2373
2374
2375
2376
2377
2378
2379
2380
2381
2382
2383
2384
2385
2386
2387
2388
2389
2390
2391
2392
2393
2394
2395
2396
2397
2398
2399
2400
2401
2402
2403
2404
2405
2406
2407
2408
2409
2410
2411
2412
2413
2414
2415
2416
2417
2418
2419
2420
2421
2422
2423
2424
2425
2426
2427
2428
2429
2430
2431
2432
2433
2434
2435
2436
2437
2438
2439
2440
2441
2442
2443
2444
2445
2446
2447
2448
2449
2450
2451
2452
2453
2454
2455
2456
2457
2458
2459
2460
2461
2462
2463
2464
2465
2466
2467
2468
2469
2470
2471
2472
2473
2474
2475
2476
2477
2478
2479
2480
2481
2482
2483
2484
2485
2486
2487
2488
2489
2490
2491
2492
2493
2494
2495
2496
2497
2498
2499
2500
2501
2502
2503
2504
2505
2506
2507
2508
2509
2510
2511
2512
2513
2514
2515
2516
2517
2518
2519
2520
2521
2522
2523
2524
2525
2526
2527
2528
2529
2530
2531
2532
2533
2534
2535
2536
2537
2538
2539
2540
2541
2542
2543
2544
2545
2546
2547
2548
2549
2550
2551
2552
2553
2554
2555
2556
2557
2558
2559
2560
2561
2562
2563
2564
2565
2566
2567
2568
2569
2570
2571
2572
2573
2574
2575
2576
2577
2578
2579
2580
2581
2582
2583
2584
2585
2586
2587
2588
2589
2590
2591
2592
2593
2594
2595
2596
2597
2598
2599
2600
2601
2602
2603
2604
2605
2606
2607
2608
2609
2610
2611
2612
2613
2614
2615
2616
2617
2618
2619
2620
2621
2622
2623
2624
2625
2626
2627
2628
2629
2630
2631
2632
2633
2634
2635
2636
2637
2638
2639
2640
2641
2642
2643
2644
2645
2646
2647
2648
2649
2650
2651
2652
2653
2654
2655
2656
2657
2658
2659
2660
2661
2662
2663
2664
2665
2666
2667
2668
2669
2670
2671
2672
2673
2674
2675
2676
2677
2678
2679
2680
2681
2682
2683
2684
2685
2686
2687
2688
2689
2690
2691
2692
2693
2694
2695
2696
2697
2698
2699
2700
2701
2702
2703
2704
2705
2706
2707
2708
2709
2710
2711
2712
2713
2714
2715
2716
2717
2718
2719
2720
2721
2722
2723
2724
2725
2726
2727
2728
2729
2730
2731
2732
2733
2734
2735
2736
2737
2738
2739
2740
2741
2742
2743
2744
2745
2746
2747
2748
2749
2750
2751
2752
2753
2754
2755
2756
2757
2758
2759
2760
2761
2762
2763
2764
2765
2766
2767
2768
2769
2770
2771
2772
2773
2774
2775
2776
2777
2778
2779
2780
2781
2782
2783
2784
2785
2786
2787
2788
2789
2790
2791
2792
2793
2794
2795
2796
2797
2798
2799
2800
2801
2802
2803
2804
2805
2806
2807
2808
2809
2810
2811
2812
2813
2814
2815
2816
2817
2818
2819
2820
2821
2822
2823
2824
2825
2826
2827
2828
2829
2830
2831
2832
2833
2834
2835
2836
2837
2838
2839
2840
2841
2842
2843
2844
2845
2846
2847
2848
2849
2850
2851
2852
2853
2854
2855
2856
2857
2858
2859
2860
2861
2862
2863
2864
2865
2866
2867
2868
2869
2870
2871
2872
2873
2874
2875
2876
2877
2878
2879
2880
2881
2882
2883
2884
2885
2886
2887
2888
2889
2890
2891
2892
2893
2894
2895
2896
2897
2898
2899
2900
2901
2902
2903
2904
2905
2906
2907
2908
2909
2910
2911
2912
2913
2914
2915
2916
2917
2918
2919
2920
2921
2922
2923
2924
2925
2926
2927
2928
2929
2930
2931
2932
2933
2934
2935
2936
2937
2938
2939
2940
2941
2942
2943
2944
2945
2946
2947
2948
2949
2950
2951
2952
2953
2954
2955
2956
2957
2958
2959
2960
2961
2962
2963
2964
2965
2966
2967
2968
2969
2970
2971
2972
2973
2974
2975
2976
2977
2978
2979
2980
2981
2982
2983
2984
2985
2986
2987
2988
2989
2990
2991
2992
2993
2994
2995
2996
2997
2998
2999
2999
3000
3001
3002
3003
3004
3005
3006
3007
3008
3009
3009
3010
3011
3012
3013
3014
3015
3016
3017
3018
3019
3019
3020
3021
3022
3023
3024
3025
3026
3027
3028
3029
3029
3030
3031
3032
3033
3034
3035
3036
3037
3038
3039
3039
3040
3041
3042
3043
3044
3045
3046
3047
3048
3049
3049
3050
3051
3052
3053
3054
3055
3056
3057
3058
3059
3059
3060
3061
3062
3063
3064
3065
3066
3067
3068
3069
3069
3070
3071
3072
3073
3074
3075
3076
3077
3078
3079
3079
3080
3081
3082
3083
3084
3085
3086
3087
3087
3088
3089
3090
3091
3092
3093
3094
3095
3096
3097
3098
3099
3099
3100
3101
3102
3103
3104
3105
3106
3107
3108
3109
3109
3110
3111
3112
3113
3114
3115
3116
3117
3118
3119
3119
3120
3121
3122
3123
3124
3125
3126
3127
3128
3129
3129
3130
3131
3132
3133
3134
3135
3136
3137
3138
3138
3139
3140
3141
3142
3143
3144
3145
3146
3147
3147
3148
3149
3149
3150
3151
3152
3153
3153
3154
3155
3156
3157
3158
3159
3159
3160
3161
3162
3163
3163
3164
3165
3166
3167
3167
3168
3169
3169
3170
3171
3172
3172
3173
3173
3174
3174
3175
3175
3176
3176
3177
3177
3178
3178
3179
3179
3180
3180
3181
3181
3182
3182
3183
3183
3184
3184
3185
3185
3186
3186
3187
3187
3188
3188
3189
3189
3190
3190
3191
3191
3192
3192
3193
3193
3194
3194
3195
3195
3196
3196
3197
3197
3198
3198
3199
3199
3200
3200
3201
3201
3202
3202
3203
3203
3204
3204
3205
3205
3206

```

1188
1189
1190 Table 11: Statistics of the 1,622 multi-logic programming problems.
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241

| | #Problems | #Test Cases | #Langs | Prob Len | Solu Len | Difficulty (E/M/H) |
|------|-----------|-------------|--------|----------|----------|--------------------|
| MLPP | 1,622 | 16,131 | 20 | 576.4 | 610.0 | 238/315/1069 |

- **Smart Code Integration:** The system automatically manages the integration of function code with testing code. It adapts to language-specific syntax, ensuring seamless execution without requiring manual intervention for code merging.
- **High Performance:** Powered by a Gunicorn multi-process architecture, the sandbox supports concurrent execution of multiple code instances, making it capable of handling a high volume of requests efficiently.
- **RESTful API:** The service provides a clean and easy-to-use HTTP-based API, allowing developers to interact with the sandbox programmatically, whether for integrating into larger applications or automating tasks.
- **Extensive Language Support:** Beyond the mainstream languages, the sandbox also supports emerging and niche languages, allowing it to cater to a wide variety of development environments and user needs.
- **Custom Execution Environments:** Users can configure specific environments for their tasks, enabling tailored execution conditions based on their unique requirements.

K PROMPTS FOR AUTOMATED WORKFLOW

The prompt of generating code solution is shown in Figure 11.

The prompt of generating test function is shown in Figure 12.

The prompt of generating programming problem is shown in Figure 13.

The prompt of LLM-as-Critic is shown in Figure 10.

The prompt of translating languages is shown in Figure 14.

L MULTI-LOGIC TASK ANALYSIS

First, we provide a clear definition of multi-logic programming tasks:

Multi-Logic Programming Problem (MLPP). MLPP is a programming task whose correct solution requires implementing and coordinating multiple core logical units—such as functions, classes, or modules. Each logical unit corresponds to an independent semantic responsibility or algorithmic objective.

We further divide Multi-Logic Problems into two categories:

- **Intra-Logic Problems:** multiple logical units collaborate to achieve a single overarching functional goal.
- **Inter-Logic Problems:** multiple logical units address distinct functional goals, each independently contributing to the full solution.

Following the structure of Table 3, we perform a detailed analysis on the 1,622 multi-logic problems in AutoCodeBench. As shown in Table 11, multi-logic problems exhibit substantially longer problem descriptions and longer canonical solutions compared to the overall dataset, reflecting richer instruction structure and greater intrinsic difficulty.

In addition, we used DeepSeek-V3.2-Exp to compute (1) the number of logical units contained in each multi-logic problem, and (2) the distribution of Inter-Logic and Intra-Logic types. On average, each problem contains **3.37** logical units. Among them, **389** are Inter-Logic, **1,223** are Intra-Logic,

```

1242 # Code Benchmark Construction Task
1243 In order to build the code benchmark, I need you to help me create a Python function, as well as two test functions. Later, I will construct programming problems based on
1244 them.
1245
1246 ## Task Overview
1247 Given a code file or code snippet in any programming language, your task is to **refine and evolve its core logic** and remove non-core logic (if necessary) to create a
1248 **self-contained**, **verifiable**, **correct** Python function. You also need to create two testing functions for testing: `demo_testing()` and `full_testing()`. Below are
1249 the detailed requirements:
1250
1251 1. **Separate Code Blocks**:
1252 - Provide **three distinct code blocks**:
1253 - One block for the rewritten Python functions that contain the core logic.
1254 - Two blocks for the test functions (`demo_testing()` and `full_testing()`).
1255 - The three code blocks are wrapped with ``python and ```.
1256
1257 2. **Test Code in `demo_testing()`**:
1258 - This function should contain **no more than 2 test cases**. It will be used to **demonstrate** the input-output format.
1259 - For each test case, **hard-code the input values** and print both the input and the corresponding output when `demo_testing()` is executed.
1260 - Avoid edge cases in this function, and focus on basic test cases for demonstration purposes.
1261
1262 3. **Test Code in `full_testing()`**:
1263 - This function should include a comprehensive set of **at least 7 test cases**.
1264 - These test cases should be designed to thoroughly test the functionality of the converted functions, including:
1265 - **Basic cases** (e.g., test cases used in `demo_testing()`).
1266 - **Boundary cases** (e.g., minimum/maximum values, empty inputs).
1267 - **Edge cases** that test the robustness of the code (e.g., abnormal value, boundary condition checks).
1268 - For each test case, **hard-code the input values** and print both the input and the corresponding output when `full_testing()` is executed.
1269 - The goal of the test function is to provide comprehensive inputs and call the function being tested to obtain the output. Input handling (e.g., validation, edge cases)
1270 should be done by the function itself, not the test.
1271
1272 4. **Use Standard and Third-Party Python Libraries**:
1273 - You are allowed to use standard Python libraries (e.g., `math`, `math` , `print` ) and third-party libraries (if required) for implementing the core logic.
1274 - Ensure that all dependencies are clearly explained in the comments if used.
1275
1276 5. **Executable Code**:
1277 - Provide **self-contained and executable Python code** with hard-coded test inputs.
1278 - Remove any non-core logic such as plots, file I/O, and unused functions.
1279 - If the original code is not directly executable (e.g., due to dependency on specific business logic or external systems), adapt the code to a self-contained, testable form
1280 by simulating the required context or simplifying the logic while preserving the core algorithm.
1281 - For overly simplistic code (like a function that simply adds two numbers), please evolve its logic to make it more challenging. For overly difficult and complex code (like
1282 multiple functions and classes containing multiple core logics that make it hard to test), please simplify its logic to make it moderately challenging.
1283
1284 6. **Code Style and Naming**:
1285 - Follow **Pythonic** naming conventions and ensure that function names are **descriptive** and **clear** in conveying their purpose.
1286 - Provide **clear and concise comments** explaining the core logic and the purpose of each function, especially the key steps.
1287
1288 ## Output Requirements
1289
1290 #### 1. Converted Python Functions
1291 - **Dependencies**: Use standard Python libraries (e.g., math) and third-party libraries when necessary.
1292 - **Adaptation for Executability**: If the original code depends on external data sources or specific business workflows, modify it to be self-contained by hardcoding inputs,
1293 simulating dependencies, or simplifying the logic, while preserving the core algorithm.
1294 - **Remove Non-Core Logic**: Eliminate all non-essential parts such as plots, file I/O, and unused functions.
1295 - **Code Style**: Follow good Python coding practices, including proper indentation, consistent naming conventions, and adherence to PEP 8 guidelines. Add comments to
1296 explain the functionality of the code and clarify key steps.
1297
1298 #### 2. Executable Test Functions
1299 - **Function Names**: Use `demo_testing()` and `full_testing()` for the test functions.
1300 - **Inputs**: Hardcode the test inputs within the test functions. Do not pass parameters into the `demo_testing()` or `full_testing()` functions.
1301 - **Output**: Directly print the test inputs and their corresponding output to the console. Ensure that outputs are in JSON-serializable types (e.g., strings, integers, floats,
1302 lists). Do not include the word "Output" in the printed results.
1303 - **Dependencies**: Use standard Python libraries (e.g., math, print) and third-party libraries, if needed, for the core logic.
1304 - **Testing Adaptation**: If the original code's logic is part of a larger system, adapt the `demo_testing()` and `full_testing()` functions to create a self-contained testing
1305 environment. This may include initializing necessary objects, setting up test data, or mocking external interactions.
1306 - **Code Style**: Follow good Python code style in the test functions, with proper indentation and clarity. Include comments to explain the test setup and key steps in the
1307 testing process.
1308 - **Test Function Format**: The structure and format of `demo_testing()` and `full_testing()` should be identical, with the only difference being the number of test cases.
1309
1310 #### 3. Code Structure
1311 .....

```

Figure 11: The prompt of generating code solution. Due to the excessive length of the prompt, we have omitted the latter part.

and 10 fall into both categories. Opus 4.1 achieved 50.6% and 51.5% performance on the Inter-Logic and Intra-Logic types, respectively.

Figure 16 provides an example of an Inter-Logic problem. The task requires implementing three separate classes, each serving a different functional purpose—clearly corresponding to independent logical goals. Figure 17 shows a more complex hybrid case. The problem requires implementing a Rational class containing multiple internal logical units that jointly support the overarching goal of “rational number representation and operations,” making them Intra-Logic. At the same time, the auxiliary function “parse_rationals” handles a separate goal of “string parsing,” forming an Inter-Logic relationship with the Rational class.

The results of Opus 4.1 and these cases indicate that, in multi-logic problems, more complex instructions, greater logical demands, and intricate relationships among logical units pose greater challenges for LLMs.

```

1296 I'll provide you with a Python code, a test function call (including inputs) that uses this Python function, and the test output obtained after executing that test function call.
1297 Please combine the provided inputs and outputs into an assert statement, and place these `assert` statements inside a new `test` function. You can only use the inputs
1298 and outputs provided by me. Please do not create your own or modify the test cases.
1299
1300 Please generate a Python test function using assert statements. You will be provided with:
1301 - Python Code: The function(s) to be tested.
1302 - Two Test Function Calls (including inputs): Python code demonstrating how the function is called with various inputs.
1303 - Test Outputs: The results obtained after executing the provided test function calls.
1304 You will receive two sets of test cases:
1305
1306 Important Considerations for assert statements:
1307 - DO NOT create or modify any test cases. Use only the inputs and outputs provided.
1308 - Avoid assertions that might differ due to floating-point precision across machines. If the original output involves floating-point numbers, and the problem context suggests
1309 it, consider using math.isclose() or asserting within a reasonable tolerance if the test cases inherently involve such comparisons and precision is a concern. However,
1310 prioritize direct equality == if the provided outputs are exact.
1311 - All assert statements must be placed within a function named test().
1312 - Please create **two separate** test functions for each of the two sets of test function calls I provide, ensuring they do not interfere with each other.
1313 - Two test function names are both "def test()". They are placed in two code blocks.
1314
1315 Here is an example:
1316 **Code**:
1317 ```python
1318 from math import sqrt
1319
1320 def is_prime(n):
1321     if n < 2:
1322         return False
1323     for i in range(2, int(sqrt(n)) + 1):
1324         if n % i == 0:
1325             return False
1326     return True
1327
1328 def nth_prime(n):
1329     if n < 1:
1330         return None
1331     count = 1
1332     i = 2
1333     while count < n:
1334         i += 1
1335         if is_prime(i):
1336             count += 1
1337     return i
1338
1339
1340 **Test Function Call 1**:
1341 ```python
1342 def demo_testing():
1343     # Hard-coded test inputs
1344     test_cases = [1]
1345
1346     for n in test_cases:
1347         # Print the input
1348         print(f"Input: {n}")
1349
1350         # Core logic
1351         result = nth_prime(n)
1352
1353         # Output (JSON-serializable)
1354         print(f"The {n}th prime number is: {result}")
1355
1356 if __name__ == "__main__":
1357     demo_testing()
1358
1359
1360 **Test Case Results1**:
1361 ```
1362 'Input: 1\nThe 1th prime number is: 2'
1363
1364
1365 **Test Function Call 2**
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
1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457
1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
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
1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619
1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673
1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727
1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781
1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835
1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889
1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943
1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997
1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051
2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105
2106
2107
2108
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159
2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199
2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213
2214
2215
2216
2217
2218
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2250
2251
2252
2253
2254
2255
2256
2257
2258
2259
2260
2261
2262
2263
2264
2265
2266
2267
2268
2269
2270
2271
2272
2273
2274
2275
2276
2277
2278
2279
2280
2281
2282
2283
2284
2285
2286
2287
2288
2289
2290
2291
2292
2293
2294
2295
2296
2297
2298
2299
2300
2301
2302
2303
2304
2305
2306
2307
2308
2309
2310
2311
2312
2313
2314
2315
2316
2317
2318
2319
2320
2321
2322
2323
2324
2325
2326
2327
2328
2329
2330
2331
2332
2333
2334
2335
2336
2337
2338
2339
2340
2341
2342
2343
2344
2345
2346
2347
2348
2349
2350
2351
2352
2353
2354
2355
2356
2357
2358
2359
2360
2361
2362
2363
2364
2365
2366
2367
2368
2369
2370
2371
2372
2373
2374
2375
2376
2377
2378
2379
2380
2381
2382
2383
2384
2385
2386
2387
2388
2389
2390
2391
2392
2393
2394
2395
2396
2397
2398
2399
2400
2401
2402
2403
2404
2405
2406
2407
2408
2409
2410
2411
2412
2413
2414
2415
2416
2417
2418
2419
2420
2421
2422
2423
2424
2425
2426
2427
2428
2429
2430
2431
2432
2433
2434
2435
2436
2437
2438
2439
2440
2441
2442
2443
2444
2445
2446
2447
2448
2449
2450
2451
2452
2453
2454
2455
2456
2457
2458
2459
2460
2461
2462
2463
2464
2465
2466
2467
2468
2469
2470
2471
2472
2473
2474
2475
2476
2477
2478
2479
2480
2481
2482
2483
2484
2485
2486
2487
2488
2489
2490
2491
2492
2493
2494
2495
2496
2497
2498
2499
2500
2501
2502
2503
2504
2505
2506
2507
2508
2509
2510
2511
2512
2513
2514
2515
2516
2517
2518
2519
2520
2521
2522
2523
2524
2525
2526
2527
2528
2529
2530
2531
2532
2533
2534
2535
2536
2537
2538
2539
2540
2541
2542
2543
2544
2545
2546
2547
2548
2549
2550
2551
2552
2553
2554
2555
2556
2557
2558
2559
2560
2561
2562
2563
2564
2565
2566
2567
2568
2569
2570
2571
2572
2573
2574
2575
2576
2577
2578
2579
2580
2581
2582
2583
2584
2585
2586
2587
2588
2589
2590
2591
2592
2593
2594
2595
2596
2597
2598
2599
2600
2601
2602
2603
2604
2605
2606
2607
2608
2609
2610
2611
2612
2613
2614
2615
2616
2617
2618
2619
2620
2621
2622
2623
2624
2625
2626
2627
2628
2629
2630
2631
2632
2633
2634
2635
2636
2637
2638
2639
2640
2641
2642
2643
2644
2645
2646
2647
2648
2649
2650
2651
2652
2653
2654
2655
2656
2657
2658
2659
2660
2661
2662
2663
2664
2665
2666
2667
2668
2669
2670
2671
2672
2673
2674
2675
2676
2677
2678
2679
2680
2681
2682
2683
2684
2685
2686
2687
2688
2689
2690
2691
2692
2693
2694
2695
2696
2697
2698
2699
2700
2701
2702
2703
2704
2705
2706
2707
2708
2709
2710
2711
2712
2713
2714
2715
2716
2717
2718
2719
2720
2721
2722
2723
2724
2725
2726
2727
2728
2729
2730
2731
2732
2733
2734
2735
2736
2737
2738
2739
2740
2741
2742
2743
2744
2745
2746
2747
2748
2749
2750
2751
2752
2753
2754
2755
2756
2757
2758
2759
2760
2761
2762
2763
2764
2765
2766
2767
2768
2769
2770
2771
2772
2773
2774
2775
2776
2777
2778
2779
2780
2781
2782
2783
2784
2785
2786
2787
2788
2789
2790
2791
2792
2793
2794
2795
2796
2797
2798
2799
2800
2801
2802
2803
2804
2805
2806
2807
2808
2809
2810
2811
2812
2813
2814
2815
2816
2817
2818
2819
2820
2821
2822
2823
2824
2825
2826
2827
2828
2829
2830
2831
2832
2833
2834
2835
2836
2837
2838
2839
2840
2841
2842
2843
2844
2845
2846
2847
2848
2849
2850
2851
2852
2853
2854
2855
2856
2857
2858
2859
2860
2861
2862
2863
2864
2865
2866
2867
2868
2869
2870
2871
2872
2873
2874
2875
2876
2877
2878
2879
2880
2881
2882
2883
2884
2885
2886
2887
2888
2889
2890
2891
2892
2893
2894
2895
2896
2897
2898
2899
2900
2901
2902
2903
2904
2905
2906
2907
2908
2909
2910
2911
2912
2913
2914
2915
2916
2917
2918
2919
2920
2921
2922
2923
2924
2925
2926
2927
2928
2929
2930
2931
2932
2933
2934
2935
2936
2937
2938
2939
2940
2941
2942
2943
2944
2945
2946
2947
2948
2949
2950
2951
2952
2953
2954
2955
2956
2957
2958
2959
2960
2961
2962
2963
2964
2965
2966
2967
2968
2969
2970
2971
2972
2973
2974
2975
2976
2977
2978
2979
2980
2981
2982
2983
2984
2985
2986
2987
2988
2989
2990
2991
2992
2993
2994
2995
2996
2997
2998
2999
2999

```

Figure 12: The prompt of generating test function. Due to the excessive length of the prompt, we have omitted the latter part.

M LLM-BASED ERROR ANALYSIS

In this section, we present an error type analysis. As shown in Appendix A, a fully automated pipeline inevitably introduces a certain amount of noise. To distinguish noisy data from genuinely difficult tasks, we designed a checklist-based template to systematically categorize error types. The checklist is as follows:

- **Model Solution Errors**
 - **Misunderstanding Problem:** The model fails to correctly interpret the problem statement, leading to a fundamentally incorrect solution.
 - **Failure to Comprehend Complex Instructions:** The model fails to correctly interpret or decompose multi-step or nuanced requirements in the problem statement, resulting in conceptually incomplete or irrelevant solutions.

```

1350
1351 You are an experienced programming tutor, adept at crafting **clear, concise, and educational programming problems**. I will supply you with a **self-contained,
1352 executable, and correct Python code**, along with one or more test functions designed to verify its correctness. Your task is to **generate a programming problem that
1353 directly corresponds to the given Python code and its test cases**.
1354
1355 Here is the code and test functions:
1356 [Python code]
1357 <<<code>>>
1358 [Python code end]
1359
1360 [test function demo]
1361 <<<demo_test>>>
1362 [test function demo end]
1363
1364 [test function]
1365 <<<full_test>>>
1366 [test function end]
1367
1368 Please ensure the problem you generate adheres to the following critical requirements:
1369 1. Language Specification: Explicitly state that solutions must be implemented in Python.
1370 2. Problem Description: Describe the problem concisely and unambiguously using plain language. Avoid technical jargon, unnecessary details, or solution hints.
1371 3. Function/Class Naming: Only mention the exact function or class names used in the test functions. Do not include implementation-specific details beyond what's in the
1372 tests.
1373 4. Input/Output Format: Define the input format (types, structure, value ranges). Define the expected output format. If necessary, specify some constraints (e.g., input size
1374 limits, allowed data types).
1375 5. Example Usage: Use the test case in the aforementioned 'test function demo' to construct example usage. The number of test cases in this is usually no more than three.
1376 Do not modify or explain the test cases—just copy them verbatim.
1377 6. No Solution Hints: The problem description must not reveal any code solution and any
1378 test cases beyond what's in the provided examples.
1379
1380 Please enclose the generated programming problem within <question> and </question> tags.
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

```

Figure 13: The prompt of generating programming problem.

Table 12: Top-10 error types.

| Type | Count |
|--|-------|
| Unreasonable Test Cases | 313 |
| Algorithmic Logic Error | 218 |
| Missing Functionality | 157 |
| Misunderstanding Problem | 151 |
| Algorithmic Logic Error + Misunderstanding Problem | 121 |
| Insufficient Fundamental Knowledge | 108 |
| Algorithmic Logic Error + Missing Functionality | 106 |
| Algorithmic Logic Error + Unreasonable Test Cases | 73 |
| Algorithmic Logic Error + Insufficient Fundamental Knowledge | 57 |
| Omission of Edge Cases | 55 |

- **Omission of Edge Cases:** The model fails to handle extreme or special input conditions.
- **Missing Functionality:** The model understands the task but fails to implement all required components or subtasks in the generated code.
- **Algorithmic Logic Error:** The algorithmic approach or control logic contains mistakes leading to incorrect results.
- **Insufficient Fundamental Knowledge:** The model demonstrates a lack of understanding of basic programming syntax, data structures, or language semantics.
- **Problem Design Errors**
 - **Logical Contradiction in the Problem:** The problem statement contains internal inconsistencies or mutually exclusive constraints that make a correct solution impossible.
 - **Ambiguous Problem Statement:** The problem description is unclear or lacks explicit constraints, leading to multiple valid interpretations.
- **Test Design Errors**
 - **Unreasonable Test Cases:** The test data deviate from the problem definition.

We applied this checklist to analyze all 1,745 tasks that Claude Opus 4.1 failed to solve, using DeepSeek-V3.2-Exp as the judge LLM. If an error was attributed to **Problem Design** or **Test Design**, the task was classified as low-quality. Otherwise, we considered it a valid high-difficulty task. The Top-10 error types are shown in Table 12.

In total, 538 tasks contained *Problem Design Errors* or *Test Design Errors*, accounting for 13.7% of all tasks. The remaining 1,207 tasks were judged to be high-quality and genuinely challenging.

```

1404 You are an experienced programming master with expertise in multiple languages, particularly in PHP. I will provide
1405 you with:
1406 1. A code generation problem description (including requirements, constraints, and examples)
1407 2. The reference solution code (in another language, typically Python/Java/C++/Go/JavaScript/Shell)
1408 3. Corresponding test cases functions include demo test cases and full test cases
1409 Your task is to:
1410 1. Accurately translate the problem description into clear English while preserving all technical details. Please emphasize that the problem needs to be solved using PHP.
1411 2. Convert the reference solution into idiomatic PHP code that:
1412   - Follows PHP conventions
1413   - Uses appropriate data structures from PHP's standard library
1414   - Follows PHP's naming conventions
1415 3. Translate the test cases into PHP's testing framework:
1416   - Using assert
1417   - Maintaining the same test coverage as original
1418   - just call the test case in main function, do not use unit test library.
1419 Additional requirements:
1420   - For algorithm problems, maintain the same time/space complexity
1421   - For system design problems, use appropriate PHP crates
1422   - Include necessary documentation (/// comments) explaining key decisions
1423 Your answer should be like:
1424 <translated_problem>
1425 new problem
1426 </translated_problem>
1427
1428 <translated_reference_solution>
1429   ``php
1430   your code
1431   ``
1432 </translated_reference_solution>
1433
1434 <demo_test_cases>
1435   ``php
1436   your code
1437   ``
1438 </demo_test_cases>
1439
1440 <full_test_cases>
1441   ``php
1442   your code
1443   ``
1444 </full_test_cases>
1445
1446 Here is the problem, reference solution and test case
1447 [problem]
1448 <<<problem>>>
1449 [problem end]
1450
1451 [reference solution]
1452 <<<code>>>
1453 [reference solution end]
1454
1455 [demo test case]
1456 <<<demo_test>>>
1457 [demo test case end]
1458
1459 [full test case]
1460 <<<full_test>>>
1461 [full test case end]

```

Figure 14: The prompt of translating languages.

Among these, the dominant failure modes were *Algorithmic Logic Errors* and *Missing Functionality*. Figures 18 and 19 illustrate representative cases. After excluding low-quality tasks, Opus 4.1 still achieves only 64.3%, demonstrating that the remaining tasks remain substantially challenging.

N THE ORDER OF I/O FORMAT

In this section, we elaborate on the design choice of generating test inputs prior to producing the programming problem description in the AutoCodeGen pipeline. Although human annotators typically define the input–output (I/O) format before writing tests, our automated workflow adopts the reverse order. We show that this alternative workflow is both feasible for LLMs and beneficial in practice. There are two valid sequences for synthesizing test cases and problem descriptions:

- Generate the test inputs first, followed by the programming problem (including its I/O specification);
- Generate the I/O specification first, followed by the test inputs.

Both workflows can be handled by modern LLMs. However, the first approach offers an additional advantage: the public test cases produced during this stage can be directly embedded into the final problem description, improving clarity and reducing ambiguity in the specification. Below is an example that illustrates why generating test inputs first is a reasonable design choice.

Consider a code solution such as:

```

1458 Please implement the following three network-related functions in Go:
1459
1460 1. ReverseIPAddress: Accepts an IPv4 address of type net.IP and returns its reversed string representation (e.g., "192.168.1.1" becomes "1.1.168.192"). If the input is invalid, return an empty string.
1461
1462 2. ReputationLookup: Accepts an IPv4 address as a string and returns its reputation value as a string (e.g., "100"). If the input is invalid, return an error.
1463
1464 3. HostsFromCIDR: Accepts a CIDR-formatted string (e.g., "192.168.1.0/24") and returns a string slice of all available host
1465 IPs within that range (excluding the network and broadcast addresses). If the input is invalid, return an error.
1466
1467 Input/Output Requirements:
1468 - All IP addresses must be in IPv4 format.
1469 - Invalid inputs include malformed or unparsable addresses.
1470 - The CIDR range must be valid and contain at least one available host address.
1471
1472 Example Usage:
1473
1474 // Test case 1: Reverse IP address
1475 ip := net.ParseIP("192.168.1.1")
1476 reversed := ReverseIPAddress(ip)
1477 // reversed should be "1.1.168.192"
1478
1479 // Test case 2: Query reputation
1480 ipStr := "8.8.8.8"
1481 rep, err := ReputationLookup(ipStr)
1482 // rep should be "100", err should be nil
1483
1484 // Test case 3: Get host list from CIDR
1485 cidr := "192.168.1.0/24"
1486 hosts, err := HostsFromCIDR(cidr)
1487 // hosts should be ["192.168.1.1", "192.168.1.2", "192.168.1.3", "192.168.1.4", "192.168.1.5", "192.168.1.6"], err should
1488 be nil
1489
1490

```

Figure 15: The prompt of generating code solution. Due to the excessive length of the prompt, we have omitted the latter part.

```

1479
1480 def cal_two_sum(a, b):
1481     return a + b
1482

```

Using explicit specifications and few-shot examples, the model may generate a test-input function such as:

```

1485 def private_test_input():
1486     test_cases = [
1487         [1, 2],
1488         [0, 0], # boundary test
1489         [9999999999999999, 3333333333] # stress test
1490     ]

```

When constructing the programming problem, the model infers the input domain from the generated test cases and produces a description such as: *Write a function named cal_two_sum that computes the sum of two non-negative integers.* If instead the model generates test cases containing floating-point values:

```

1495
1496 def private_test_input():
1497     test_cases = [
1498         [1, 2],
1499         [0, 0],
1500         [4.243, 4.222], # stress test
1501         [-0.45, 888], # boundary test
1502     ]

```

The model naturally adapts the problem description accordingly: *Write a function named cal_two_sum that computes the sum of two floating-point numbers.* Therefore, When test inputs are generated first, the model derives the I/O format by analyzing the input domain represented in the test cases. When the I/O format is specified first, the model constructs matching test cases based on the declared constraints. Both approaches are logically consistent, but generating test inputs first leverages the model’s ability to generalize from concrete examples and simplifies the subsequent problem construction stage.

Finally, regardless of the workflow, AutoCodeGen employs an LLM-as-Critic verification stage to ensure strict alignment between the generated programming problem and the test function. Any inconsistent, ambiguous, or underspecified cases are filtered out during this procedure. This ensures the final benchmark maintains high quality and internal coherence.

```

1512
1513
1514
1515
1516
1517
1518
1519
1520 # 3D Camera System Implementation
1521
1522 ## Problem Description
1523 You are tasked with implementing a 3D camera system for a game engine. The system should include:
1524 1. A Vector class for 3D mathematical operations
1525 2. An AABB (Axis-Aligned Bounding Box) class for collision detection
1526 3. A Camera class that represents a view frustum in 3D space
1527
1528 The camera should be able to:
1529 - Track its position and orientation in 3D space
1530 - Calculate its view frustum's bounding box
1531 - Determine if objects (represented as AABBs) are within its view frustum
1532
1533 ## Class Requirements
1534
1535 ### Vector Class
1536 ````cpp
1537 class Vector {
1538 public:
1539     float x, y, z;
1540
1541     Vector(float x=0, float y=0, float z=0);
1542     Vector operator+(const Vector& other) const;
1543     Vector operator-(const Vector& other) const;
1544     Vector operator*(float scalar) const;
1545     float dot(const Vector& other) const;
1546     Vector cross(const Vector& other) const;
1547     float magnitude() const;
1548     Vector normalized() const;
1549     void print() const;
1550 };
1551
1552 ### AABB Class
1553 ````cpp
1554 class AABB {
1555 public:
1556     Vector min, max;
1557
1558     AABB(Vector min=Vector(), Vector max=Vector());
1559     bool contains(const Vector& point) const;
1560     void print() const;
1561 };
1562
1563 ### Camera Class
1564 ````cpp
1565 class Camera {
1566 public:
1567     Camera(const string& name="DefaultCamera",
1568           const Vector3 pose=Vector(),
1569           const Vector3 rotator=Vector(),
1570           float fov=60.0f,
1571           float zFar=1000.0f,
1572           float zNear=0.1f,
1573           float width=800.0f,
1574           float height=600.0f);
1575
1576 };
1577
1578 ## Example Usage
1579 ````cpp
1580 ...
1581
1582 ## Problem Specifications
1583 ...
1584
1585 ## Constraints
1586 ...
1587
1588 ## Notes
1589 ...

```

Figure 16: Example of an inter-logic problem. The full content is partially omitted due to length.

```

1566
1567
1568
1569
1570
1571
1572
1573
1574 # Rational Number Operations
1575 Implement a class 'Rational' to represent and manipulate rational numbers (fractions) with the following capabilities:
1576 ## Class Requirements
1577
1578 ## Private Members
1579 - 'long long numerator': The numerator of the fraction
1580 - 'long long denominator': The denominator of the fraction
1581 - 'long long gcd(long a, long long b) const': A helper function to compute the greatest common divisor of two numbers
1582 - 'void normalize()': A helper function to reduce the fraction to its simplest form and ensure the denominator is positive
1583
1584 ## Public Members
1585 ## Constructor
1586 - 'Rational(long num = 0, long long denom = 1)': Creates a Rational number with given numerator and denominator
1587 (defaults to 0/1). Throws 'invalid_argument' if denominator is zero.
1588
1589 ## Arithmetic Operators
1590 - 'Rational operator+(const Rational& other) const': Adds two Rational numbers
1591 - 'Rational operator-(const Rational& other) const': Subtracts two Rational numbers
1592 - 'Rational operator*(const Rational& other) const': Multiplies two Rational numbers
1593 - 'Rational operator/(const Rational& other) const': Divides two Rational numbers (throws 'domain_error' if dividing by
1594 zero)
1595
1596 ## Comparison Operators
1597 - 'bool operator==(const Rational& other) const': Checks equality of two Rational numbers
1598 - 'bool operator<(const Rational& other) const': Checks if this Rational is less than another
1599 - 'bool operator>(const Rational& other) const': Checks if this Rational is greater than another
1600
1601 ## Conversion Functions
1602 - 'double to_double() const': Converts the Rational number to a double
1603 - 'string to_string(bool show_parens = false) const': Returns a string representation of the Rational number. When
1604 'show_parens' is true, negative numbers should be enclosed in parentheses. The string should be in mixed number form when
1605 appropriate (e.g., "2 1/3" for 7/3).
1606
1607 ## Getters
1608 - 'long long get_numerator() const': Returns the numerator
1609 - 'long long get_denominator() const': Returns the denominator
1610
1611 ## Helper Function
1612 - 'vector<Rational> parse_rationals(const string& input)': Parses a space-separated string of rational numbers (either in
1613 "a/b" form or whole numbers) and returns a vector of Rational numbers.
1614
1615 ## Example Usage
1616
1617
1618
1619
1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673
1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727
1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781
1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835
1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889
1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943
1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997
1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051
2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105
2106
2107
2108
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159
2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199
2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213
2214
2215
2216
2217
2218
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2250
2251
2252
2253
2254
2255
2256
2257
2258
2259
2260
2261
2262
2263
2264
2265
2266
2267
2268
2269
2270
2271
2272
2273
2274
2275
2276
2277
2278
2279
2280
2281
2282
2283
2284
2285
2286
2287
2288
2289
2290
2291
2292
2293
2294
2295
2296
2297
2298
2299
2300
2301
2302
2303
2304
2305
2306
2307
2308
2309
2310
2311
2312
2313
2314
2315
2316
2317
2318
2319
2320
2321
2322
2323
2324
2325
2326
2327
2328
2329
2330
2331
2332
2333
2334
2335
2336
2337
2338
2339
2340
2341
2342
2343
2344
2345
2346
2347
2348
2349
2350
2351
2352
2353
2354
2355
2356
2357
2358
2359
2360
2361
2362
2363
2364
2365
2366
2367
2368
2369
2370
2371
2372
2373
2374
2375
2376
2377
2378
2379
2380
2381
2382
2383
2384
2385
2386
2387
2388
2389
2390
2391
2392
2393
2394
2395
2396
2397
2398
2399
2400
2401
2402
2403
2404
2405
2406
2407
2408
2409
2410
2411
2412
2413
2414
2415
2416
2417
2418
2419
2420
2421
2422
2423
2424
2425
2426
2427
2428
2429
2430
2431
2432
2433
2434
2435
2436
2437
2438
2439
2440
2441
2442
2443
2444
2445
2446
2447
2448
2449
2450
2451
2452
2453
2454
2455
2456
2457
2458
2459
2460
2461
2462
2463
2464
2465
2466
2467
2468
2469
2470
2471
2472
2473
2474
2475
2476
2477
2478
2479
2480
2481
2482
2483
2484
2485
2486
2487
2488
2489
2490
2491
2492
2493
2494
2495
2496
2497
2498
2499
2500
2501
2502
2503
2504
2505
2506
2507
2508
2509
2510
2511
2512
2513
2514
2515
2516
2517
2518
2519
2520
2521
2522
2523
2524
2525
2526
2527
2528
2529
2530
2531
2532
2533
2534
2535
2536
2537
2538
2539
2540
2541
2542
2543
2544
2545
2546
2547
2548
2549
2550
2551
2552
2553
2554
2555
2556
2557
2558
2559
2560
2561
2562
2563
2564
2565
2566
2567
2568
2569
2570
2571
2572
2573
2574
2575
2576
2577
2578
2579
2580
2581
2582
2583
2584
2585
2586
2587
2588
2589
2590
2591
2592
2593
2594
2595
2596
2597
2598
2599
2600
2601
2602
2603
2604
2605
2606
2607
2608
2609
2610
2611
2612
2613
2614
2615
2616
2617
2618
2619
2620
2621
2622
2623
2624
2625
2626
2627
2628
2629
2630
2631
2632
2633
2634
2635
2636
2637
2638
2639
2640
2641
2642
2643
2644
2645
2646
2647
2648
2649
2650
2651
2652
2653
2654
2655
2656
2657
2658
2659
2660
2661
2662
2663
2664
2665
2666
2667
2668
2669
2670
2671
2672
2673
2674
2675
2676
2677
2678
2679
2680
2681
2682
2683
2684
2685
2686
2687
2688
2689
2690
2691
2692
2693
2694
2695
2696
2697
2698
2699
2700
2701
2702
2703
2704
2705
2706
2707
2708
2709
2710
2711
2712
2713
2714
2715
2716
2717
2718
2719
2720
2721
2722
2723
2724
2725
2726
2727
2728
2729
2730
2731
2732
2733
2734
2735
2736
2737
2738
2739
2740
2741
2742
2743
2744
2745
2746
2747
2748
2749
2750
2751
2752
2753
2754
2755
2756
2757
2758
2759
2760
2761
2762
2763
2764
2765
2766
2767
2768
2769
2770
2771
2772
2773
2774
2775
2776
2777
2778
2779
2780
2781
2782
2783
2784
2785
2786
2787
2788
2789
2790
2791
2792
2793
2794
2795
2796
2797
2798
2799
2800
2801
2802
2803
2804
2805
2806
2807
2808
2809
2810
2811
2812
2813
2814
2815
2816
2817
2818
2819
2820
2821
2822
2823
2824
2825
2826
2827
2828
2829
2830
2831
2832
2833
2834
2835
2836
2837
2838
2839
2840
2841
2842
2843
2844
2845
2846
2847
2848
2849
2850
2851
2852
2853
2854
2855
2856
2857
2858
2859
2860
2861
2862
2863
2864
2865
2866
2867
2868
2869
2870
2871
2872
2873
2874
2875
2876
2877
2878
2879
2880
2881
2882
2883
2884
2885
2886
2887
2888
2889
2890
2891
2892
2893
2894
2895
2896
2897
2898
2899
2900
2901
2902
2903
2904
2905
2906
2907
2908
2909
2910
2911
2912
2913
2914
2915
2916
2917
2918
2919
2920
2921
2922
2923
2924
2925
2926
2927
2928
2929
2930
2931
2932
2933
2934
2935
2936
2937
2938
2939
2940
2941
2942
2943
2944
2945
2946
2947
2948
2949
2950
2951
2952
2953
2954
2955
2956
2957
2958
2959
2960
2961
2962
2963
2964
2965
2966
2967
2968
2969
2970
2971
2972
2973
2974
2975
2976
2977
2978
2979
2980
2981
2982
2983
2984
2985
2986
2987
2988
2989
2990
2991
2992
2993
2994
2995
2996
2997
2998
2999
3000
3001
3002
3003
3004
3005
3006
3007
3008
3009
3010
3011
3012
3013
3014
3015
3016
3017
3018
3019
3020
3021
3022
3023
3024
3025
3026
3027
3028
3029
3030
3031
3032
3033
3034
3035
3036
3037
3038
3039
3040
3041
3042
3043
3044
3045
3046
3047
3048
3049
3050
3051
3052
3053
3054
3055
3056
3057
3058
3059
3060
3061
3062
3063
3064
3065
3066
3067
3068
3069
3070
3071
3072
3073
3074
3075
3076
3077
3078
3079
3080
3081
3082
3083
3084
3085
3086
3087
3088
3089
3090
3091
3092
3093
3094
3095
3096
3097
3098
3099
3100
3101
3102
3103
3104
3105
3106
3107
3108
3109
3110
3111
3112
3113
3114
3115
3116
3117
3118
3119
3120
3121
3122
3123
3124
3125
3126
3127
3128
3129
3130
3131
3132
3133
3134
3135
3136
3137
3138
3139
3140
3141
3142
3143
3144
3145
3146
3147
3148
3149
3150
3151
3152
3153
3154
3155
3156
3157
3158
3159
3160
3161
3162
3163
3164
3165
3166
3167
3168
3169
3170
3171
3172
3173
3174
3175
3176
3177
3178
3179
3180
3181
3182
3183
3184
3185
3186
3187
3188
3189
3190
3191
3192
3193
3194
3195
3196
3197
3198
3199
3200
3201
3202
3203
3204
3205
3206
3207
3208
3209
3210
3211
3212
3213
3214
3215
3216
3217
3218
3219
3220
3221
3222
3223
3224
3225
3226
3227
3228
3229
3230
3231
3232
3233
3234
3235
3236
3237
3238
3239
3240
3241
3242
3243
3244
3245
3246
3247
3248
3249
3250
3251
3252
3253
3254
3255
3256
3257
3258
3259
3260
3261
3262
3263
3264
3265
3266
3267
3268
3269
3270
3271
3272
3273
3274
3275
3276
3277
3278
3279
3280
3281
3282
3283
3284
3285
3286
3287
3288
3289
3290
3291
3292
3293
3294
3295
3296
3297
3298
3299
3300
3301
3302
3303
3304
3305
3306
3307
3308
3309
3310
3311
3312
3313
3314
3315
3316
3317
3318
3319
3320
3321
3322
3323
3324
3325
3326
3327
3328
3329
3330
3331
3332
3333
3334
3335
3336
3337
3338
3339
3340
3341
3342
3343
3344
3345
3346
3347
3348
3349
3350
3351
3352
3353
3354
3355
3356
3357
3358
3359
3360
3361
3362
3363
3364
3365
3366
3367
3368
3369
3370
3371
3372
3373
3374
3375
3376
3377
3378
3379
3380
3381
3382
3383
3384
3385
3386
3387
3388
3389
3390
3391
3392
3393
3394
3395
3396
3397
3398
3399
3400
3401
3402
3403
3404
3405
3406
3407
3408
3409
3410
3411
3412
3413
3414
3415
3416
3417
3418
3419
3420
3421
3422
3423
3424
3425
3426
3427
3428
3429
3430
3431
3432
3433
3434
3435
3436
3437
3438
3439
3440
3441
3442
3443
3444
3445
3446
3447
3448
3449
3450
3451
3452
3453
3454
3455
3456
3457
3458
3459
3460
3461
3462
3463
3464
3465
3466
3467
3468
3469
3470
3471
3472
3473
3474
3475
3476
3477
3478
3479
3480
3481
3482
3483
3484
3485
3486
3487
3488
3489
3490
3491
3492
3493
3494
3495
3496
3497
3498
3499
3500
3501
3502
3503
3504
3505
3506
3507
3508
3509
3510
3511
3512
3513
3514
3515
3516
3517
3518
3519
3520
3521
3522
3523
3524
3525
3526
3527
3528
3529
3530
3531
3532
3533
3534
3535
3536
3537
3538
3539
3540
3541
3542
3543
3544
3545
3546
3547
3548
3549
3550
3551
3552
3553
3554
3555
3556
3557
3558
3559
3560
3561
3562
3563
3564
3565
3566
3567
3568
3569
3570
3571
3572
3573
3574
3575
3576
3577
3578
3579
3580
3581
3582
3583
3584
3585
3586
3587
3588
3589
3590
3591
3592
3593
3594
3595
3596
3597
3598
3599
3600
3601
3602
3603
3604
3605
3606
3607
3608
3609
3610
3611
3612
3613
3614
3615
3616
3617
3618
3619
3620
3621
3622
3623
3624
3625
3626
3627
3628
3629
3630
3631
3632
3633
3634
3635
3636
3637
3638
3639
3640
3641
3642
3643
3644
3645
3646
3647
3648
3649
3650
3651
3652
3653
3654
3655
3656
3657
3658
3659
3660
3661
3662
3663
3664
3665
3666
3667
3668
3669
3670
3671
3672
3673
3674
3675
3676
3677
3678
3679
3680
3681
3682
3683
3684
3685
3686
3687
3688
3689
3690
3691
3692
3693
3694
3695
3696
3697
3698
3699
3700
3701
3702
3703
3704
3705
3706
3707
3708
3709
3710
3711
3712
3713
3714
3715
3716
3717
3718
3719
3720
3721
3722
3723
3724
3725
3726
3727
3728
3729
3730
3731
3732
3733
3734
3735

```

```

1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630 [programming problem]
1631 Write a JavaScript function 'octalToBinary' that converts a given octal number string into its corresponding binary string.
1632
1633 Input requirements:
1634 - The input is a string representing an octal number.
1635 - The string may be empty or contain non-octal digit characters (characters other than 0-7).
1636 - Valid octal numbers may include leading zeros (e.g., "0000").
1637
1638 Output requirements:
1639 - If the input is a valid octal string, return the corresponding binary string representation.
1640 - If the input is an empty string, return "Error: Empty input".
1641 - If the input contains non-octal digit characters, return "Error: Invalid octal string".
1642 - The binary output should not omit leading zeros (e.g., octal "7" should return "111" and not "111").
1643
1644 Example usage:
1645 ```javascript
1646 const assert = require('assert');
1647
1648 const demoTesting = () => {
1649     assert.strictEqual(octalToBinary("123"), "1010011");
1650     assert.strictEqual(octalToBinary("7"), "111");
1651 }
1652
1653 [Private Test Function]
1654 const assert = require('assert');
1655
1656 const fullTesting = () => {
1657     assert.strictEqual(octalToBinary("123"), "1010011");
1658     assert.strictEqual(octalToBinary("7"), "111");
1659     assert.strictEqual(octalToBinary("0"), "0");
1660     assert.strictEqual(octalToBinary("777"), "11111111");
1661     assert.strictEqual(octalToBinary("100"), "1000000");
1662     assert.strictEqual(octalToBinary("12345678"), "101001110010111011000");
1663     assert.strictEqual(octalToBinary(""), "Error: Empty input");
1664     assert.strictEqual(octalToBinary("abc"), "Error: Invalid octal string");
1665     assert.strictEqual(octalToBinary("89"), "Error: Invalid octal string");
1666     assert.strictEqual(octalToBinary("0000"), "0000000000");
1667 }
1668
1669 [Code Solution from Opus 4.1]
1670 ```javascript
1671 function octalToBinary(octalStr) {
1672     // Check for empty input
1673     if (octalStr === '') {
1674         return "Error: Empty input";
1675     }
1676
1677     // Check if all characters are valid octal digits (0-7)
1678     for (let char of octalStr) {
1679         if (char < '0' || char > '7') {
1680             return "Error: Invalid octal string";
1681         }
1682     }
1683
1684     // Convert each octal digit to 3-bit binary representation
1685     let binaryStr = '';
1686     for (let digit of octalStr) {
1687         const octalDigit = parseInt(digit);
1688         const binaryDigit = octalDigit.toString(2).padStart(3, '0');
1689         binaryStr += binaryDigit;
1690     }
1691
1692     // Remove leading zeros, but keep at least one digit
1693     while (binaryStr.length > 1 && binaryStr[0] === '0') {
1694         binaryStr = binaryStr.substring(1);
1695     }
1696
1697     return binaryStr;
1698 }
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727
1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781
1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835
1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889
1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943
1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997
1998
1999
1999

```

Logic Error!

Figure 18: A case of Algorithmic logic errors.

```

1674
1675
1676
1677
1678
1679
1680
1681
1682
1683 [programming problem]
1684 Write a Go function to find the lowest common ancestor of two nodes in a binary tree.
1685 Given the root of a binary tree and two nodes p and q in the tree, implement the function 'lowestCommonAncestor' to return
1686 the lowest common ancestor (LCA) of these two nodes. The lowest common ancestor is defined as the lowest node that has both
1687 p and q as descendants (where a node can be a descendant of itself).
1688 The function signature is as follows:
1689 ```go
1690 func lowestCommonAncestor(root, p, q *TreeNode) *TreeNode
1691 ````
1692 The 'TreeNode' is defined as:
1693 ```go
1694 type TreeNode struct {
1695     Val   int
1696     Left  *TreeNode
1697     Right *TreeNode
1698 }
1699 ````
1700 Input constraints:
1701 - All node values in the tree are unique.
1702 - Both p and q exist in the given binary tree.
1703 - The number of nodes in the tree ranges from [2, 10^5].
1704 Example usage:
1705 ```go
1706 import "testing"
1707
1708 func TestDemo(t *testing.T) {
1709     // Create the tree structure for test case 1
1710     node7 := &TreeNode{Val: 7, Left: nil, Right: nil}
1711     node4 := &TreeNode{Val: 4, Left: nil, Right: nil}
1712     node2 := &TreeNode{Val: 2, Left: node7, Right: node4}
1713     node6 := &TreeNode{Val: 6, Left: nil, Right: nil}
1714
1715 }
1716 ````
1717 Notes:
1718 - Ensure your solution handles various cases, including but not limited to: p and q in different subtrees, p being an
1719 ancestor of q, q being an ancestor of p, p and q being the same node, etc.
1720 - Do not modify the tree structure.
1721 - The time complexity should be as efficient as possible.
1722 [Code Solution from Opus 4.1]
1723 ```go
1724 func lowestCommonAncestor(root, p, q *TreeNode) *TreeNode {
1725     // Base case: if root is nil or root is one of the target nodes
1726     if root == nil || root == p || root == q {
1727         return root
1728     }
1729
1730     // Recursively search for p and q in left and right subtrees
1731     left := lowestCommonAncestor(root.Left, p, q)
1732     right := lowestCommonAncestor(root.Right, p, q)
1733
1734     // If both left and right are non-nil, current node is the LCA
1735     // This means p and q are in different subtrees
1736     if left != nil && right != nil {
1737         return root
1738     }
1739
1740     // If only one subtree contains both nodes, return that result
1741     // If left is non-nil, both nodes are in left subtree
1742     if left != nil {
1743         return left
1744     }
1745     // Otherwise, both nodes are in right subtree (or both are nil)
1746     return right
1747 }
1748 ````
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781
1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835
1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889
1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943
1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997
1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051
2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105
2106
2107
2108
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159
2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199
2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213
2214
2215
2216
2217
2218
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2250
2251
2252
2253
2254
2255
2256
2257
2258
2259
2260
2261
2262
2263
2264
2265
2266
2267
2268
2269
2270
2271
2272
2273
2274
2275
2276
2277
2278
2279
2280
2281
2282
2283
2284
2285
2286
2287
2288
2289
2290
2291
2292
2293
2294
2295
2296
2297
2298
2299
2300
2301
2302
2303
2304
2305
2306
2307
2308
2309
2310
2311
2312
2313
2314
2315
2316
2317
2318
2319
2320
2321
2322
2323
2324
2325
2326
2327
2328
2329
2330
2331
2332
2333
2334
2335
2336
2337
2338
2339
2340
2341
2342
2343
2344
2345
2346
2347
2348
2349
2350
2351
2352
2353
2354
2355
2356
2357
2358
2359
2360
2361
2362
2363
2364
2365
2366
2367
2368
2369
2370
2371
2372
2373
2374
2375
2376
2377
2378
2379
2380
2381
2382
2383
2384
2385
2386
2387
2388
2389
2390
2391
2392
2393
2394
2395
2396
2397
2398
2399
2400
2401
2402
2403
2404
2405
2406
2407
2408
2409
2410
2411
2412
2413
2414
2415
2416
2417
2418
2419
2420
2421
2422
2423
2424
2425
2426
2427
2428
2429
2430
2431
2432
2433
2434
2435
2436
2437
2438
2439
2440
2441
2442
2443
2444
2445
2446
2447
2448
2449
2450
2451
2452
2453
2454
2455
2456
2457
2458
2459
2460
2461
2462
2463
2464
2465
2466
2467
2468
2469
2470
2471
2472
2473
2474
2475
2476
2477
2478
2479
2480
2481
2482
2483
2484
2485
2486
2487
2488
2489
2490
2491
2492
2493
2494
2495
2496
2497
2498
2499
2500
2501
2502
2503
2504
2505
2506
2507
2508
2509
2510
2511
2512
2513
2514
2515
2516
2517
2518
2519
2520
2521
2522
2523
2524
2525
2526
2527
2528
2529
2530
2531
2532
2533
2534
2535
2536
2537
2538
2539
2540
2541
2542
2543
2544
2545
2546
2547
2548
2549
2550
2551
2552
2553
2554
2555
2556
2557
2558
2559
2560
2561
2562
2563
2564
2565
2566
2567
2568
2569
2570
2571
2572
2573
2574
2575
2576
2577
2578
2579
2580
2581
2582
2583
2584
2585
2586
2587
2588
2589
2590
2591
2592
2593
2594
2595
2596
2597
2598
2599
2600
2601
2602
2603
2604
2605
2606
2607
2608
2609
2610
2611
2612
2613
2614
2615
2616
2617
2618
2619
2620
2621
2622
2623
2624
2625
2626
2627
2628
2629
2630
2631
2632
2633
2634
2635
2636
2637
2638
2639
2640
2641
2642
2643
2644
2645
2646
2647
2648
2649
2650
2651
2652
2653
2654
2655
2656
2657
2658
2659
2660
2661
2662
2663
2664
2665
2666
2667
2668
2669
2670
2671
2672
2673
2674
2675
2676
2677
2678
2679
2680
2681
2682
2683
2684
2685
2686
2687
2688
2689
2690
2691
2692
2693
2694
2695
2696
2697
2698
2699
2700
2701
2702
2703
2704
2705
2706
2707
2708
2709
2710
2711
2712
2713
2714
2715
2716
2717
2718
2719
2720
2721
2722
2723
2724
2725
2726
2727
2728
2729
2730
2731
2732
2733
2734
2735
2736
2737
2738
2739
2740
2741
2742
2743
2744
2745
2746
2747
2748
2749
2750
2751
2752
2753
2754
2755
2756
2757
2758
2759
2760
2761
2762
2763
2764
2765
2766
2767
2768
2769
2770
2771
2772
2773
2774
2775
2776
2777
2778
2779
2780
2781
2782
2783
2784
2785
2786
2787
2788
2789
2790
2791
2792
2793
2794
2795
2796
2797
2798
2799
2800
2801
2802
2803
2804
2805
2806
2807
2808
2809
2810
2811
2812
2813
2814
2815
2816
2817
2818
2819
2820
2821
2822
2823
2824
2825
2826
2827
2828
2829
2830
2831
2832
2833
2834
2835
2836
2837
2838
2839
2840
2841
2842
2843
2844
2845
2846
2847
2848
2849
2850
2851
2852
2853
2854
2855
2856
2857
2858
2859
2860
2861
2862
2863
2864
2865
2866
2867
2868
2869
2870
2871
2872
2873
2874
2875
2876
2877
2878
2879
2880
2881
2882
2883
2884
2885
2886
2887
2888
2889
2890
2891
2892
2893
2894
2895
2896
2897
2898
2899
2900
2901
2902
2903
2904
2905
2906
2907
2908
2909
2910
2911
2912
2913
2914
2915
2916
2917
2918
2919
2920
2921
2922
2923
2924
2925
2926
2927
2928
2929
2930
2931
2932
2933
2934
2935
2936
2937
2938
2939
2940
2941
2942
2943
2944
2945
2946
2947
2948
2949
2950
2951
2952
2953
2954
2955
2956
2957
2958
2959
2960
2961
2962
2963
2964
2965
2966
2967
2968
2969
2970
2971
2972
2973
2974
2975
2976
2977
2978
2979
2980
2981
2982
2983
2984
2985
2986
2987
2988
2989
2990
2991
2992
2993
2994
2995
2996
2997
2998
2999
3000
3001
3002
3003
3004
3005
3006
3007
3008
3009
3010
3011
3012
3013
3014
3015
3016
3017
3018
3019
3020
3021
3022
3023
3024
3025
3026
3027
3028
3029
3030
3031
3032
3033
3034
3035
3036
3037
3038
3039
3040
3041
3042
3043
3044
3045
3046
3047
3048
3049
3050
3051
3052
3053
3054
3055
3056
3057
3058
3059
3060
3061
3062
3063
3064
3065
3066
3067
3068
3069
3070
3071
3072
3073
3074
3075
3076
3077
3078
3079
3080
3081
3082
3083
3084
3085
3086
3087
3088
3089
3090
3091
3092
3093
3094
3095
3096
3097
3098
3099
3100
3101
3102
3103
3104
3105
3106
3107
3108
3109
3110
3111
3112
3113
3114
3115
3116
3117
3118
3119
3120
3121
3122
3123
3124
3125
3126
3127
3128
3129
3130
3131
3132
3133
3134
3135
3136
3137
3138
3139
3140
3141
3142
3143
3144
3145
3146
3147
3148
3149
3150
3151
3152
3153
3154
3155
3156
3157
3158
3159
3160
3161
3162
3163
3164
3165
3166
3167
3168
3169
3170
3171
3172
3173
3174
3175
3176
3177
3178
3179
3180
3181
3182
3183
3184
3185
3186
3187
3188
3189
3190
3191
3192
3193
3194
3195
3196
3197
3198
3199
3200
3201
3202
3203
3204
3205
3206
3207
3208
3209
3210
3211
3212
3213
3214
3215
3216
3217
3218
3219
3220
3221
3222
3223
3224
3225
3226
3227
3228
3229
3230
3231
3232
3233
3234
3235
3236
3237
3238
3239
3240
3241
3242
3243
3244
3245
3246
3247
3248
3249
3250
3251
3252
3253
3254
3255
3256
3257
3258
3259
3260
3261
3262
3263
3264
3265
3266
3267
3268
3269
3270
3271
3272
3273
3274
3275
3276
3277
3278
3279
3280
3281
3282
3283
3284
3285
3286
3287
3288
3289
3290
3291
3292
3293
3294
3295
3296
3297
3298
3299
3300
3301
3302
3303
3304
3305
3306
3307
3308
3309
3310
3311
3312
3313
3314
3315
3316
3317
3318
3319
3320
3321
3322
3323
3324
3325
3326
3327
3328
3329
3330
3331
3332
3333
3334
3335
3336
3337
3338
3339
3340
3341
3342
3343
3344
3345
3346
3347
3348
3349
3350
3351
3352
3353
3354
3355
3356
3357
3358
3359
3360
3361
3362
3363
3364
3365
3366
3367
3368
3369
3370
3371
3372
3373
3374
3375
3376
3377
3378
3379
3380
3381
3382
3383
3384
3385
3386
3387
3388
3389
3390
3391
3392
3393
3394
3395
3396
3397
3398
3399
3400
3401
3402
3403
3404
3405
3406
3407
3408
3409
3410
3411
3412
3413
3414
3415
3416
3417
3418
3419
3420
3421
3422
3423
3424
3425
3426
3427
3428
3429
3430
3431
3432
3433
3434
3435
3436
3437
3438
3439
3440
3441
3442
3443
3444
3445
3446
3447
3448
3449
3450
3451
3452
3453
3454
3455
3456
3457
3458
3459
3460
3461
3462
3463
3464
3465
3466
3467
3468
3469
3470
3471
3472
3473
3474
3475
3476
3477
3478
3479
3480
3481
3482
3483
3484
3485
3486
3487
3488
3489
3490
3491
3492
3493
3494
3495
3496
3497
3498
3499
3500
3501
3502
3503
3504
3505
3506
3507
3508
3509
3510
3511
3512
3513
3514
3515
3516
3517
3518
3519
3520
3521
3522
3523
3524
3525
3526
3527
3528
3529
3530
3531
3532
3533
3534
3535
3536
3537
3538
3539
3540
3541
3542
3543
3544
3545
3546
3547
3548
3549
3550
3551
3552
3553
3554
3555
3556
3557
3558
3559
3560
3561
3562
3563
3564
3565
3566
3567
3568
3569
3570
3571
3572
3573
3574
3575
3576
3577
3578
3579
3580
3581
3582
3583
3584
3585
3586
3587
3588
3589
3590
3591
3592
3593
3594
3595
3596
3597
3598
3599
3600
3601
3602
3603
3604
3605
3606
3607
3608
3609
3610
3611
3612
3613
3614
3615
3616
3617
3618
3619
3620
3621
3622
3623
3624
3625
3626
3627
3628
3629
3630
3631
3632
3633
3634
3635
3636
3637
3638
3639
3640
3641
3642
3643
3644
3645
3646
3647
3648
3649
3650
3651
3652
3653
3654
3655
3656
3657
3658
3659
3660
3661
3662
3663
3664
3665
3666
3667
3668
3669
3670
3671
3672
3673
3674
3675
3676
3677
3678
3679
3680
3681
3682
3683
3684
3685
3686
3687
3688
3689
3690
3691
3692
3693
3694
3695
3696
3697
3698
3699
3700
3701
3702
3703
3704
3705
3706
3707
3708
3709
3710
3711
3712
3713
3714
3715
3716
3717
3718
3719
3720
3721
3722
3723
3724
3725
3726
3727
3728
3729
3730
3731
3732
3733
3734
3735
3736
3737
3738
3739
3740
3741
3742
3743
3744
3745
3746
3747
3748
3749
3750
3751
3752
3753
3754
3755
3756
3757
3758
3759
3760
3761
3762
3763
3764
3765
3766
3767
3768
3769
3770
3771
3772
3773
3774
3775
3776
3777
3778
3779
3780
3781
3782
3783
3784
3785
3786
3787
3788
3789
3790
3791
3792
3793
3794
3795
3796
3797
3798
3799
3800
3801
3802
3803
3804
3805
3806
3807
3808
3809
3810
3811
3812
3813
3814
3815
3816
3817
3818
3819
3820
3821
```