

Evaluating and Improving LLM-based Competitive Program Generation

Minnan Wei^{✉*}, Ziming Li[✉], Xiang Chen^{✉*}, Menglin Zheng[✉], Ziyuan Qu[✉], Cheng Yu[✉], Siyu Chen[✉], Xiaolin Ju[✉]

^aSchool of Artificial Intelligence and Computer Science, Nantong University, Nantong, China

Abstract

Context: Due to the demand for strong algorithmic reasoning, complex logic implementation, and strict adherence to input/output formats and resource constraints, competitive programming generation by large language models (LLMs) is considered the most challenging problem in current LLM-based code generation. However, previous studies often evaluate LLMs using simple prompts and benchmark datasets prone to data leakage. Moreover, prior work has limited consideration of the diversity in algorithm types and difficulty levels.

Objective: In this study, we aim to evaluate and improve LLMs in solving real-world competitive programming problems.

Methods: We initially collect 117 problems from nine regional ICPC/CCPC contests held in 2024 and design four filtering criteria to construct a curated benchmark consisting of 80 problems. Leveraging DeepSeek-R1 as the LLM, we evaluate its competitive program generation capabilities through the online judge (OJ) platforms, guided by a carefully designed basic prompt. For incorrect submissions, we construct a fine-grained error taxonomy and then propose a targeted improvement framework by combining a multi-turn dialogue-based repair phase and an information-augmented regeneration phase.

Results: Experimental results show that only 5 out of 80 problems are fully accepted when using basic prompts. For the unsolved problems, we construct the error taxonomy, including general errors (such as design, boundary, condition, data type, syntax, and input/output errors) and specialized errors (such as those in mathematical problems, greedy algorithms, and graph theories). After applying our proposed improvement strategies, we substantially increased the number of correct solutions, with 46 out of 80 problems successfully accepted.

Conclusion: Our study highlights the current limitations of LLM-based competitive program generation and outlines promising directions for improving the performance.

Keywords: Competitive program generation, Large language model, Prompt engineering, Error taxonomy, Empirical study

1. Introduction

Background. The primary goal of LLM-based code generation [1] is to automate software development by generating functional, efficient, and context-aware code from natural language descriptions or partial inputs. This technology improves developer productivity by reducing manual coding effort and accelerates software prototyping and maintenance. Among various code generation tasks, competitive program generation is regarded as the most challenging task. Unlike traditional software development tasks, competitive programming requires LLMs to demonstrate precise algorithmic reasoning. The LLMs should fully understand the problem description and independently design and implement efficient algorithms. These problems frequently involve complex logic and extensive edge case handling, further increasing the difficulty. Additionally, solutions should strictly conform to specified input/output formats

and adhere to given performance constraints, such as time and memory limits, posing a substantial challenge to the LLM's ability to generate correct and optimized code.

Research Motivation. Although an increasing number of studies [2, 3, 4, 5] have focused on leveraging LLMs for competitive program generation, current research still suffers from the following limitations.

First, existing research on LLM-based competitive program generation often lacks task-specific prompt design. Without prompts tailored to the unique requirements of competitive programming, the generated solutions frequently fall short in terms of correctness, efficiency, or adherence to problem constraints [6, 7].

Second, existing datasets for competitive program generation often suffer from the issue of data leakage [8], as problems or code snippets may have appeared during the pre-training phase of LLMs [9]. This issue is particularly evident in widely used benchmarks such as HumanEval [10] and MBPP [11], which include problems collected from platforms like Codeforces and LeetCode [12]. As a result, it becomes challenging to determine whether LLMs are genuinely solving problems or simply recalling memorized data.

Third, existing datasets rarely incorporate problems from real-world competitive programming contests, resulting in a lim-

*Corresponding author

Email addresses: minnanwei@gmail.com (Minnan Wei[✉]),
dzycd53@gmail.com (Ziming Li[✉]), xchencs@ntu.edu.cn (Xiang Chen[✉]),
2230110478@stmail.ntu.edu.cn (Menglin Zheng[✉]),
quziyanu@gmail.com (Ziyuan Qu[✉]), zydy1227@gmail.com (Cheng Yu[✉]),
chensiuy043@gmail.com (Siyu Chen[✉]), ju.xl@ntu.edu.cn (Xiaolin Ju[✉])

ited diversity of problem types and difficulty levels [13]. Therefore, evaluations of previous studies may fail to comprehensively assess how LLMs perform across the wide range of reasoning and implementation challenges presented by different difficulty levels and problem types.

Benchmark. To fill this gap, we collected real-world competitive programming problems from five Asian regional contests of the International Collegiate Programming Contest (ICPC) and four regional contests of the China Collegiate Programming Contest (CCPC) held in 2024, thereby helping to mitigate the data leakage issue. The initial pool consists of 117 problems. After applying four carefully designed filtering criteria (such as excluding “gold-level” problems, geometry problems with visual images, and problem descriptions with difficult content), we ultimately selected 80 problems that better align with the current capabilities of LLMs. In summary, our constructed benchmark contains 13 warm-up problems, 36 bronze-level problems, and 31 silver-level problems, covering a diverse range of algorithm types and difficulty levels.

Methodology. Our empirical study methodology is illustrated in Figure 1. Specifically, we first generate initial solutions using our designed basic prompt tailored for the competitive algorithm generation task, leveraging DeepSeek-R1 as the LLM. The correctness of the generated programs is evaluated on the VJudge¹ and Codeforces² online judge (OJ) platforms. The results are categorized as Accept (AC), Wrong Answer (WA), Runtime Error (RE), Time Limit Exceeded (TLE), and Compile Error (CE). We then conduct error analysis through LLM-assisted detection and taxonomy-based classification. Finally, we propose a targeted improvement framework based on multi-turn dialogue and information-augmented strategies.

LLM Evaluation. To evaluate the performance of LLM for competitive program generation, we want to answer the following two research questions.

RQ1: Given our designed basic prompt, how correct is the competitive program generated by LLM?

Experimental results show that, under basic prompting, the LLM DeepSeek-R1 fully solves only 5 out of 80 problems. The failures consist of 63 cases of WA, 5 cases of TLE, and 7 cases of CE. In particular, the LLM performs relatively better on warm-up and single-algorithm problems, it struggles significantly with more challenging tasks, such as bronze- and silver-level problems and those requiring multiple algorithmic techniques. These findings highlight the limitations of current LLMs in generating correct and efficient solutions for the competitive program generation task when relying on basic prompts.

RQ2: What types of errors commonly occur in incorrect programs generated by LLM?

We first constructed a comprehensive hierarchical error taxonomy to categorize the common failures of LLM-generated programs systematically. Based on this taxonomy, in the general error category, error analysis reveals that design-related errors (28.6%) and boundary-related Errors (15.5%) are the predominant issues, followed by condition-related errors, data type

errors, syntax errors, and input/output errors. Typical subcategories include incorrect algorithm (12.4%), misunderstanding problem requirements (10.6%), and Incorrect Handling of edge cases or input boundaries (8.7%). In the algorithm-specific error category, we summarize several types of common errors that occur in specific algorithm types, such as mathematical problems and greedy algorithms.

LLM Improvement. To improve LLMs for competitive program generation, we design an improvement framework that combines multi-turn dialogue repair and information-augmented regeneration to enhance program correctness systematically. The framework consists of three phases: **Phase 1** performs error identification and taxonomy classification to guide subsequent repair; **Phase 2** applies targeted, multi-turn dialogue prompts to iteratively revise the faulty code; and **Phase 3** conducts regeneration using structured, problem-specific information when earlier fixes fail. This staged design ensures both precision in fixing localized issues and completeness in regenerating complete solutions.

Improvement Results: After applying our designed improvement framework, the number of AC solutions increased from 5 to 46 out of 80 problems. Specifically, 33 WA cases, 2 TLE cases, and 6 CE cases were successfully accepted.

Impact of Our Study: Our study systematically evaluates the limitations of the current generation of competitive programs based on LLM using a constructed benchmark gathered from recent real-world algorithmic contests. Based on our analysis, we propose potential strategies to improve LLM performance. Promising results can provide practical insights to guide future research in this task.

The main contributions of our study can be summarized as follows:

- **Challenging Benchmark.** We construct a challenging benchmark for competitive program generation using 80 problems from the 2024 ICPC and CCPC regional contests. Compared to existing datasets, our benchmark helps mitigate the data leakage issue. Furthermore, it can provide broad coverage across diverse problem types and difficulty levels.
- **Initial Evaluation.** For the competitive program generation task, we design task-specific basic prompts. Based on the DeepSeek-R1 LLM, we evaluate its performance on our challenging benchmark. The results show that the generation capability is still limited. To better understand the causes of failure, we construct a hierarchical error taxonomy and conduct a detailed analysis for each category.
- **Improvement.** Based on our error analysis, we propose a targeted improvement framework that combines multi-turn dialogue repair with information-augmented regeneration. Evaluation results show a substantial performance improvement, with the number of AC solutions increasing from 5 to 46 out of 80 problems, demonstrating the practical effectiveness of our proposed strategies.

¹<https://vjudge.net/>

²<https://codeforces.com/submissions/xxx>

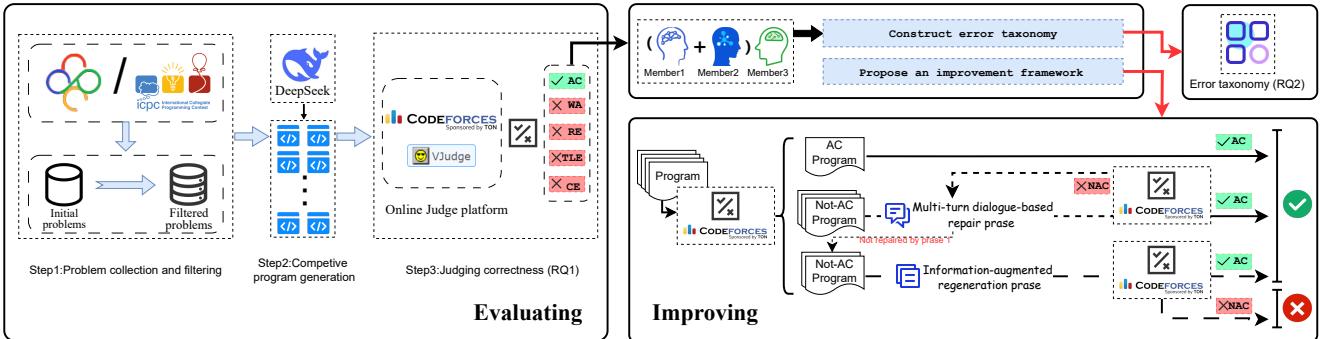


Figure 1: Our empirical study methodology for evaluating and improving LLM-based competitive program generation

Open science. To facilitate reproducibility and encourage other researchers to follow our study, we share the benchmark, scripts, and detailed results on GitHub³.

Paper Organization. The remainder of this paper is organized as follows. Section 1 provides an overview of the study background and motivation. Section 2 describes the benchmark construction and experimental setup. Section 3 formulates the key research questions guiding our study. Section 4 presents the empirical evaluation results and error analysis. Section 5 proposes improvement strategies based on the identified errors. Finally, Section 9 summarizes the key findings and discusses directions for future work.

2. Experimental Setup of LLM Evaluation

In this section, we describe the experimental setup used to evaluate LLM performance on competitive program generation. Specifically, we first show the construction process of our benchmark. Then, we introduce the LLM used in our experiments. Finally, we show the details of the basic prompt tailored for this task.

2.1. Benchmark Construction

In this subsection, we first present the competitive contests used to collect competitive programming problems. Then, we introduce our filtering criteria for selecting problems. Finally, we categorize the selected problems based on their types and difficulty levels.

2.1.1. Competitive Problem Collection

To ensure diversity in problem types and difficulty levels, and to mitigate the potential data leakage issue, we selected regional competitive contests from ICPC and CCPC held in 2024.

For the ICPC, we select the following five Asian regional contests and provide their corresponding URLs:

- **ICPC-CD.** ICPC Asia Chengdu Regional Contest⁴

- **ICPC-NJ.** ICPC Asia Nanjing Regional Contest⁵
- **ICPC-SY.** ICPC Asia Shenyang Regional Contest⁶.
- **ICPC-KMR.** ICPC Asia Kunming Regional Contest⁷.
- **ICPC-KMI.** ICPC Kunming Invitational Contest⁸

For the CCPC, we select the following four regional contests and provide their corresponding URLs:

- **CCPC-HB.** CCPC Harbin Regional Contest⁹
- **CCPC-JN.** CCPC Jina Regional Contest¹⁰
- **CCPC-ZZ.** CCPC Zhengzhou Regional Contest¹¹
- **CCPC-SF.** CCPC Contest Special for Female¹²

2.1.2. Competitive Problem Filtering

To ensure alignment with the current LLM capability, we designed the following filtering rules.

Filter 1: Golden problem. In contests, the “Golden problem” typically refers to the most challenging and complex problem. These problems often involve advanced algorithms, intricate mathematical theories, or require highly innovative problem-solving approaches [3]. Recently, LLMs have encountered significant difficulties in solving such high-difficulty problems due to their limited complex reasoning capabilities and insufficient deep mathematical understanding [13, 14]. Based on these considerations, we filter out this type of problem from our benchmark. For example, in ICPC-CD, problem H (Friendship is Magic)¹³ is a golden problem, and only two competitors can pass all the test cases on Codeforces.

Filter 2: Geometry problems with visual images. Some competitive programming problems, particularly those related

⁵<https://codeforces.com/gym/105484>

⁶<https://codeforces.com/gym/105578>

⁷<https://codeforces.com/gym/105588>

⁸<https://codeforces.com/gym/105386>

⁹<https://codeforces.com/gym/105459>

¹⁰<https://codeforces.com/gym/105540>

¹¹<https://codeforces.com/gym/105632>

¹²<https://codeforces.com/gym/105487>

¹³<https://codeforces.com/gym/105486/problem/H>

Table 1: The number of problems filtered out by each criterion and the final number of selected problems

	ICPC-CD	ICPC-NJ	ICPC-SY	ICPC-KMR	ICPC-KMI	CCPC-HB	CCPC-JN	CCPC-ZZ	CCPC-SF	TOTAL
Initial Number	13	13	13	13	13	13	13	13	13	117
Filter 1	-3	-3	-3	-1	-1	-1	-3	-4	-2	21
Filter 2	/	/	-1	/	/	/	-2	-1	-1	5
Filter 3	-1	/	-2	-2	-1	/	-1	/	-2	9
Filter 4	/	/	/	/	/	-2	/	/	/	2
Final Number	9	10	7	10	11	10	7	8	8	80

to computational geometry, include diagrams or rely heavily on visual representations to convey essential geometric information. Current LLMs struggle to interpret and reason about such visual content, as they primarily process textual data and cannot extract spatial relationships from images. Based on these considerations, we filter out this type of problem from our benchmark. For example, in CCPC-JN, problem G (The Wheel of Fortune)¹⁴ requires competitors to understand the accompanying image before solving the problem, which is also challenging for LLMs.

Filter 3: The problem description contains difficult content to deal with. In algorithmic contests, the problem description may contain a large amount of content that is not suitable for LLM-based competitive program generation, such as excessive tables or overly long background stories. Including such content in the prompt may not only distract the model during the initial generation but also significantly reduce the effectiveness of subsequent strategies in our improvement framework, such as multi-turn dialogue or regeneration phases. Therefore, we filter out this type of problem from our benchmark. For example, in ICPC-KMR, the problem description of the problem K (Key Recovery)¹⁵ is difficult to process and challenging to incorporate effectively into the prompt.

Filter 4: Submission is not possible due to platform restrictions. Some problems on the Codeforces platform are only available via downloadable PDF files. For these problems, the platform does not provide a submission interface due to platform restrictions, making it impossible to submit solutions or obtain evaluation results. Therefore, we exclude such problems from our benchmark. For example, in CCPC-HB, problem B (Concave Hull)¹⁶ and problem M (Weird Ceiling)¹⁷ belong to this category.

We present the filtering results in Table 1, which shows the number of problems removed by each filtering criterion and the final number of selected problems. Specifically, the initial problem pool consists of 117 problems collected from nine competitive programming contests. After applying Filter 1, we removed 21 problems; Filter 2 removed 5 problems; Filter 3 removed 9 problems; and Filter 4 removed 2 problems. Finally, 80 problems were selected to construct our benchmark.

2.1.3. Competitive Problem Classification

In our empirical study, we classify the selected problems from two different perspectives: problem types and problem difficulty levels.

Problem types. By following the suggestion¹⁸, we categorize the problems in our benchmark according to their problem types and report the corresponding proportions, as shown in Figure 2. Note that a single problem may belong to multiple types. For example, in ICPC-CD, problem D (Disrupting Communications)¹⁹, the solution requires the combined use of three algorithmic techniques: graph theory, dynamic programming, and data structures. Next, we provide a brief introduction to each problem type.

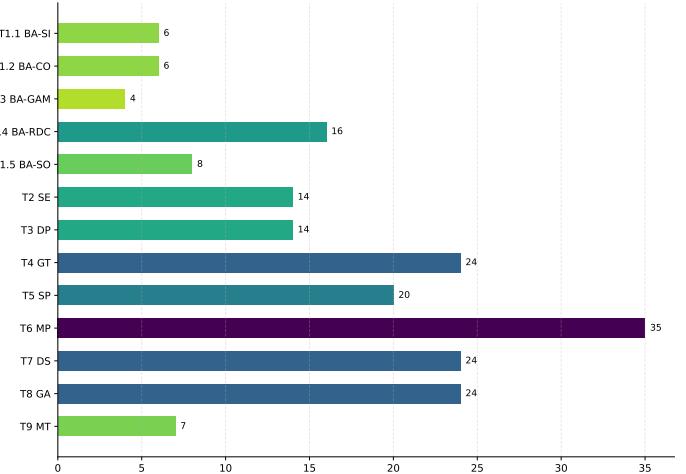


Figure 2: statistics of different problem types

T1 Basic Algorithms (BA). Basic algorithms are widely applicable across various problems and are typically fundamental, concise, and conceptually straightforward.

- **T1.1 Simulation (SI).** Simulates specific processes or operations described in the problem.
- **T1.2 Construction (CO).** Builds outputs directly by exploiting problem patterns or structures.

¹⁴<https://codeforces.com/gym/105540/problem/G>

¹⁵<https://codeforces.com/gym/105588/problem/K>

¹⁶<https://codeforces.com/gym/105459/problem/B>

¹⁷<https://codeforces.com/gym/105459/problem/M>

¹⁸<https://oi-wiki.org/>

¹⁹<https://codeforces.com/gym/105486/problem/E>

- **T1.3 Game (GAM).** Involves decision-making under competitive or adversarial settings, often based on game theory.
- **T1.4 Recursion & Divide-and-Conquer (RDC).** Solves problems by recursively breaking them into smaller parts.
- **T1.5 Sorting (SO).** Rearranges data in a specific order using classic or customized sorting methods.

T2 Search (SE). Explores the solution space to find optimal or valid solutions, often via exhaustive enumeration.

T3 Dynamic Programming (DP). Solves problems by breaking them into overlapping subproblems with optimal substructure.

T4 Graph Theory (GT). Involves problems modeled as graphs, focusing on node and edge relationships.

T5 String Problems (SP). Deals with operations or pattern matching on character sequences.

T6 Mathematical Problems (MP). Requires algorithmic solutions rooted in number theory, combinatorics, or probability.

T7 Data Structure (DS). Involves efficient data organization for fast access, updates, or queries.

T8 Greedy Algorithm (GA). Builds solutions step-by-step by always choosing the locally optimal option.

T9 Miscellaneous Technique (MT). Covers specialized or hard-to-classify techniques.

Problem difficulties. We categorize problem difficulty into four levels: warm-up, bronze, silver, and gold. This classification is based on multiple factors, including Codeforces difficulty ratings²⁰, official post-contest solutions, and the number of successful submissions. To ensure consistency and fairness, we also refer to established mappings between Codeforces ratings and USACO²¹ levels, and cross-check these with the typical distribution of difficulty levels in regional contests.

Table 2 presents the distribution of problem difficulties in our benchmark. Specifically, the benchmark comprises 13 warm-up problems, 36 bronze-level problems, and 31 silver-level problems.

- **Warm-up problems.** Focus on basic programming concepts such as variables, loops, and conditions, intended to help participants quickly get familiar with the contest environment.
- **Bronze-level problems.** Involve fundamental algorithms and data structures such as sorting and simulation. They are suitable for beginners and emphasize accurate implementation.
- **Silver-level problems.** Require more advanced techniques such as graph traversal and prefix sums, focusing on algorithm efficiency and complexity analysis.

- **Gold-level problems.** Involve complex algorithms like dynamic programming and computational geometry, targeting experienced participants with strong problem-solving skills. Notice in Section 2.1.2, gold-level problems are excluded from our benchmark.

2.2. LLM Selection

In our empirical study, we use DeepSeek-R1 as the representative LLM to evaluate competitive program generation performance. We select DeepSeek-R1 due to its strong performance on programming benchmarks, effective instruction-following ability, and cost-efficiency for large-scale inference. According to its official documentation [15], DeepSeek-R1 achieves overall code generation performance comparable to that of ChatGPT, making it a practical choice for this task.

For DeepSeek-R1, we set the temperature hyperparameter to 0.7. This choice is inspired by recent insights into code-specific temperature adaptation presented by Zhu et al. [16], which suggest that a moderate temperature can effectively balance the diversity and accuracy of generated programs.

2.3. Basic Prompt Design

Following the previous study [17], we design a basic prompt tailored to the characteristics of the competitive program generation task. Specifically, our prompt consists of three components: (1) a natural language (NL) description part that includes role-playing instructions and task instructions; (2) a question stem (QS) part that provides the problem background along with input and output specifications; and (3) a test case (TC) part that presents an example input-output pair to illustrate the expected behavior. In the following, we provide a detailed description of each part.

NL Part: To optimize prompt design, this component incorporates two widely adopted types of instructions: (1) role-playing instructions: By assigning the LLM a specific role (e.g., “You are a participant in a competitive contest”), by providing a problem description above and a test case below, the model is encouraged to behave as a contestant and generate a corresponding solution. (2) explicit task instructions: Specifying the desired outcome (e.g. “Please use C++ programming language to solve this competitive problem”) helps guide the model toward the intended task, reducing ambiguity and improving the quality of the generated code.

QS Part: In alignment with prior work on LLM-based code generation [18], this component includes the following elements: (1) a comprehensive introduction to the problem background, including clarifications of any domain-specific terminology; (2) relevant notes or constraints specified within the problem description; (3) a detailed specification of the input format; and (4) a detailed specification of the output format.

TC Part: This component contains corresponding test cases. Test cases play a crucial role in evaluating the correctness of generated programs. They not only verify whether a program meets the expected functional requirements but also help uncover potential bugs for edge cases.

²⁰<https://codeforces.com/blog/entry/78825>

²¹<https://www.vplanetcoding.com/blog/usaco-codeforces>

Table 2: Statistics of different difficulty levels across different competitive contests

	ICPC-CD	ICPC-NJ	ICPC-SY	ICPC-KMR	ICPC-KMI	CCPC-HB	CCPC-JN	CCPC-ZZ	CCPC-SF	TOTAL
Warm-up	1	2	1	2	1	1	2	1	2	13
Bronze-level	4	4	2	6	5	5	3	4	3	36
Silver-level	4	3	4	2	5	5	2	3	3	31

We use an illustrative example shown in Figure 3 to demonstrate our designed basic prompt. This example corresponds to Problem C from the CCPC-SF contest²². Notice we utilize the “deep thinking” feature of DeepSeek when sending the prompt to enhance the reasoning process of the model.

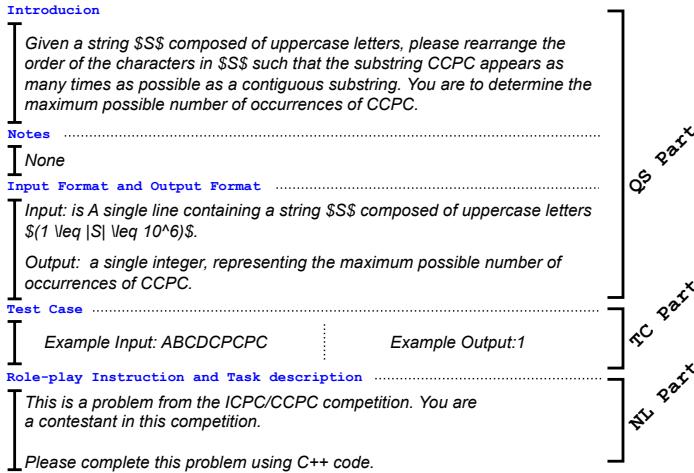


Figure 3: An example to illustrate our basic prompt design

Figure 4 shows the corresponding output results of this illustrative example. The response from DeepSeek is composed of four main sections: deep thinking, approach, C++ code, and explanation. The **deep thinking** section captures the model’s internal reasoning process when interpreting the prompt. The **approach** section outlines the proposed solution strategy. The **C++ code** section presents the generated algorithm implementation. Finally, the **explanation** section provides a detailed interpretation of key code segments.

3. Research Questions

To evaluate the LLM’s ability for competitive problem generation, we want to answer the following two research questions.

RQ1: Given our designed basic prompt, how correct is the competitive program generated by LLM?

Motivation. In RQ1, we aim to assess the basic code generation ability of LLMs on real-world competitive programming problems. Using our designed prompt, we aim to evaluate how well DeepSeek-R1 performs without additional guidance. To reveal strengths and limitations, we analyze results by problem

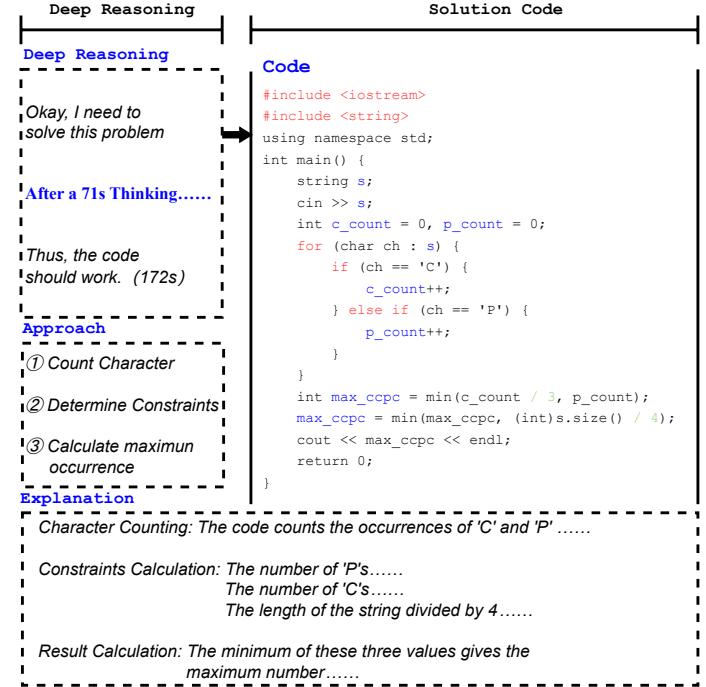


Figure 4: An example to illustrate the output of LLM

difficulty and problem types, offering a clear foundation for understanding LLMs’ capabilities on this task.

RQ2: What types of errors commonly occur in incorrect programs generated by LLM?

Motivation. In RQ2, we aim to conduct an in-depth analysis of the incorrect solutions generated by the LLM to uncover the underlying causes of failure. By systematically examining the failed cases from RQ1, we seek to identify recurring error patterns and classify them into a comprehensive hierarchical error taxonomy. This taxonomy will help to better understand the types and distribution of errors (such as algorithm design flaws, logic mistakes, or performance bottlenecks) that limit the effectiveness of LLMs in competitive program generation. The insights gained from this analysis are crucial for designing the targeted optimization strategies for improvement.

4. LLM Evaluation Results

4.1. RQ1: Correctness Analysis

Approach.

We utilize OJ platforms to assess the correctness of generated programs. This fully automated process enables fair, standardized, and efficient verification of both functional correct-

²²<https://codeforces.com/gym/105487/problem/C>

ness and runtime performance. When a generated program is submitted to an OJ platform, it may return one of the following statuses:

- **AC (Accepted).** The program produces correct outputs for all test cases within the required constraints.
- **WA (Wrong Answer).** The program runs but generates incorrect results for one or more test cases.
- **RE (Runtime Error).** The program crashes during execution, commonly due to issues such as segmentation faults or division by zero.
- **TLE (Time Limit Exceeded).** The program fails to complete execution within the prescribed time limit.
- **CE (Compile Error).** The program fails to compile due to syntax errors or unsupported language feature usage.

Results. To evaluate the effectiveness of DeepSeek-R1 under our designed basic prompt, we submitted the generated solutions for all 80 competitive programming problems in our benchmark. The results show that only 5 solutions achieved AC status, successfully passing all test cases. In contrast, 63 submissions resulted in WA, 5 in TLE, and 7 in CE, indicating that most of the generated programs failed to produce correct or efficient outputs. This distribution highlights the substantial limitations of current LLMs in handling real-world algorithmic challenges when only using the basic prompt. To gain deeper insights, we further analyze the detailed results from two different perspectives: problem difficulty levels and problem types.

Analysis in terms of difficulty levels.

The results by difficulty level are shown in Table 3. Specifically, among the 13 warm-up problems, 4 submissions achieved AC, 1 resulted in TLE, and the remaining 8 were WA. For the 36 bronze-level problems, there were 0 AC, 1 TLE, 3 CE, and 32 WA. Among the 31 Silver-level problems, only 1 achieved AC, while 3 resulted in TLE, 4 in CE, and 23 in WA. These results suggest that DeepSeek-R1 performs better on low-difficulty problems, while its effectiveness drops sharply on more challenging ones, highlighting its limitations in handling high-difficulty problems.

Table 3: Program generation result analysis in terms of problem difficulty level with basic prompt

Difficulty Level	AC	Other Status
Warm-up	4/13	9/13 (8 WA, 1 TLE)
Bronze	0/36	36/36 (32 WA, 1 TLE, 3 CE)
Silver	1/31	30/31 (23 WA, 3 TLE, 4 CE)
Summary	5/80	75/80 (63 WA, 5 TLE, 7 CE)

Analysis in terms of problem types.

Notice to facilitate analysis, we briefly categorize the problems into single-algorithm and multi-algorithm types based on the number of algorithms involved. The results are presented in Table 4. Among the 49 single-algorithm problems, 4 submissions achieved AC, 37 resulted in WA, 4 in TLE, and 37 were

CE. In contrast, for the 31 multi-algorithm problems, only 1 submission achieved AC, with 26 WA, 1 TLE, and 3 CE. These results indicate that LLMs, such as DeepSeek-R1, perform significantly better on single-algorithm problems compared to multi-algorithm ones, suggesting that solving problems requiring the integration of multiple algorithms poses a greater challenge for competitive code generation.

Table 4: Program generation result analysis in terms of problem types (i.e., single-algorithm type and multi-algorithm type) with basic prompt

Problem Type	AC	Other Status
Single algorithm	4/49	45/49 (37 WA, 4 TLE, 4 CE)
Multi-algorithm	1/31	30/31 (26 WA, 1 TLE, 3 CE)
Summary	5/80	75/80 (63 WA, 5 TLE, 7 CE)

Summary for RQ1:

LLMs demonstrate limited effectiveness on competitive programming tasks when using basic prompts, with only 5 out of 80 submissions achieving AC. Submission correctness declines with increasing difficulty, and multi-algorithm problems result in more errors, indicating LLMs struggle with both algorithmic complexity and integration of multiple algorithms.

4.2. RQ2: Error Taxonomy Construction

Approach. To construct a fine-grained taxonomy of errors, we randomly selected 80% of the non-AC (not accepted) programs (i.e., 60 out of 75 incorrect submissions) as our development set. This taxonomy was constructed by following the previous study [19, 20]. The annotation was conducted collaboratively by two annotators (i.e., the fifth and sixth authors), both of whom are top performers from a university algorithm training team. They carefully examined each program’s structure, including function calls, entry points, output behavior, and algorithmic logic, to perform error diagnosis on competitive algorithmic programs generated by LLMs.

We adopted an open coding procedure [21], wherein each annotator reviewed all development set programs at least twice. During this analysis, they identified and labeled all observable errors, summarizing each one using concise and descriptive phrases. Some programs exhibited multiple distinct errors contributing to test failures, for example, a single solution might simultaneously suffer from “misunderstanding problem requirements” and “failure to handle boundary conditions.” For these cases, the program was first debugged and manually corrected into an AC version; suspected faulty regions were then re-inserted and validated incrementally to pinpoint all the errors leading to failure.

The annotators iteratively refined the taxonomy by processing each erroneous program and updating error descriptions to improve clarity and granularity. If a program contained multiple types of errors, each was independently assigned to the corresponding category. In cases of disagreement, a third adjudicator (i.e., the third author), who serves as the leader of

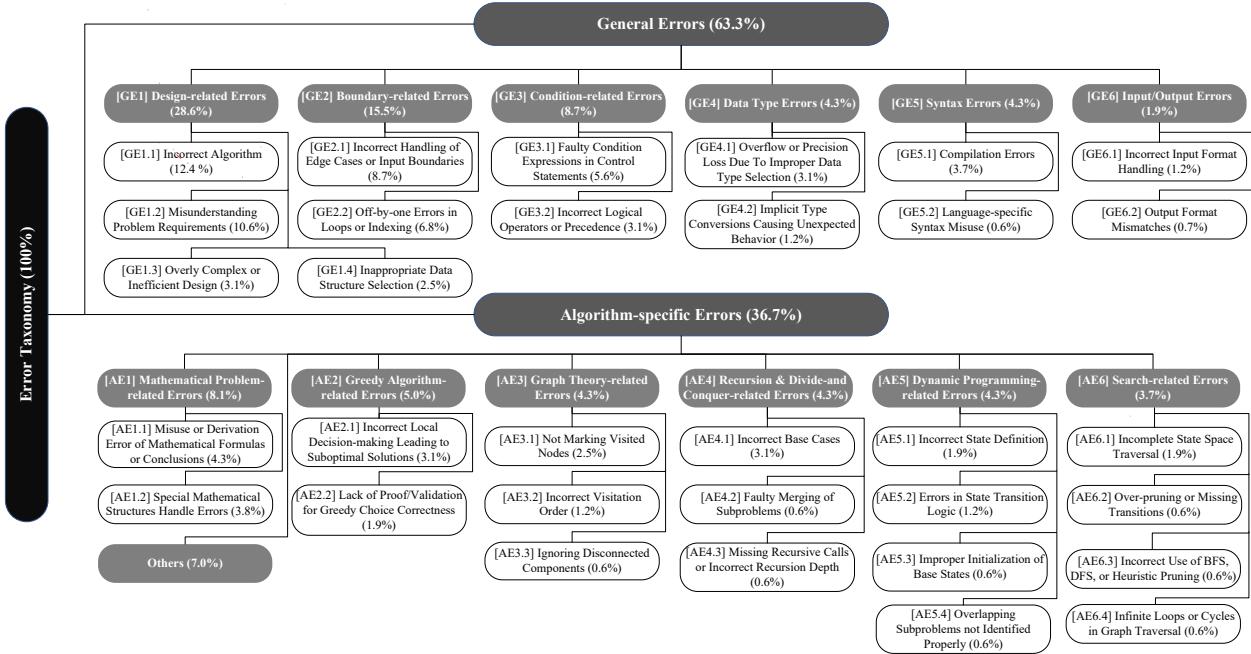


Figure 5: Hierarchical error taxonomy of LLM-generated competitive programs.

the algorithm training team, resolved conflicts through discussion. After thorough inspection and agreement among all participants, we finalized a preliminary error taxonomy covering the full spectrum of LLM-generated failure modes.

Subsequently, we applied the preliminary taxonomy to the remaining 20% of incorrect programs as a validation set. The same two annotators independently labeled these 15 programs by inspecting the code, identifying faults via backfilling, and assigning each error to its corresponding category. If an error could not be classified, it was temporarily labeled as “Pending.” We used Cohen’s Kappa (κ) [22] to evaluate inter-annotator agreement, yielding a score of 0.86, which indicates excellent reliability [23] and confirms the robustness of our open coding procedure.

For “Pending” cases, we brought in a fourth annotator (i.e., the fourth author) to mediate and assist in determining the appropriate categories. If a previously unobserved error type emerged during this step, we added new categories accordingly. In total, six algorithm-specific error categories were introduced to extend the taxonomy. The final taxonomy thus comprehensively captures the diverse error patterns present in LLM-generated competitive programs, with consensus reached among all annotators. Detailed results are reported in Section 4.2.

Result.

We use Figure 5 to present our constructed hierarchical error taxonomy. Specifically, it consists of two main categories: general errors and algorithm-specific errors. For each category, we report its corresponding proportion. Note that a single erroneous program may involve multiple types of errors, and thus can be associated with multiple categories. To better illustrate this, we provide a representative example (Problem K²³ from

CCPC-ZZ) in Figure 6. Then we provide detailed explanations for each category.

General Errors. General errors capture various aspects of program correctness, encompassing design-related errors, syntax errors, input/output errors, boundary-related errors, condition-related errors, and data type errors. Then we provide detailed descriptions of each top-level category along with its corresponding subcategories.

Design-related Errors (GE1). These errors stem from flaws in the overall solution strategy or algorithm design, often caused by a misunderstanding of the problem requirements or the selection of an inappropriate algorithm. Fixing such errors typically requires substantial revisions to the algorithm design, rather than simple code-level modifications. This category consists of four subcategories.

- **Incorrect Algorithm (GE1.1).** This type of error arises when the model selects an inappropriate algorithm for the given problem, leading to a fundamental flaw in the overall solution strategy. A common example is the use of a greedy algorithm for problems that inherently require dynamic programming or exhaustive search. While a greedy approach might pass simple test cases, it often fails on edge cases due to its inability to account for all of the problem’s constraints.
- **Misunderstanding Problem Requirements (GE1.2).** This type of error stems from the model’s incorrect semantic understanding of the problem description. Such misunderstandings typically occur during the initial problem analysis phase, which leads to the subsequent algorithm design in the wrong direction. As a result, even if the implementation is logically consistent, the final program

²³<https://codeforces.com/gym/105459/problem/C>

```

01. #include <iostream>
02. #include <vector>
03. #include <queue>
04. using namespace std;
05.
06. int main() {
07.     int n, k;
08.     double p;
09.     cin >> n >> k >> p;
10.
11.     double inv = 1.0 / (1 - p);
12.     vector<double> e(n + 1);
13.     e[0] = 1.0;
14.     for (int i = 1; i <= n; i++) {
15.         e[i] = e[i - 1] * inv;
16.     }
17.     vector<int> s(n + 1, 0);
18.     priority_queue<pair<double, int>, pq; GE2.1
19.     !
20.     ifor (int i = 1; i <= n; i++) { GE1.1
21.         double marginal = e[i] / ((s[i] + 1) * (s[i] + 2));
22.         pq.push({marginal, i});
23.     }
24.     ifor (int i = 0; i < k; i++) { !
25.         if (pq.empty()) break;
26.         auto top = pq.top();
27.         pq.pop();
28.         int idx = top.second;
29.         s[idx]++;
30.         double new_marginal = e[idx] / ((s[idx] + 1) * (s[idx] + 2));
31.         pq.push({new_marginal, idx});
32.     }
33.     double total = 0.0;
34.     ifor (int i = 1; i <= n; i++) {
35.         total += e[i] / (s[i] + 1);
36.     }
37.     printf("%.10f\n", total); GE1.1
38.     return 0;
39. }

```

GE1.1 : Incorrect Algorithm.
The use of local greedy rather than global optimal DP/mathematical derivation leads to a relatively large overall expectation

GE2.1 : Incorrect Handling of Edge Cases or Input Boundaries
The particular approximation of $n \leq 20$ and $k \leq 200$ was not utilized

Figure 6: An example of the generated solution involving multiple error categories.

fails to fulfill the actual requirements of the task. An illustrative example can be found in Figure 7, which features Problem G from ICPC-KMI²⁴. In the figure, the upper section displays the problem description, while the lower section presents an incorrect code snippet accompanied by a brief analysis of the corresponding error on the side.

Find the lexicographically smallest permutation $p_0, p_1, p_2, \dots, p_{n-1}$ of $0, 1, 2, \dots, (n-1)$ satisfying the following constraint: For all $0 \leq i < n$, $p_0 \oplus p_1 \oplus \dots \oplus p_i > 0$. Here \oplus is the bitwise exclusive or operation.

We say an integer sequence p_0, p_1, \dots, p_{n-1} of length n is lexicographically smaller than another integer sequence q_0, q_1, \dots, q_{n-1} of length n , if there exists an integer $0 \leq k < n$ such that $p_i = q_i$ for all $0 \leq i < k$ and $p_k < q_k$.

```

01. vector<int> permutation;
02. for (int i = 1; i <= n; ++i) { !
03.     if (i % 2 == 1) { 1
04.         permutation.push_back(1);
05.     } else {
06.         permutation.push_back(2);
07.     }
08. }

```

ERROR : did not understand the meaning of "permutation" in the title, the elements do not repeat each other

The 2024 ICPC Kunming Invitational Contest Problem G
Be Positive

Figure 7: An example of GE1.2 error

• Overly Complex or Inefficient Design (GE1.3). This type of error refers to situations where the solver constructs a solution that, while logically correct, incorporates unnecessary layers of computation, overly complex control flow, or inefficient algorithms. Such designs tend

²⁴<https://codeforces.com/gym/105386/problem/G>

to increase time or space complexity, making the solution susceptible to time-limit or memory-limit violations, especially when handling large-scale inputs. An illustrative example can be found in Figure 8, which showcases Problem M from ICPC-KMI²⁵.

BaoBao has prepared a pizza for you! This pizza is a convex polygon with stuffed crust along all of its edges. However, the crust is delicate, allowing you to cut only through its vertices and not into the middle of the edges. Unfortunately, there's a sizable circular piece of pineapple on the pizza that you absolutely want to avoid.

Calculate the largest single piece of pineapple-free pizza that you can obtain with a single straight cut and output its size. A piece is considered pineapple-free when no part of the pineapple falls strictly within it. That is, the area of intersection of the pineapple and the pizza is 0

```

01. ll largestPineappleFreePiece(const vector<Point>& polygon, Point center, ll r) {
02.     int n = polygon.size();
03.     ll maxArea = 0; !
04.     for (int i = 0; i < n; ++i) {
05.         for (int j = i + 1; j < n; ++j) {
06.             if (!isLineIntersection(polygon[i], polygon[j], center, r)) {
07.                 vector<Point> part1, part2;
08.                 for (int k = 0; k < n; ++k) {
09.                     ll cross = (polygon[j] - polygon[i]).cross(polygon[k] - polygon[i]);
10.                     if (cross > 0) part1.push_back(polygon[k]);
11.                     if (cross <= 0) part2.push_back(polygon[k]);
12.                 }
13.                 ll area1 = area(part1); ERROR : Brute force. The correct solution is similar to the rotating jam
14.                 ll area2 = area(part2);
15.                 maxArea = max(maxArea, max(area1, area2));
16.             }
17.         }
18.     }
19.     return maxArea;
20. }

```

2024 ICPC Kunming Invitation Contest Problem M
Italian Cuisine

Figure 8: An example of GE1.3 error

• Inappropriate Data Structure Selection (GE1.4). This type of error arises when the selected data structure does not align with the operational requirements or constraints of the problem. Such a mismatch can lead to low time or space efficiency and may even prevent the correct implementation of the intended algorithm, particularly in problems with strict performance constraints.

Boundary-related Errors (GE2). Improper handling of edge values can compromise the program's logical correctness at the boundaries of the input domain. These errors often remain undetected until exposed by adversarial or hidden test cases.

• Incorrect Handling of Edge Cases or Input Boundaries (GE2.1). This type of error occurs when the program fails to properly handle edge cases or input boundaries, leading to incorrect behavior under extreme or uncommon conditions. An example (Problem A²⁶ from CCPC-JN) can be found in Figure 9.

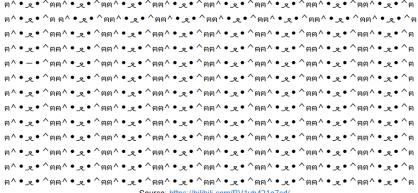
• Off-by-one Errors in Loops or Indexing (GE2.2). This type of error arises when loop boundaries or array indices are set incorrectly, typically due to using $<$ instead of \leq , or starting from an off-position such as 1 instead of 0. These errors often lead to out-of-bounds access, skipped elements, or incomplete processing of edge cases, especially in problems involving arrays, strings, or intervals. An example (Problem A²⁷ from ICPC-KMI) can be found in Figure 10.

²⁵<https://codeforces.com/gym/105386/problem/M>

²⁶<https://codeforces.com/gym/105540/problem/A>

²⁷<https://codeforces.com/gym/105386/problem/A>

Hailing from a long lost tribe of vastaya, Neeko can blend into any crowd by borrowing the appearances of others, even absorbing something of their emotional state to tell friend from foe in an instant. No one is ever sure where — or who — Neeko might be, but those who intend to do her harm will soon witness her true colors revealed, and feel the full power of her primordial spirit magic unleashed upon them.



Source: <https://beta.compet.ly/VtUo421e7cf>

Neeko has infiltrated a grid of characters with dimensions $n \times m$, each cell in the grid is represented by a string of length k . She can mimic the other cells; however, Neeko's cell has at least one character different from the other cells while all the other cells are the same. Your task is to find Neeko.

```

01. if (grid[0][0] == grid[0][1]) {
02.     majority = grid[0][0];
03. } else {
04.     if (grid[0][0] == grid[0][2]) {
05.         majority = grid[0][0];
06.     } else {
07.         majority = grid[0][1];
08.     }
09. }

```

ERROR :
The boundary case n,m starts from 2, while the case m equals 2 is not considered in the error scheme

! 2024 China Collegiate Programming Contest (CCPC) Jinan Site Problem A The Fool

Figure 9: An example of GE2.1 error

Condition-related Errors (GE3). Faulty boolean expressions (such as incorrect operators, misinterpreted precedence, or misplaced negations) can misdirect control flow or prematurely terminate loops. These errors typically result in partially correct outputs or missed corner cases.

- **Faulty Condition Expressions in Control Statements (GE3.1).** This type of error occurs when the conditions within control structures are incorrectly formulated, causing essential execution paths to be skipped or irrelevant code blocks to be executed unintentionally. An example (Problem C²⁸ from CCPC-HB) can be found in Figure 11.

In some regions, people are more accustomed to giving directions using cardinal directions, such as go south to the second intersection, then head east to the first intersection. However, due to the complex road network planning in Harbin, many streets do not align perfectly with cardinal directions. Thus, if you provide directions using absolute directions to someone who has lived in Harbin for a long time, they may struggle to understand your intended route.

In Harbin, people are more accustomed to using relative directions to give guidance. For the same location, a Harbin resident might first instruct you to face south, and then say walk straight along the road to the second intersection, then turn left, and then straight to the first intersection.

To address this difference, you decide to write a program that converts the direction-giving style using cardinal directions into the style preferred by Harbin residents. Of course, using a real map of Harbin would be too complicated, so in this problem, you can assume the map is an infinitely large grid.

```

01. // Function to get the relative direction based on current facing and next direction
02. char getRelativeDirection(char currentFacing, char nextDirection) {
03.     if (currentFacing == 'S') {
04.         if (nextDirection == 'S') return 'Z';
05.         if (nextDirection == 'E') return 'R';
06.         if (nextDirection == 'W') return 'L';
07.     } else if (currentFacing == 'N') {
08.         if (nextDirection == 'N') return 'Z';
09.         if (nextDirection == 'W') return 'R';
10.        if (nextDirection == 'E') return 'L';
11.    } else if (currentFacing == 'E') {
12.        if (nextDirection == 'E') return 'Z';
13.        if (nextDirection == 'N') return 'R';
14.        if (nextDirection == 'S') return 'L';
15.    } else if (currentFacing == 'W') {
16.        if (nextDirection == 'W') return 'Z';
17.        if (nextDirection == 'S') return 'R';
18.        if (nextDirection == 'N') return 'L';
19.    }
20.    return 'Z'; // Default to straight
21. }

```

! 2024 China Collegiate Programming Contest Harbin Onsite Problem C Giving Directions in Harbin

Figure 11: An example of GE3.1 error

• Incorrect Logical Operators or Precedence (GE3.2).

This type of error arises when logical expressions use incorrect operators or depend on operator precedence without proper use of parentheses. As a result, control flow may follow unintended branches, leading to missed conditions or incorrect outputs in test cases.

Data Type Errors (GE4). This error occurs when an inappropriate numeric range or precision is used, or when implicit type conversions alter the intended semantics. Such issues can lead to overflow, underflow, or subtle rounding errors. Selecting an appropriate data type or applying explicit casts can effectively resolve the problem.

- **Overflow or Precision Loss Due To Improper Data Type Selection (GE4.1).** This error arises when a variable is declared with a data type that lacks sufficient range or precision to represent the required values. For instance, using `int` where `long long` is necessary can cause integer overflow and lead to incorrect results.

- **Implicit Type Conversions Causing Unexpected Behavior (GE4.2).** This error occurs when operations involve mixed data types, causing implicit conversions that

²⁸<https://codeforces.com/gym/105459/problem/C>

alter the program's intended semantics. Such conversions can lead to loss of precision, sign mismatches, or unintended value truncation, particularly in arithmetic or conditional expressions.

Syntax Errors (GE5). Purely lexical or grammatical mistakes that prevent compilation or produce immediate runtime faults, such as missing tokens, mismatched delimiters, or language-specific misuse of keywords. They reflect lapses in basic code authoring rather than conceptual misunderstanding.

- **Compilation Errors (GE5.1).** Syntax errors occur when the program fails to compile due to missing semicolons, mismatched brackets, undeclared variables, or other programming language syntax rule violations.
- **Language-specific Syntax Misuse (GE5.2).** This refers to the incorrect use of programming language features that behave differently depending on the language environment. For example, in C++, mixing `cin/cout` with `scanf/printf` without disabling synchronization (via `std::ios::sync_with_stdio(false)`) can lead to unexpected I/O behavior or significant performance degradation. Such misuse does not cause compilation errors but can result in inefficient execution.

Input/Output Errors (GE6). The program mishandles the prescribed I/O contract: it misparses input, formats output incorrectly, or neglects special cases like empty streams, leading to WA status despite correct internal logic.

- **Incorrect Input Format Handling (GE6.1).** This error occurs when the program does not adhere to the input format specified in the problem statement. As a result, it may read input incorrectly or prematurely, leading to incorrect values being processed by the subsequent logic.
- **Output Format Mismatches (GE6.2).** This type of error occurs when the program's output format does not strictly conform to the problem's specifications. Even if the underlying logic and computed values are correct, such formatting issues can still lead to a WA judgment. An example (Problem D²⁹ from ICPC-SY) can be found in Figure 12.

Algorithm-specific Errors. Algorithm-specific errors refer to logic-specific mistakes tied to particular algorithm types (e.g., dynamic programming, graph algorithms, greedy algorithms)

Mathematical Problem-related Errors (AE1). Errors related to mathematical problems.

- **Misuse or Derivation Error of Mathematical Formulas or Conclusions (AE1.1).** These errors occur when the solver applies mathematical formulas incorrectly, derives flawed expressions, or ignores the necessary preconditions for applying known results. Common issues

A permutation of 1 to n in this problem is a 1-based sequence of n integers in which every integer from 1 to n appears exactly once. Alice and Bob are playing a game on $A = [a_1, a_2, \dots, a_n]$ and $B = [b_1, b_2, \dots, b_n]$, two permutations of 1 to n . They take turns to perform operations, with Alice going first, and the one with no possible operation loses.

In each turn, Alice can only operate on the permutation A , while Bob can only operate on the permutation B . A valid operation consists of choosing two indices i and j on the operable permutation ($1 \leq i, j \leq n$) and swapping their corresponding elements, under the constraint that $\sum_{k=1}^n a_k b_k = a_1 b_1 + a_2 b_2 + \dots + a_n b_n$, the dot product of these two permutations, must strictly increase after the swap. Namely, a player loses if he or she cannot perform any swap to increase the dot product, and the other player wins the game.

As Alice and Bob are both clever enough to adopt the best strategy, the winner of such a game can be determined from the start. Therefore, they decide to play the game n times with modifications to make it less boring. Please help them to predict the winners of these n games.

More specifically, two initial permutations $A = [a_1, a_2, \dots, a_n]$, $B = [b_1, b_2, \dots, b_n]$, and $(n - 1)$ modifications $\{[(t_1, l_1, r_1, d_1), (t_2, l_2, r_2, d_2), \dots, (t_{n-1}, l_{n-1}, r_{n-1}, d_{n-1})]\}$ will be given. The first game will start from the given permutations A and B , and for $k = 1, 2, \dots, n - 1$, the $(k + 1)$ -th game will start from the permutations for the k -th game after shifting the interval between indices t_k and r_k of the permutation t_k left d_k times.

Please note that shifting an interval $[p_1, p_2, \dots, p_r]$ left once will give $[p_{t+1}, p_{t+2}, \dots, p_r, p_t]$. For instance, shifting the interval $[2, 4]$ of $[1, 2, 3, 4, 5]$ left once will give $[1, 3, 4, 2, 5]$, and twice will give $[1, 4, 2, 3, 5]$. Besides, shifting the interval $[1, 5]$ of $[3, 1, 4, 5, 2]$ left 4 times will give $[2, 3, 1, 4]$.

```

01. for (int i = 0; i < n - 1; ++i) {
02.     char t;
03.     int l, r, d;
04.     cin >> t >> l >> r >> d;
05.     if (t == 'A') {
06.         shiftLeft(A, l, r, d);
07.     } else {
08.         shiftLeft(B, l, r, d);
09.     }
10.     result += getWinner(A, B);
11. }
12. cout << result << '\n';
13. 
```

The 2024 ICPC Asia Shenyang Regional Contest Problem D Dot Product Game

ERROR : Incorrect output format

`sync_with_stdio(false)`

Figure 12: An example of GE6.2 error

include misusing modular arithmetic identities, applying combinatorial formulas without handling edge conditions, or invoking theorems under invalid assumptions.

- **Special Mathematical Structure Handle Errors (AE1.2).** When dealing with data with special mathematical properties (such as prime sieve, factorization, GCD, modular inverses), the model often makes mistakes in implementation, has logical structure defects, or fails to cover all special cases.

Greedy Algorithm-related Errors (AE2). Errors related to the greedy algorithm.

- **Incorrect Local Decision-making Leading to Suboptimal Solutions (AE2.1).** Greedy algorithms make choices based on local optimality with the hope that this leads to a global optimum. When the local decision rule is flawed, such as picking the largest item without considering constraints, it may yield a feasible but non-optimal solution.

- **Lack of Proof/Validation for Greedy Choice Correctness (AE2.2).** Greedy algorithms require a correctness guarantee (e.g., the property of greedy choice or the optimal substructure). Failing to validate that the greedy strategy works for all inputs formally can lead to solutions that pass sample cases but fail on edge cases, especially where greedy behavior breaks down.

Graph Theory-related Errors (AE3). Errors related to graph-theory algorithms.

- **Not Marking Visited Nodes (AE3.1).** Failing to mark nodes as visited during traversal may lead to infinite loops or repeated visits, particularly in graphs with cycles. This often breaks DFS or BFS logic and affects termination.

²⁹<https://codeforces.com/gym/105578/problem/D>

- **Incorrect Visitation Order (AE3.2).** Marking nodes too late or processing them in the wrong sequence can cause incorrect behavior in traversal-based algorithms (i.e., topological sort, cycle detection). Ensuring proper visitation timing is key to correctness.
- **Ignoring Disconnected Components (AE3.3).** Only exploring from a single start node and neglecting other unvisited nodes can cause entire components to be missed. This results in incomplete outputs for problems involving connectivity or component counting.

Recursion & Divide-and-Conquer-related Errors (AE4).

Errors related to recursion and divide-and-conquer.

- **Incorrect Base Cases (AE4.1).** In divide-and-conquer algorithms, recursion must terminate at well-defined base cases. Errors in base case logic, such as incorrect stopping conditions or failure to handle minimal input sizes, can result in infinite recursion or wrong answers on trivial subproblems.
- **Faulty Merging of Subproblems (AE4.2).** After dividing the problem and solving each part recursively, the merge step is responsible for combining the sub-results into a complete solution. If the merging logic is incorrect, the final output may fail to reflect the correct solution to the original problem accurately.

- **Missing Recursive Calls or Incorrect Recursion Depth (AE4.3).** Divide-and-conquer relies on recursive decomposition of the input. Omitting necessary recursive calls or failing to reach sufficient depth can cause parts of the problem to remain unsolved, leading to incomplete or incorrect results.

Dynamic Programming-related Errors (AE5).

Errors related to the dynamic programming.

- **Incorrect State Definition (AE5.1).** Dynamic programming relies on defining problem states that capture sufficient information for recurrence. If the state is too coarse (i.e., loses key distinctions) or too fine (i.e., leads to overcomplexity), the DP formulation will fail to represent the full problem correctly, resulting in incorrect or incomplete solutions.
- **Errors in State Transition Logic (AE5.2).** State transitions describe how a larger problem is built from smaller subproblems. Mistakes in this recurrence, such as using the wrong indices, conditions, or transition direction, cause the DP table to be filled incorrectly, producing wrong final answers even when the states are defined properly.
- **Improper Initialization of Base States (AE5.3).** Dynamic programming solutions depend on correctly initializing base cases (e.g., $dp[0]$, $dp[1]$) from which all other values are derived. If base values are missing or set incorrectly, subsequent transitions accumulate errors, propagating incorrect values throughout the DP table.

- **Overlapping Subproblems not Identified Properly (AE5.4).**

A key idea of Dynamic programming is recognizing and reusing solutions to overlapping subproblems. If the problem is solved repeatedly for the same input (e.g., in a naive recursive form), it indicates that the overlapping substructure has not been exploited. This often leads to redundant computation and inefficiency.

Search-related Errors (AE6).

Errors related to the search algorithm.

- **Incomplete State Space Traversal (AE6.1).** This error occurs when the search algorithm fails to explore all reachable states due to early termination or limited expansion logic. It leads to missing valid solutions, especially in problems requiring full coverage or global optimality. An example (Problem B³⁰ from CCPC-ZZ) can be found in Figure 13.

Bobo has been playing a puzzle called *Rolling Stones*, which takes place on an equilateral triangular board consisting of n ($n \geq 2$) rows and n^2 cells. Each cell on the board is labeled with a number from 1 to 4. Bobo also has a tetrahedral stone, with each face numbered from 1 to 4 (a tetrahedral die), initially placed at the first cell in the first row of the board. The position of the stone is as follows: the face with the number 1 is towards the left, the face with the number 2 is towards the next row, the face with the number 3 is towards the right, and the face with the number 4 is on the bottom side.

The goal of the puzzle is to roll the stone to a target cell under the following rules:

- **Matching Numbers:** When the stone rests on a cell, the number on the cell must match the number on the stone's bottom face.
- **Single Visit:** Each cell can only be visited once throughout the journey, including the starting and target cells.

The stone rolls by tipping along an edge that touches the board, moving it to a neighboring cell. Given the board layout, the target cell, and the stone's initial orientation, Bobo wants to know: is it possible to reach the target cell following the rules? If possible, what is the minimum number of rolls required to reach the target?

The illustration for a solution of the first sample test is given as follows.

```

01. // Move Left
02. if (j > 1) {
03.     int ni = i;
04.     int nj = j - 1;
05.     ...
06. }
07.
08. // Move right
09. if (j < (2 * i - 1)) {
10.     int ni = i;
11.     int nj = j + 1;
12.     ...
13. }
14.
15. // Move down-Left
16. if (i < n) {
17.     int ni = i + 1;
18.     int nj = j - 1;
19.     ...
20. }
21.

```



ERROR :
 ① misinterpreting the regular triangle chessboard as a quadrilateral chessboard
 ② misinterpreting the position of the chessboard

2024 China Collegiate Programming Contest Zhengzhou
 Onsite Problem B
 Rolling Stones

Figure 13: An example of AE4.1 error

- **Over-pruning or Missing Transitions (AE6.2).** Search algorithms often include pruning to improve efficiency, but overly aggressive or incorrect pruning can remove valid paths. Similarly, neglecting certain transitions in the state graph results in an incomplete or incorrect traversal.

- **Incorrect Use of BFS, DFS, or Heuristic Pruning (AE6.3).** Each search strategy has a specific application scenario.

³⁰<https://codeforces.com/gym/105632/problem/B>

For example, BFS guarantees the shortest path in unweighted graphs, while DFS does not. Applying the wrong strategy or misusing heuristics (e.g., in A*) can produce logically incorrect or suboptimal results.

- **Infinite Loops or Cycles in Graph Traversal (AE6.4).**

Failing to track visited states or cycles can cause the algorithm to revisit the same nodes indefinitely. This often results in non-terminating execution or redundant computation, especially in graphs with cycles or bidirectional edges.

Summary for RQ2:

The error taxonomy system consists of two main categories: general errors (63.3%) and algorithm-specific errors (36.7%). Within general errors, the most prevalent types are design-related errors (28.6%), including incorrect algorithm selection (12.4%) and misunderstanding problem requirements (10.6%), followed by boundary-related errors (15.5%) and condition-related errors (8.7%). These errors typically arise from flawed logic, misunderstanding of the problem constraints, or inadequate control flow structures.

In contrast, algorithm-specific errors are distributed across typical algorithmic paradigms. In particular, mathematical problem-related errors (8.1%) are the most frequent in this category, followed by greedy algorithm errors (5.0%). Other algorithm-specific issues, such as those involving graph theory (4.3%), recursion & divide-and-conquer (4.3%), dynamic programming (4.3%), and search algorithms (3.7%), each contribute modest portions.

a structured taxonomy for guiding subsequent repair. It serves as the foundation of the improvement framework by ensuring that each failure is clearly understood before attempting correction.

We first identify the root cause of failure using a structured inspection process, as illustrated in Figure 15. This process integrates both manual and automated analysis to locate faults at various levels of the program, from high-level algorithm design to low-level implementation details. Key steps include checking functional behavior, boundary conditions, algorithm correctness, and resource efficiency. These checks help uncover common issues such as logical contradictions, format violations, memory errors, and time complexity bottlenecks.

After locating the errors, they are categorized using our fine-grained hierarchical taxonomy, which distinguishes between GE and AE. This classification not only standardizes the understanding of error types but also enables precise mapping to targeted repair strategies in the repair phases.

5.2. Phase 2: Multi-turn Dialogue-based Repair.

This phase aims to iteratively repair non-AC programs by guiding the LLM through a sequence of error-aware prompts. Based on the identified error types, the framework dynamically selects appropriate repair strategies from a modular set of seven options.

In this phase, the LLM is prompted through multiple concise and targeted instructions to correct errors step-by-step. Each repair strategy is designed to address a specific class of failure. The detailed prompt templates and the mapping between error types and repair strategies are provided on GitHub³¹.

- **Strategy 1: Full Algorithm Regeneration** Reconstruct the solution from scratch when the original algorithm paradigm is fundamentally flawed or misaligned with the problem requirements.
- **Strategy 2: Logic Completion** Fill in missing steps or incomplete branches in control logic, especially in conditionals or loops.
- **Strategy 3: Prompt-Level Clarification** Clarify ambiguous problem descriptions or constraints at the prompt level to better align model understanding with problem intent.
- **Strategy 4: I/O Format Fix** Standardize input/output handling to conform with expected formats (e.g., newline characters, spacing, or value delimiters).
- **Strategy 5: Syntax-Level Repair** Resolve compilation or syntax issues such as missing semicolons, incorrect function declarations, or unmatched brackets.
- **Strategy 6: Calculation Precision / Memory Safety Enhancement** Address issues related to integer overflow, floating-point errors, or unsafe memory access patterns.

5. Improvement Framework

The evaluation results in Section 4 reveal that most LLM-generated programs fail due to algorithm design flaws, logic errors, and improper handling of edge cases. These common failure patterns highlight the need for a more systematic approach to enhance model correctness. Motivated by these findings, we simulate the reasoning process of contestants and propose a taxonomy-driven improvement framework for LLM-generated competitive programs. The overall framework is shown in Figure 14, consisting of three main phases: Specifically, **Phase 1** identifies the root causes of failure and classifies errors based on a hierarchical taxonomy. **Phase 2** uses multi-turn dialogue with targeted prompts to iteratively fix the code. If the program is still incorrect, **Phase 3** regenerates a new solution using a structured prompt that incorporates key problem insights. Together, these phases enable systematic and effective improvement of LLM-generated competitive programs.

5.1. Phase 1: Error Diagnosis

This phase aims to diagnose why a basic prompt-generated program fails to achieve AC and to map the observed errors into

³¹<https://github.com/minnanWei/LLMs-Competitive-Program-Generation/blob/main/README.md?plain=1>

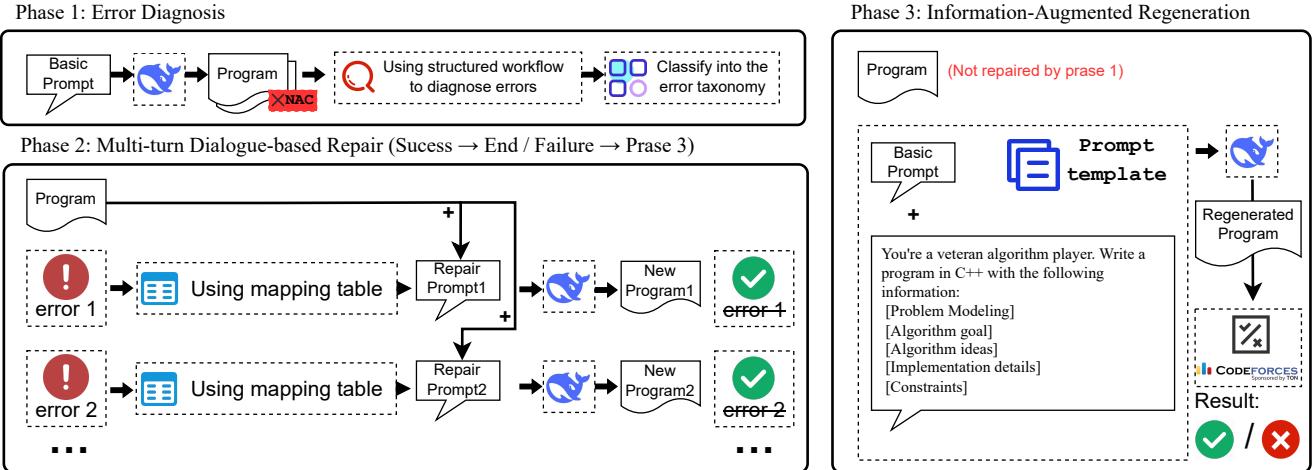


Figure 14: Improvement framework for LLM-based competitive program generation

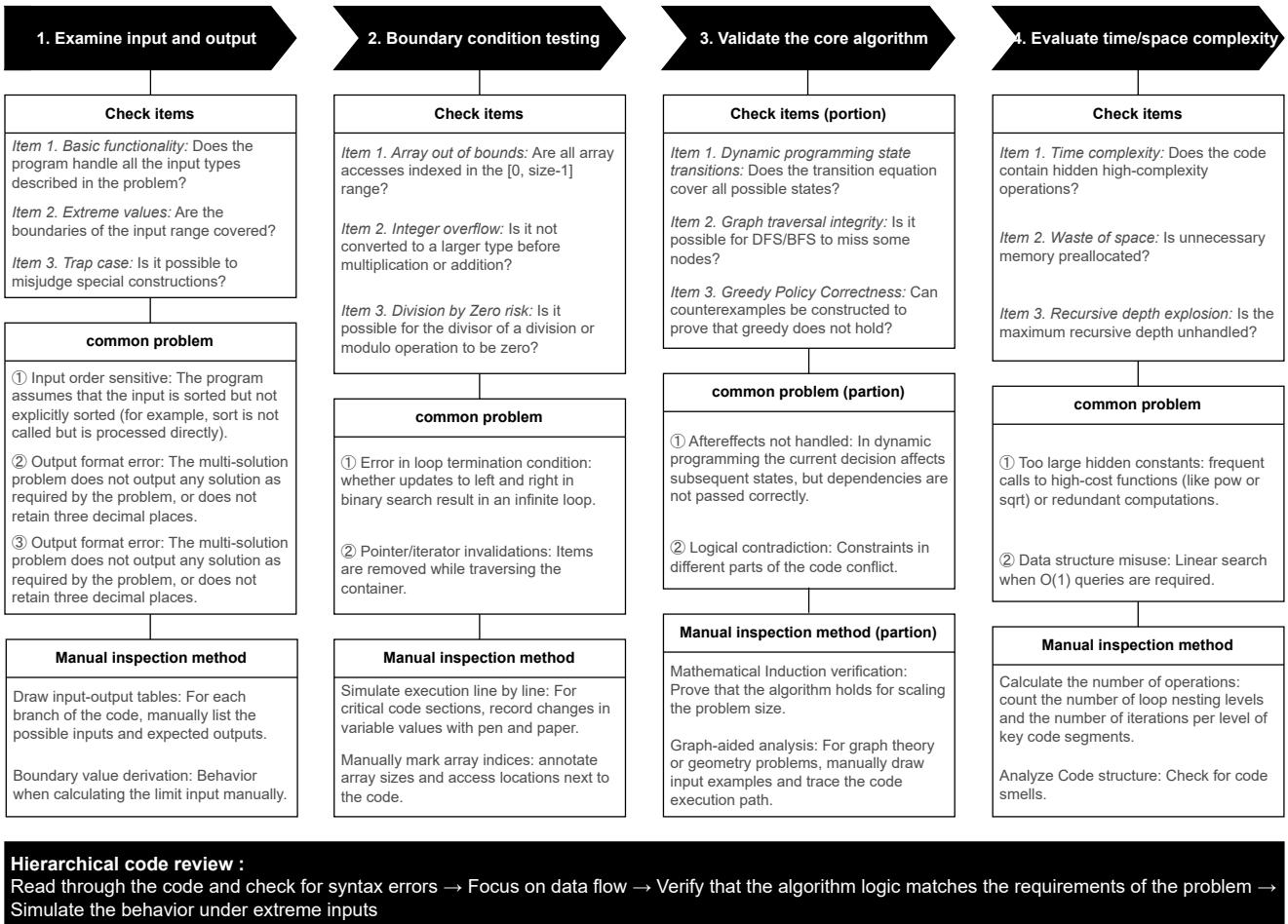


Figure 15: A structured workflow for identifying and diagnosing errors in non-AC programs.

- **Strategy 7: Termination Assurance** Modify loop boundaries or recursive calls to ensure the program halts under all input conditions.

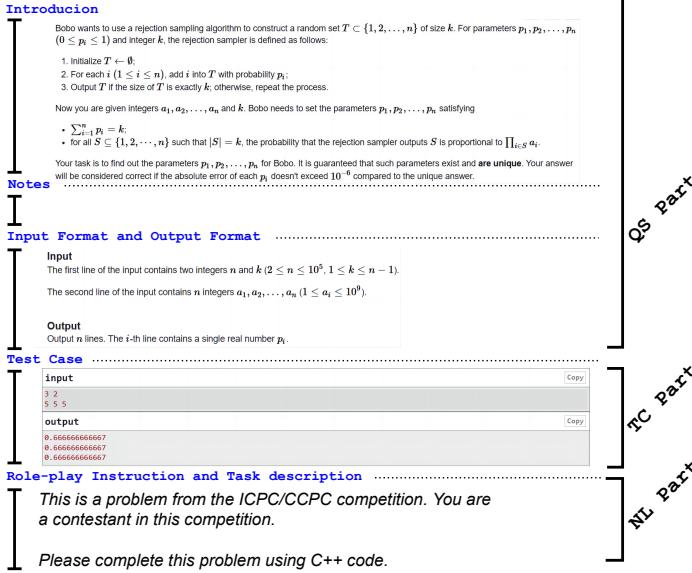


Figure 16: Basic prompt of 2024 China Collegiate Programming Contest (CCPC) Zhengzhou Onsite, problem M (Rolling Stones).

An example (Problem M³² from CCPC-ZZ) is used to demonstrate the repair process. The basic prompt for this example can be found in Figure 16, which serves as the foundation for subsequent multi-turn repair. As shown in Figure 17, the dialogue-based repair phrase progressively enhances the correctness of LLM-generated programs: it first targets higher-level logic errors (GE2.2 – incorrect condition judgements) through mathematical correction and algorithm-refinement techniques, and then resolves lower-level implementation issues (GE3.3 - data-type and precision errors) via precision-enhancement and memory safety strategies. This staged workflow systematically improves both algorithm design and implementation quality, enabling the competitive program to evolve from an initial WA5 status, to CE in the second round, and finally to AC after all test cases pass.

5.3. Phase 3: Information-Augmented Regeneration.

Programs that remain incorrect after phase 1 are regenerated from scratch with a prompt template. Following Liu et al. [24], lengthy conversations can induce forgetting and hallucination; therefore, we inject structured, task-specific scaffolding (i.e., Problem Modeling, Algorithm Objectives, Algorithm Ideas, Implementation Details, and Constraints) into the new prompt (prompt template is shown on Github³³). This holistic context enables the LLM to design a fresh, coherent solution instead of patching isolated defects.

Give an example to demonstrate the information-augmented regeneration prompt template: CCPC-FM, Problem F (Perfect Square)³⁴. The program generated by the LLM based on the basic prompt produced an incorrect result on the first test case. However, after applying the basic prompt combined with the information-augmented regeneration prompt template (shown in Figure 18), the regenerated program successfully passed all test cases and achieved an AC state.

6. Improvement Results

After applying our improvement framework, the number of fully correct (AC) solutions increases from 5 to 46 out of 80 problems, while significantly reducing incorrect outputs (i.e., WA decreased from 63 to 30, TLE from 5 to 3, and CE from 7 to 1). The complete transition of program statuses across the two phases is illustrated in the Sankey diagram, as shown in Figure 19. These results demonstrate the effectiveness of structured error categorization and targeted repair in enhancing LLM-based competitive programming capabilities. Notably, the proposed approach proves effective not only across different difficulty levels but also generalizes well to problems requiring either single or multiple algorithms. The improvement process followed a phased strategy: Phase 2 focused on correcting problems with relatively few errors and simple fixes (i.e., cases where the core logic was sound but minor implementation issues remained). In contrast, Phase 3 addressed more complex cases that Phase 2 could not resolve (i.e., programs with multiple interdependent errors or requiring substantial restructuring).

Analysis in terms of difficulty levels. Shown in Figure 20, after applying the taxonomy-driven improvement framework, the generation results changed as follows. For warm-up problems, the results improved from 4 AC, 8 WA, 1 TLE under the basic prompt to 8 AC, 5 WA after Phase 2, and finally to 11 AC, 2 WA after Phase 3. For bronze-level problems, the outcomes progressed from 0 AC, 32 WA, 1 TLE, and 3 CE to 17 AC, 19 WA after Phase 2, and ultimately to 27 AC, 9 WA. For silver-level problems, the improvement trajectory was from 1 AC, 23 WA, 3 TLE, and 4 CE to 4 AC, 22 WA, 3 TLE, 2 CE after Phase 2, and then to 8 AC, 19 WA, 3 TLE, 1 CE after Phase 3.

Analysis in terms of the problem types. Shown in Figure 21, after applying the taxonomy-driven improvement framework, the generation results also improved significantly in terms of problem types. For single-algorithm-type problems, the results improved from 4 AC, 37 WA, 4 TLE, and 4 CE under the basic prompt to 24 AC, 23 WA, 1 TLE, and 1 CE after Phase 2, and further to 33 AC, 13 WA, 2 TLE, and 1 CE after Phase 3. For multi-algorithm-type problems, the results progressed from 1 AC, 26 WA, 1 TLE, and 3 CE to 6 AC, 24 WA, 1 TLE after Phase 2, and ultimately to 13 AC, 17 WA, 1 TLE after Phase 3.

³⁴<https://codeforces.com/gym/105487/problem/F>

³²<https://codeforces.com/gym/105632/problem/M>

³³<https://github.com/minnanWei/LLMs-Competitive-Program-Generation/blob/main/README.md?plain=1>

```

Strategy II (Logic Completion)
1. Derive the expression of p_i correctly.
2. Implement the bisection method to solve c to ensure the correct judgment of monotonicity.
3. Reasonably select the bisection range and the number of iterations to ensure accuracy.

Round 1 improving
Handling GR2.2: Incorrect Judgment of Conditions
N_Accuracy C++17 (GCC 7-32) Wrong answer on test 5 281 ms 100 KB
1. double low = 0.8L, high = 1e20;
2. for (int iter = 0; iter < 100; ++iter) {
3.     double mid = (low + high) * 0.5;
4.     double sum = 0.0;
5.     for (int i = 0; i < n; ++i) {
6.         if (product > 1e15) {
7.             if (product > 1e15) {
8.                 if (product > 1e15) {
9.                     if (product > 1e15) {
10.                        if (product > 1e15) {
11.                            if (product > 1e15) {
12.                                if (product > 1e15) {
13.                                    if (product > 1e15) {
14.                                        if (product > 1e15) {
15.                                            if (product > 1e15) {
16.                                                if (product > 1e15) {
17.                                                    if (product > 1e15) {
18.                                                        if (product > 1e15) {
19.                                                            if (product > 1e15) {
20.                                                                if (product > 1e15) {
21.                                                                    if (product > 1e15) {
22.                                                                        if (product > 1e15) {
23.                                                                            if (product > 1e15) {
24.                                                                                if (product > 1e15) {
25.                                                                                    if (product > 1e15) {
26.                                                                                        if (product > 1e15) {
27.                                                                                            if (product > 1e15) {
28.                                                                                                if (product > 1e15) {
29.                                                                if (product > 1e15) {
30.                                                                if (product > 1e15) {
31.                                                                if (product > 1e15) {
32.                                                                if (product > 1e15) {
33.                                                                if (product > 1e15) {
34.                                                                if (product > 1e15) {
35.                                                                if (product > 1e15) {
36.                                                                if (product > 1e15) {
37.                                                                if (product > 1e15) {
38.                                                                if (product > 1e15) {
39.                                                                if (product > 1e15) {
40.                                                                if (product > 1e15) {
41. } }
```

GR2.2 - Lambda averaging leads to illegal interval shaping and affects correctness

```

Strategy VI: Calculation Precision/Memory Safety Enhancement
Your code crashes due to calculation precision.
Below is the failed test case:
Input:
4 2 3.4
Correct output:
Your output:
0.31003569758
0.31003569758
0.574115
0.574114879981
0.642525378623

Round 2 improving
Handling GR3.3: Data type and precision error
N_Accuracy C++20 (GCC 13-64) Accepted 0 ms 0 KB
1. long double low = 0.8L, high = 1e15L;
2. for (int iter = 0; iter < 200; ++iter) {
3.     long double mid = (low + high) * 0.5L;
4.     long double sum = 0.0L;
5.     for (int i = 0; i < n; ++i) {
6.         sum += (mid * a[i]) / (1.0L + mid * a[i]);
7.     }
8.     if (sum < k) {
9.         low = mid;
10.    } else {
11.        high = mid;
12.    }
13. }
```

GR3.3 - Data accuracy does not meet the requirement

```

Strategy VI: Calculation Precision/Memory Safety Enhancement
Your code crashes due to calculation precision.
Below is the failed test case:
Input:
4 2 3.4
Correct output:
Your output:
0.31003569758
0.31003569758
0.574115
0.574114879981
0.642525378623

Round 2 improving
Handling GR3.3: Data type and precision error
N_Accuracy C++20 (GCC 13-64) Accepted 0 ms 0 KB
1. long double low = 0.8L, high = 1e15L;
2. for (int iter = 0; iter < 200; ++iter) {
3.     long double mid = (low + high) * 0.5L;
4.     long double sum = 0.0L;
5.     for (int i = 0; i < n; ++i) {
6.         sum += (mid * a[i]) / (1.0L + mid * a[i]);
7.     }
8.     if (sum < k) {
9.         low = mid;
10.    } else {
11.        high = mid;
12.    }
13. }
```

GR3.3 - Data accuracy does not meet the requirement

Figure 17: A multi-turn dialogue-based repair strategy example.

You are a veteran algorithm player. Write a program in C++ with the following information:

[Problem Modeling]:
Given n positive integers a_i (each $\leq 10^6$), select a divisor d_i for each a_i such that the product $\prod_{i=1}^n d_i$ is a perfect square y^2 . Compute the sum of all such valid y modulo $10^9 + 7$.

[Algorithm goal]:
Time Complexity: $O(n \log a_i)$
Space Complexity: $O(n + 10^6)$

[Algorithm ideas]:
1. Use a linear sieve to get minimum prime factors. Factorize each a_i.
2. For each prime p, maintain two DP states:
- f[0]: sum of y when the total exponent of p is even
- f[1]: sum of y when the total exponent is odd
Precompute:
- s0 = sum of p^k for even k
- s1 = sum of p^k for odd k
Transitions:
- new_f0 = f0 * s0 + f1 * s1 * p
- new_f1 = f1 * s0 + f0 * s1
3. Final answer = product of all f[0] values modulo $10^9 + 7$

[Implementation details]:
- Use sieve up to 10^6 for fast factorization
- For each a_i, factorize using the sieve
- For any remaining large prime x > 1, treat separately
- Handles edge cases: n = 1, all a_i = 1, all primes, or mixed values

[Constraints]:
\$1 \leq n \leq 10^{65}\$, \$1 \leq a_i \leq 10^{68}\$
Output: sum of all valid \$y \bmod 10^9 + 7\$
Covers boundary cases: \$n = 1\$, all \$a_i = 1\$, all \$a_i\$ are primes, and mixed values
Input is guaranteed to be valid

Figure 18: A prompt template example for CCP-C-FM, Problem F.

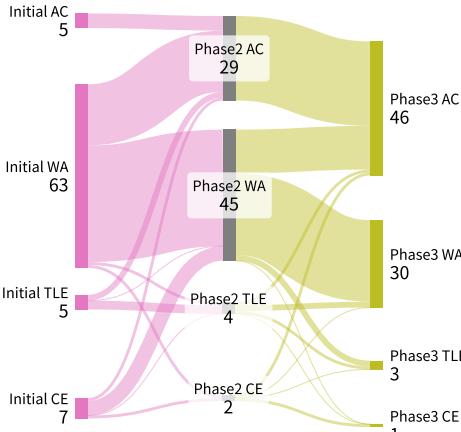


Figure 19: Comparison of program generation results before and after improvement

7. Threats to Validity

7.1. Internal Threats

The first internal threat concerns the design of the basic prompts used in our evaluation. To mitigate this threat, we followed established practices from prior studies [25, 26, 27] and carefully crafted our basic prompts to align with the specific characteristics and requirements of competitive program generation tasks.

The second internal threat refers to the quality of the error taxonomy. Since the labeling process is performed manually, it can introduce subjectivity and inconsistency [28]. To mitigate this risk, we followed established taxonomy construction methodologies [17], where annotators independently labeled each instance. Disagreements were resolved through discussion, with a third party introduced to mediate and ensure consistency.

The final internal threat lies in the evaluation of program correctness. In our empirical study, we primarily rely on online judge platforms for validation. Typically, each problem is accompanied by comprehensive test cases, which not only cover normal functional scenarios but also include various boundary cases.

7.2. External Threats

The first external threat relates to the construction of our benchmark. To mitigate the risk of data leakage, we selected 80 problems from the 2024 ICPC and CCPC contests, which were released after the pre-training cutoff date of DeepSeek [29]. In addition, choosing problems from these contests helps ensure broader coverage across diverse problem types and difficulty levels.

The second external threat stems from the choice of the LLM. We selected DeepSeek-R1 due to its strong performance in code generation tasks and its accessibility as a state-of-the-art open-source model [30].

The final external threat concerns the choice of programming language during code generation. In our empirical study, we primarily focused on C++, as it is the dominant language in algorithmic competitions due to its high execution efficiency,

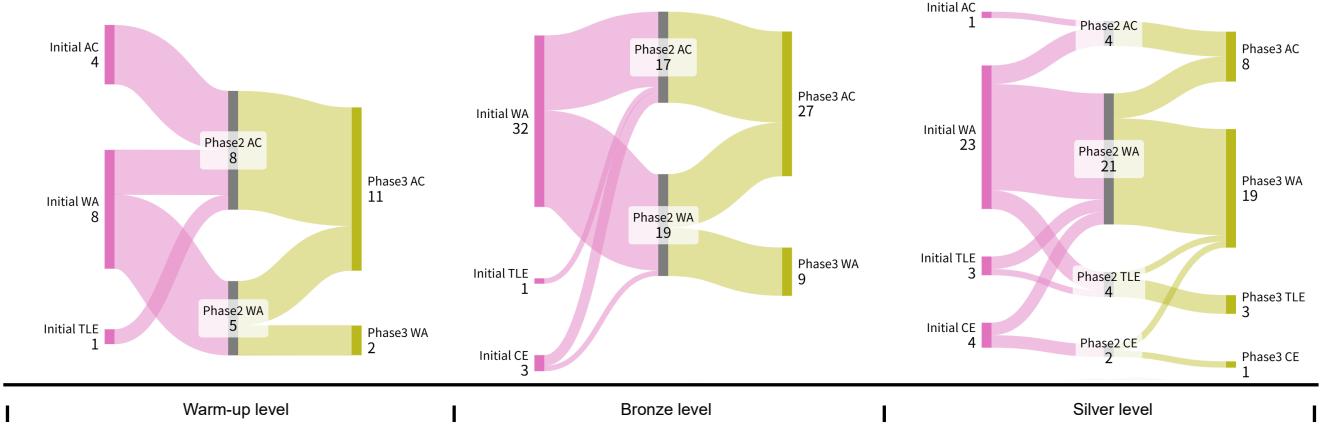


Figure 20: Comparison of program generation results before and after improvement in terms of difficulty levels.

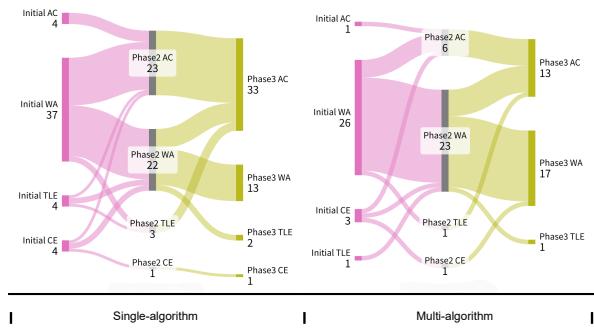


Figure 21: Comparison of program generation results before and after improvement in terms of problem types

rich standard library (e.g., STL), and widespread adoption in platforms like ICPC and Codeforces.

8. Related Work

Competitive program generation focuses on automatically generating solutions to algorithmic problems from natural language descriptions. Compared with general-purpose coding tasks, this setting imposes stricter demands on algorithmic reasoning, correctness, and compliance with time/memory constraints and input/output specifications. As a result, it has become a critical benchmark for assessing the reasoning and planning abilities of LLMs.

Benchmarks and datasets. Early evaluations of program synthesis relied on datasets like HumanEval [10] and MBPP [11], which focus on small-scale functional problems and often lack coverage of complex algorithmic logic. More recent datasets like APPS [31], CodeContests [32] aim to better reflect competitive programming environments. However, many of these benchmarks suffer from data leakage due to overlap with training corpora, limited algorithm type diversity, and insufficient stratification across difficulty levels. This hinders the reliable evaluation of model generalization and robustness.

Algorithmic code generation methods. The release of AlphaCode[32] marked a significant milestone by generating

competitive programming solutions using a sampling-and-filtering pipeline. Follow-up work, such as **Self-Edit**[33] adopted runtime feedback and automatic correction to improve faulty outputs, while **AlphaCodium**[7] proposed test-driven generation combined with iterative regeneration. Recent frameworks like **MapCoder**[34] introduce a multi-agent structure for retrieval, planning, coding, and self-debugging. Lightweight models [35] and instruction-tuned systems like **Magicoder**[5] further explore efficiency and domain specialization. In terms of evaluation methodology, **ProBench**[36] presents an online-judge-based evaluation pipeline with fine-grained feedback, but lacks structured analysis of error types and model failure modes.

Motivation and novelty of our study. Despite recent progress, current datasets and evaluation settings still suffer from key limitations: (i) potential data contamination due to overlaps with training data, (ii) incomplete coverage of algorithm types (e.g., greedy, DP, graph), and (iii) lack of stratification across difficulty levels. To address these issues, we construct a new benchmark by collecting and filtering problems from nine official ICPC and CCPC contests held in 2024. We then evaluate the latest LLM **DeepSeek-R1** on this benchmark using carefully designed basic prompts. The results show that only 5 out of 80 problems are fully solved. Motivated by this finding, we perform a systematic manual error analysis and construct a hierarchical error taxonomy to categorize common failure patterns. Finally, based on this taxonomy, we propose a **taxonomy-driven improvement framework** that combines multi-turn repair and information-augmented regeneration. Experiments show that this approach significantly improves solution correctness, increasing the number of accepted programs from 5 to 46, and providing concrete guidance for future work in LLM-based algorithmic code generation.

9. Conclusion

Our study addresses critical shortcomings in existing competitive programming datasets, such as data leakage and limited diversity in problem types and difficulty. To overcome these challenges, we developed a comprehensive benchmark sourced from recent ICPC and CCPC contests. Our evaluation

of the DeepSeek-R1 model on this benchmark revealed significant performance limitations, motivating an in-depth error analysis and the creation of a detailed error taxonomy. Building upon these insights, we introduced an improvement framework that substantially enhances code generation accuracy across a broader spectrum of problems. These findings not only highlight current gaps in LLM capabilities but also offer a promising foundation for advancing future research in automated competitive programming code generation.

In future work, we will continue to expand our benchmark by incorporating problems from top-tier competitions such as Codeforces, the ICPC series, and the IOI series. We also plan to evaluate more advanced LLMs, such as GPT-4, Claude, and Gemini. Furthermore, we aim to design and implement more effective optimization strategies specifically tailored to address the most challenging problems.

Acknowledgments

This research was partially supported by the National Natural Science Foundation of China (Grant no. 61202006) and the Postgraduate Research & Practice Innovation Program of Jiangsu Province (Grant no. SJCX24_2022).

Declaration of Competing Interests

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

CRediT Authorship Contribution Statement

Minnan Wei: Data curation, Writing -review & editing, Software, Validation. **Ziming Li:** Data curation, Software, Validation. **Xiang Chen:** Software, Conceptualization, Methodology, Writing -review & editing, Supervision. **Menglin Zheng:** Data curation, Writing-review & editing, Validation. **Ziyang Qu:** Data curation, Software, Validation. **Cheng Yu:** Data curation, Software, Validation. **Siyu Chen:** Data curation, Software, Validation. **Xiaolin Ju:** Methodology, Writing -review & editing.

References

- [1] J. Jiang, F. Wang, J. Shen, S. Kim, S. Kim, A survey on large language models for code generation, arXiv preprint arXiv:2406.00515 (2024).
- [2] J. Zhang, D. Li, J. C. Kolesar, H. Shi, R. Piskac, Automated feedback generation for competition-level code, in: Proceedings of the 37th IEEE/ACM International Conference on Automated Software Engineering, 2022, pp. 1–13.
- [3] M. S. Hossain, A. Tabassum, M. F. Arefin, T. S. Zaman, Llm-pros: Analyzing large language models' performance in competitive problem solving, arXiv preprint arXiv:2502.04355 (2025).
- [4] S. Li, F. Yang, J. Liu, Y. Li, Z. Liu, H. Sun, Y. Wang, Y. Chen, Y. Fu, W. Shi, Z. Liu, W. Tang, Y. Shen, B. Tang, M. Ding, J. Tang, J. Song, Competition-level code generation with alphacontest, arXiv preprint arXiv:2312.05820 (2023).
- [5] Y. Lu, Y. Zhang, G. Li, S. Huang, Z. Liu, Y. Chen, X. Han, Y. Xu, F. Wei, Magicoder: The rise of generative ai for code understanding and generation, arXiv preprint arXiv:2403.05530 (2024).
- [6] T. Wang, N. Zhou, Z. Chen, Enhancing computer programming education with llms: A study on effective prompt engineering for python code generation, arXiv preprint arXiv:2407.05437 (2024).
- [7] T. Ridnik, D. Kredo, I. Friedman, Code generation with alphacodium: From prompt engineering to flow engineering, arXiv preprint arXiv:2401.08500 (2024).
- [8] Z. Wang, M. Shao, J. Bhandari, L. Mankali, R. Karri, O. Sinanoglu, M. Shafique, J. Knechtel, Vericontaminated: Assessing llm-driven verilog coding for data contamination, arXiv preprint arXiv:2503.13572 (2025).
- [9] A. Matton, T. Sherborne, D. Aumiller, E. Tommasone, M. Alizadeh, J. He, R. Ma, M. Voisin, E. Gilsenan-McMahon, M. Gallé, On leakage of code generation evaluation datasets, arXiv preprint arXiv:2407.07565 (2024).
- [10] M. Chen, J. Tworek, H. Jun, Q. Yuan, H. P. D. O. Pinto, J. Kaplan, H. Edwards, Y. Burda, N. Joseph, G. Brockman, et al., Evaluating large language models trained on code, arXiv preprint arXiv:2107.03374 (2021).
- [11] J. Austin, A. Odena, M. Nye, M. Bosma, H. Michalewski, D. Dohan, E. Jiang, C. Cai, M. Terry, Q. Le, et al., Program synthesis with large language models, arXiv preprint arXiv:2108.07732 (2021).
- [12] M. Riddell, A. Ni, A. Cohan, Quantifying contamination in evaluating code generation capabilities of language models, arXiv preprint arXiv:2403.04811 (2024).
- [13] Y. Huang, Z. Lin, X. Liu, Y. Gong, S. Lu, F. Lei, Y. Liang, Y. Shen, C. Lin, N. Duan, et al., Competition-level problems are effective llm evaluators, arXiv preprint arXiv:2312.02143 (2023).
- [14] Z. Zheng, Z. Cheng, Z. Shen, S. Zhou, K. Liu, H. He, D. Li, S. Wei, H. Hao, J. Yao, et al., Livecodebench pro: How do olympiad medalists judge llms in competitive programming?, arXiv preprint arXiv:2506.11928 (2025).
- [15] D. Guo, D. Yang, H. Zhang, J. Song, R. Zhang, R. Xu, Q. Zhu, S. Ma, P. Wang, X. Bi, et al., Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, arXiv preprint arXiv:2501.12948 (2025).
- [16] Y. Zhu, J. Li, G. Li, Y. Zhao, Z. Jin, H. Mei, Hot or cold? adaptive temperature sampling for code generation with large language models, in: Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 38, 2024, pp. 437–445.
- [17] Z. Liu, Y. Tang, X. Luo, Y. Zhou, L. F. Zhang, No need to lift a finger anymore? assessing the quality of code generation by chatgpt, IEEE Transactions on Software Engineering (2024) 1–35.
- [18] C. Liu, X. Bao, H. Zhang, N. Zhang, H. Hu, X. Zhang, M. Yan, Guiding chatgpt for better code generation: An empirical study, in: 2024 IEEE International Conference on Software Analysis, Evolution and Reengineering (SANER), IEEE, 2024, pp. 102–113.
- [19] J. Wen, Z. Chen, Y. Liu, Y. Lou, Y. Ma, G. Huang, X. Jin, X. Liu, An empirical study on challenges of application development in serverless computing, in: Proceedings of the 29th ACM joint meeting on European software engineering conference and symposium on the foundations of software engineering, 2021, pp. 416–428.
- [20] X. Chen, C. Gao, C. Chen, G. Zhang, Y. Liu, An empirical study on challenges for llm application developers, ACM Transactions on Software Engineering and Methodology 1 (1) (2025) 1–35. doi:10.1145/3715007.
- [21] C. B. Seaman, Qualitative methods in empirical studies of software engineering, IEEE Transactions on software engineering 25 (4) (1999) 557–572.
- [22] J. Cohen, A coefficient of agreement for nominal scales, Educational and psychological measurement 20 (1) (1960) 37–46.
- [23] J. R. Landis, G. G. Koch, The measurement of observer agreement for categorical data, biometrics (1977) 159–174.
- [24] N. F. Liu, K. Lin, J. Hewitt, A. Paranjape, M. Bevilacqua, F. Petroni, P. Liang, Lost in the middle: How language models use long contexts, arXiv preprint arXiv:2307.03172 (2023).
- [25] L. Reynolds, K. McDonell, Prompt programming for large language models: Beyond the few-shot paradigm, in: Extended abstracts of the 2021 CHI conference on human factors in computing systems, 2021, pp. 1–7.
- [26] L. Beurer-Kellner, M. Fischer, M. Vechev, Prompting is programming: A query language for large language models, Proceedings of the ACM on Programming Languages 7 (PLDI) (2023) 1946–1969.
- [27] B. Chen, Z. Zhang, N. Langrené, S. Zhu, Unleashing the potential of prompt engineering in large language models: a comprehensive review, arXiv preprint arXiv:2310.14735 (2023).
- [28] Y. Yan, S. Wang, J. Huo, H. Li, B. Li, J. Su, X. Gao, Y.-F. Zhang, T. Xu,

- Z. Chu, et al., Errorradar: Benchmarking complex mathematical reasoning of multimodal large language models via error detection, arXiv preprint arXiv:2410.04509 (2024).
- [29] H. Zhang, K. Zhang, Z. Li, J. Li, Y. Li, Y. Zhao, Y. Zhu, F. Liu, G. Li, et al., Deep learning for code generation: a survey, *Science China Information Sciences* 67 (9) (2024) 191101.
- [30] R. Shakya, F. Vadiee, M. Khalil, A showdown of chatgpt vs deepseek in solving programming tasks, in: 2025 International Conference on New Trends in Computing Sciences (ICTCS), IEEE, 2025, pp. 413–418.
- [31] D. Hendrycks, S. Basart, S. Kadavath, M. Mazeika, A. Arora, E. Guo, C. Burns, S. Puranik, H. He, D. Song, et al., Measuring coding challenge competence with apps, arXiv preprint arXiv:2105.09938 (2021).
- [32] Y. Li, D. Choi, J. Chung, N. Kushman, J. Schriftwieser, R. Leblond, T. Eccles, J. Keeling, F. Gimeno, A. Dal Lago, et al., Competition-level code generation with alphacode, *Science* 378 (6624) (2022) 1092–1097.
- [33] K. Zhang, Z. Li, J. Li, G. Li, Z. Jin, Self-edit: Fault-aware code editor for code generation, in: Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), 2023, pp. 769–787.
- [34] M. A. Islam, M. E. Ali, M. R. Parvez, Mapcoder: Multi-agent code generation for competitive problem solving, arXiv preprint arXiv:2405.11403 (2024).
- [35] D. Souza, R. Gheyi, L. Albuquerque, G. Soares, M. Ribeiro, Code generation with small language models: A deep evaluation on codeforces, arXiv preprint arXiv:2504.07343 (2025).
- [36] L. Yang, R. Jin, L. Shi, J. Peng, Y. Chen, D. Xiong, Probbench: Benchmarking large language models in competitive programming, arXiv preprint arXiv:2502.20868 (2025).

Minnan Wei is currently pursuing a Master's degree at the School of Artificial Intelligence and Computer Science, Nantong University. His research interests include competitive program generation, and vulnerability detection.

Ziming Li is currently pursuing his Bachelor's degree at the School of Artificial Intelligence and Computer Science, Nantong University. His research interests include competitive program generation.

Xiang Chen received the B.Sc. degree in information management and systems from Xi'an Jiaotong University, China, in 2002. Then he received his M.Sc. and Ph.D. degrees in computer software and theory from Nanjing University, China, in 2008 and 2011, respectively. He is currently an Associate Professor at the School of Artificial Intelligence and Computer Science, Nantong University. He has authored or co-authored more than 160 papers in refereed journals or conferences, such as IEEE Transactions on Software Engineering, ACM Transactions on Software Engineering and Methodology, Empirical Software Engineering, Information and Software Technology, Journal of Systems and Software, Software Testing, Verification and Reliability, Journal of Software: Evolution and Process, International Conference on Software Engineering (ICSE), International Conference on the Foundations of Software Engineering (FSE), International Symposium on Software Testing and Analysis (ISSTA), International Conference Automated Software Engineering (ASE), International Conference on Software Maintenance and Evolution (ICSME), International Con-

ference on Program Comprehension (ICPC), and International Conference on Software Analysis, Evolution and Reengineering (SANER). His research interests include software engineering, in particular large language models for software engineering, software testing and maintenance, software repository mining, and empirical software engineering. He received two ACM SIGSOFT distinguished paper awards in ICSE 2021 and ICPC 2023. He is the editorial board member of Information and Software Technology. More information can be found at:

<https://xchen.cs.github.io/index.html>.

Menglin Zheng is currently pursuing a Bachelor's degree in Software Engineering at the School of Artificial Intelligence and Computer Science, Nantong University, with a focus on artificial intelligence and multimodal analysis.

Ziyan Qu is currently pursuing his Bachelor's degree at the School of Artificial Intelligence and Computer Science, Nantong University. His research interests include competitive program generation.

Cheng Yu is currently pursuing his Bachelor's degree at the School of Artificial Intelligence and Computer Science, Nantong University. His research interests include competitive program generation.

Siyu Chen is currently pursuing a Master's degree at the School of Artificial Intelligence and Computer Science, Nan-tong University. Her research interests include software vulnerability analysis.

Xiaolin Ju was born in April 1976. He received the B.S. degree in information science from Wuhan University in 1998, the M.Sc. degree in computer science from Southeast University in 2004, and the Ph.D. degree in computer science from the Chinese University of Mining Technology in 2014. He is currently an associate professor at the School of Information Science and Technology, Nantong University, Nantong, China. His current research interests include software testing, such as collective intelligence, deep learning testing and optimization, and software defects analysis.