

# MACS CODER: A MULTI-AGENT CODING FRAMEWORK FOR SMALL LMS — FROM FAST THINKING TO DEEP PLANNING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Large Language Models (LLMs) have made significant strides in code generation, yet solving complex programming tasks remains a major challenge. Current state-of-the-art (SOTA) multi-agent frameworks, while powerful, often use a resource-intensive, **one-size-fits-all** strategy. We introduce **MACS-Coder** (Multi-Agent Adaptive Coding Structure), a novel dual-process framework designed for high efficiency on personal computers (macOS and Windows PCs). Inspired by human cognition, it comprises two systems: a **Fast Thinking System** for rapid, low-cost code generation, and a **Deep Planning System** for methodical, deliberative problem-solving. This dual architecture allows small models to achieve performance comparable to much larger proprietary models while consuming far less energy and producing lower CO<sub>2</sub> emissions. MACS-Coder dynamically adapts its strategy, employing its Fast Thinking System for simpler tasks and activating its Deep Planning System—composed of planning, structured templating, and fine-grained debugging agents—for complex challenges. Extensive experiments across multiple benchmarks, including the highly challenging LiveCodeBench, demonstrate that MACS-Coder achieves new SOTA pass@1 results. Using the gpt-oss-20B model, it attains accuracies of **99.4%** on HumanEval, **93.2%** on MBPP, and **83.2%** on LiveCodeBench V5, consistently outperforming prior methods such as CodeSIM and MapCoder in both accuracy and computational efficiency. When scaled to a larger open-source backbone (e.g., gpt-oss-120B), MACS-Coder achieves SOTA performance on live-coding benchmarks, surpassing earlier SOTA models. The primary contribution of our work is to bridge the performance gap between compact open-source models and elite closed-source systems: we show that an open-source gpt-oss-20B model empowered by MACS-Coder can achieve performance comparable to top-tier models such as o4-Mini (High) and Gemini 2.5 Pro. By making SOTA-level performance more accessible and resource-efficient, MACS-Coder represents a significant step toward democratizing advanced AI-assisted programming. We will open-source the framework and evaluation code to facilitate future research. Explore our code and video demo at [here](#).

## 1 INTRODUCTION

In recent years, the rapid rise of large language models (LLMs) has driven major advances in AI-assisted programming. These models fundamentally change how developers generate, reason about, and debug code, enabling tasks such as maintaining codebases, adding new features, and even building complete applications without programming expertise, a phenomenon sometimes referred to as “**vibe coding**.”

While flagship proprietary models like OpenAI’s GPT-5 line (OpenAI, 2025) and Anthropic’s CLAUDE-4 (Anthropic, 2025) continue to push the state of the art in coding, a rapidly expanding ecosystem of open-source families (e.g., Meta’s LLAMA-4 (MetaAI, 2025), Alibaba’s QWEN-3 (Yang et al., 2025)) is narrowing performance gaps and broadening accessibility for researchers.

Despite these gains, two practical limitations remain. First, top-performing proprietary models typically depend on very large parameter counts and heavy compute budgets, restricting routine

054 access for individuals, small laboratories, and many production systems. Second, modern multi-agent  
055 and pipeline-based approaches (for instance, staged planning  $\rightarrow$  coding  $\rightarrow$  verification systems  
056 such as MapCoder (Islam et al., 2024a) and CodeSIM (Islam et al., 2025a) improve robustness by  
057 decomposition, but commonly adopt a “one-size-fits-all” strategy that invokes a complete, resource-  
058 intensive stack for every problem instance. This leads to wasted computation on simple tasks and can  
059 still be suboptimal for truly hard problems because the pipeline does not adaptively match reasoning  
060 effort or specialist models to instance difficulty.

061 Small models (the model parameters  $\lesssim 20$  B) remain valuable for research and personalized work-  
062 flows due to their low cost and ease of fine-tuning, but often lag behind larger models in core coding  
063 and mathematical reasoning benchmarks. The central problem we address is whether it is possible to  
064 retain the reliability of cooperative multi-agent pipelines while greatly reducing average compute by  
065 tailoring the amount and kind of reasoning performed per instance.

066 To fill this gap, we propose **MACS-Coder**. MACS-Coder integrates three core ideas: (1) a *Fast-  
067 and-Deep Planning* dual-system architecture that dynamically selects between a low-cost fast path  
068 for easy instances and a deep-planning pipeline for harder ones; (2) *structured generation* through  
069 code templates produced by an STD-IO Tool to improve output stability and make verification easier;  
070 and (3) a *fine-grained debugging* mechanism that isolates and repairs different classes of errors  
071 precisely. The entire architecture is illustrated in Figure 1. In Figure 2, we illustrate how MACS-  
072 Coder operates during the Deep-Planning phase, highlighting its efficiency and comprehensiveness.  
073 The workflow demonstrates a reasoning process that systematically explores solutions, integrates  
074 contextual knowledge, and refines intermediate decisions.

075 We evaluated MACS-Coder on the challenging LiveCodeBench V5 code generation benchmark  
076 (Jain et al., 2024), with its performance on basic tasks such as HumanEval and MBPP detailed  
077 in the Appendix B.1. Our experiments primarily feature the gpt-oss-20B (Agarwal et al., 2025)  
078 and Qwen3-series (Yang et al., 2025) models, while extensive evaluations of a wider variety of  
079 open-source small language models (SLMs) can be found in Appendix B.2. The results show that  
080 MACS-Coder improves computational efficiency while matching or exceeding strong baselines such  
081 as CodeSIM and MapCoder on accuracy. In particular, our experiments demonstrate that targeted  
082 orchestration and structured guidance can substantially amplify the effective performance of midsize  
083 open models, enabling competitive results against lighter proprietary variants while dramatically  
084 reducing average compute.

## 085 CONTRIBUTIONS

- 087 • We analyze cost performance trade-offs across single-agent, fixed multi-agent, and adaptive  
088 orchestration strategies on a spectrum of coding tasks.
- 089 • We introduce MACS-Coder, an adaptive orchestration framework that combines instance-  
090 wise difficulty estimation, progressive delegation, and early exit decisions.
- 091 • We empirically demonstrate that MACS-Coder preserves high success rates on hard tasks  
092 while significantly reducing average computation on easy/medium tasks and that it improves  
093 the utility of small/medium open models when used as targeted specialists.
- 094 • We will release the MACS-Coder implementation and experimental configurations to support  
095 reproducibility and future research.

## 098 2 RELATED WORK

### 100 2.1 LLMs FOR CODE

101 Program synthesis and code generation have long been fundamental challenges in artificial intelli-  
102 gence. The evolution of LLMs has significantly transformed the landscape of automated code-related  
103 tasks, including code completion, code translation, code summarization, and code repair. General-  
104 purpose LLMs are pre-trained on large-scale text, code, and mathematical data to strengthen logical  
105 reasoning and code understanding.

106 Notable proprietary models include Anthropic’s Claude series (Anthropic; 2025), OpenAI’s GPT  
107 family (e.g., GPT-4 (Achiam et al., 2023) and the more recent GPT-5 line (OpenAI, 2025), including

108 the "o" series reasoning models (OpenAI, 2025)), and Google's Gemini series (notably Gemini 2.5  
109 (Comanici et al., 2025)). At the same time, the open-source ecosystem has matured, producing  
110 powerful models that democratize code generation. Representative open-source examples include the  
111 Llama series and Meta's Code Llama variants (Grattafiori et al., 2024; MetaAI, 2025; Rozière et al.,  
112 2023), the Mistral family and Mistral Coder (MistralAI, 2024; Mistral AI team, 2024), DeepSeek and  
113 DeepSeek Coder (Shao et al., 2024; DeepSeek-AI et al., 2024; Guo et al., 2024), and Alibaba's Qwen  
114 series including Qwen-2.5/3 Coder variants (Yang et al., 2024; 2025; Hui et al., 2024; Qwen Team,  
115 2025). These models have demonstrated strong capabilities on many programming problems and  
116 form the foundational core of current code-generation agents.

## 117 118 2.2 MULTI-AGENT CODE GENERATION

119  
120 Although LLM-based code generation techniques can produce standalone programs, their single-  
121 response mode exposes significant limitations for complex engineering-oriented development. Native  
122 LLMs cannot autonomously decompose tasks, interact with real development environments, validate  
123 generated code, or implement continuous self-correction, so they struggle with cross-file context,  
124 dynamic debugging, and iterative optimization. To address these gaps, research has shifted from  
125 prompt engineering toward multi-agent frameworks that decompose problem solving into stages  
126 handled by specialized agents.

127 **MapCoder** (Islam et al., 2024b) proposes a multi-agent framework that mimics the human develop-  
128 ment cycle, incorporating agents for planning, coding, and debugging. Although this approach shows  
129 progress, its workflow is fixed and lacks adaptability to problem difficulty, leading to significant  
130 resource consumption on all problems and limiting its efficiency.

131 **CodeSIM** (Islam et al., 2025b) introduces a novel verification method inspired by human problem-  
132 solving, using input-output simulation to validate generated plans and perform internal debugging.  
133 CodeSIM has achieved SOTA performance in accuracy, but similar to MapCoder, it employs a  
134 resource-intensive, fixed pipeline that fails to dynamically adjust its strategy based on problem  
135 complexity, leaving room for improvement in computational efficiency.

## 136 137 138 2.3 CODE DEBUGGING

139  
140 **LDB** (LLM Debugger) (Zhong et al., 2024) is an external debugging framework that refines programs  
141 using runtime execution information. By segmenting code into basic blocks and tracking intermediate  
142 variable values, LDB enables an LLM to perform block-by-block verification against the task  
143 description, allowing it to accurately pinpoint and correct errors.

144 **LPW** (LLM Programming Workflow) (Lei et al., 2025) features a key innovation called "plan  
145 verification," where a natural language plan is validated against tests to generate a detailed solution  
146 with expected intermediate outputs. This allows for precise debugging by comparing a failed  
147 program's execution trace to the pre-verified outputs, creating a structured signal for refinement.

## 148 149 150 3 MACS-CODER

151  
152 Our objective is to propose an efficient code agent through a two-stage "Fast and Deep Planning"  
153 design. The Fast-Thinking stage addresses simple problems that do not require extensive thought,  
154 while the Deep-Planning stage guides the LLM through step-by-step reasoning based on the processes  
155 of human engineers. By leveraging the concept of test-time compute, it allocates more resources  
156 to enhance the LLM's final performance. Inspired by software engineering practices, we designed  
157 a four-step LLM agent process for the Deep-Planning stage to simulate an engineer's workflow:  
158 analyzing the problem (Planning Agent), generating structured code (Coding Agent), debugging  
159 (Debugging Agent), and handling standard I/O (STD IO Tool). A detailed ablation study, presented in  
160 Appendix C, quantifies the significant impact of each of these components on both Pass@1 accuracy  
161 and token consumption. Drawing inspiration from recent works like MapCoder, CodeSIM, and LDB,  
we developed MACS-Coder to achieve a balance of performance and efficiency.

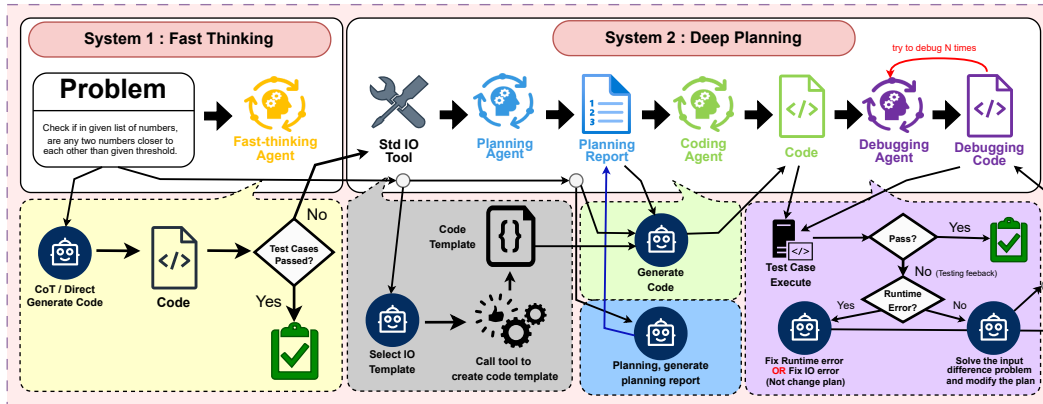


Figure 1: Overview of MACS-Coder

### 3.1 FAST THINKING SYSTEM

The initial stage of our framework is the **Fast-Thinking System**, which serves as a rapid and resource-efficient filter. This stage is designed to handle problems for which the LLMs already exhibit high proficiency. The **Fast-Thinking Agent** generates a code solution with minimal guidance and evaluates its correctness against a set of unit tests.

This process acts as a simple yet effective heuristic for difficulty assessment:

- If the generated code passes all unit tests, the solution is deemed correct and is returned immediately, minimizing computational expenditure.
- If any test fails, the problem is considered non-trivial and is escalated to the Deep-Planning System.

Our implementation adapts its prompting strategy based on the problem domain; for complex competitive programming tasks (LiveCodeBench V5), the agent attempts direct code generation. For benchmarks with clearer solution patterns (HumanEval, MBPP), a lightweight planning-style prompt is used to improve initial accuracy.

### 3.2 DEEP PLANNING SYSTEM

This stage is designed for difficult problems or those where the LLM’s responses are unstable. By increasing resource consumption, it enables the LLM to think more comprehensively and perform debugging to enhance its final performance.

#### 3.2.1 ENVIRONMENT PREPARATION: STD IO TOOL

The STD IO Tool is specifically designed for competitive programming problems, for which we created several input and output templates based on common formats. The LLM uses a tool-calling mechanism to select an appropriate template, ultimately producing a complete **code template**. This allows subsequent agents to modify the template directly when generating code rather than starting from scratch. Furthermore, our designed code templates improve the readability of the final code; details are provided in Appendix D.1. Note that this tool is not activated for simpler benchmarks like HumanEval and MBPP, which lack complex I/O handling. The role of the STD IO Tool within our overall workflow is illustrated in Figure 2 (highlighted as **Step 1**).

#### 3.2.2 STRATEGY FORMULATION: PLANNING AGENT

In the strategy formulation stage, the Planning Agent prompts the LLM to generate a **Planning Report** from the problem description. This report outlines the core algorithm, analyzes potential edge cases, and establishes a coding plan, providing a robust foundation for the implementation phase. To select an optimal algorithm, the LLM first enumerates all viable solutions and then chooses the most

216 stable and straightforward strategy that prioritizes robustness and increases the likelihood of a correct  
 217 solution. This stage corresponds to **Step 2** in our workflow, as depicted in Figure 2. (Full details of  
 218 the instructional prompts are in Appendix D.2.)

### 220 3.2.3 CODE SYNTHESIS: CODING AGENT

221 The Coding Agent (**Step 3**, Figure 2) synthesizes the final code. It takes the problem description, the  
 222 Planning Report, and the Code Template as input, generating an implementation that conforms to the  
 223 template’s structure. To enforce adherence to the plan, the LLM is also prompted to generate code  
 224 comments corresponding to each step in the Planning Report. The resulting code is then validated  
 225 against a suite of unit tests. Code that passes all tests is considered successful; otherwise, it is passed  
 226 to the Debugging Agent for refinement. (The prompts for this stage are detailed in Appendix D.3 and  
 227 Appendix D.4.)

### 229 3.2.4 AUTOMATED DEBUGGING: DEBUGGING AGENT

230 In the Debugging Agent stage (**Step 4**, Figure 2), the agent receives the original problem, the Planning  
 231 Report, the generated code, and the execution log from unit testing as input to debug the code. The  
 232 Debugging Agent consists of three modules: the STD IO Error Block, the Runtime Error Block, and  
 233 the Wrong Answer Block. We first perform a string comparison on the execution log to determine  
 234 if the code executed correctly. If it did, the process moves to the Output Error Block; otherwise, it  
 235 is routed to either the Runtime Error Block or the STD IO Error Block. If this is the first execution  
 236 failure, it enters the STD IO Error Block; subsequent failures are handled by the Runtime Error  
 237 Block.

238 **STD IO Error Block** This block is entered on the first execution failure. Our experiments revealed  
 239 that many execution failures are related to input parsing. If data is not read correctly, an error occurs.  
 240 Therefore, this stage prompts the LLM to regenerate the code for reading and displaying input, while  
 241 the core algorithmic code is preserved.

242 **Runtime Error Block** If the code fails to execute for a second time or more, this block is entered. It  
 243 uses a more general prompt, asking the LLM to analyze the error itself and identify the cause of the  
 244 failure to modify the code, again without altering the core problem-solving logic.

245 **Wrong Answer Block** If the code executes correctly but the output does not match the expected  
 246 answer, this block is entered. Here, the LLM is prompted to analyze the incorrect output, simulate  
 247 the failing test case, and generate an improved algorithm to modify the code.

248 (The specific debugging prompts for each block are presented in Appendix D.5.)

## 251 4 EXPERIMENTAL SETUP

### 253 4.1 EVALUATION BENCHMARKS AND BASELINES

254 We evaluate our method, MACS-Coder, on a diverse set of programming benchmarks. Our primary  
 255 evaluation is on **LiveCodeBench V5** (Jain et al., 2024), a contamination-free benchmark of 880  
 256 competitive problems (279 Easy, 331 Medium, 270 Hard). For a comprehensive assessment of its  
 257 performance on basic tasks, we present detailed results for **HumanEval** (Chen et al., 2021) (164  
 258 problems), **MBPP** (Austin et al., 2021) (974 problems), and their extended-test (-ET) versions (Dong  
 259 et al., 2023) in Appendix B.1.

260 We compare MACS-Coder against several strong baselines: **Direct** generation (Chen et al., 2021),  
 261 **Chain of Thought (CoT)** (Wei et al., 2022), the multi-agent **MapCoder** (Islam et al., 2024b), and a  
 262 state-of-the-art debugging framework, **CodeSIM** (Islam et al., 2025b). We also include **LDB** (Zhong  
 263 et al., 2024), although its use is confined to the fundamental benchmarks (HumanEval/MBPP), as its  
 264 framework is incompatible with the LiveCodeBench test format.

### 267 4.2 IMPLEMENTATION DETAILS

268 Our core evaluation metric is **pass@1**. Our primary experiments utilize three models: **Qwen3 (8B,**  
 269 **14B)** and **gpt-oss (20B)**. To validate generalization, we present further experiments on a broader

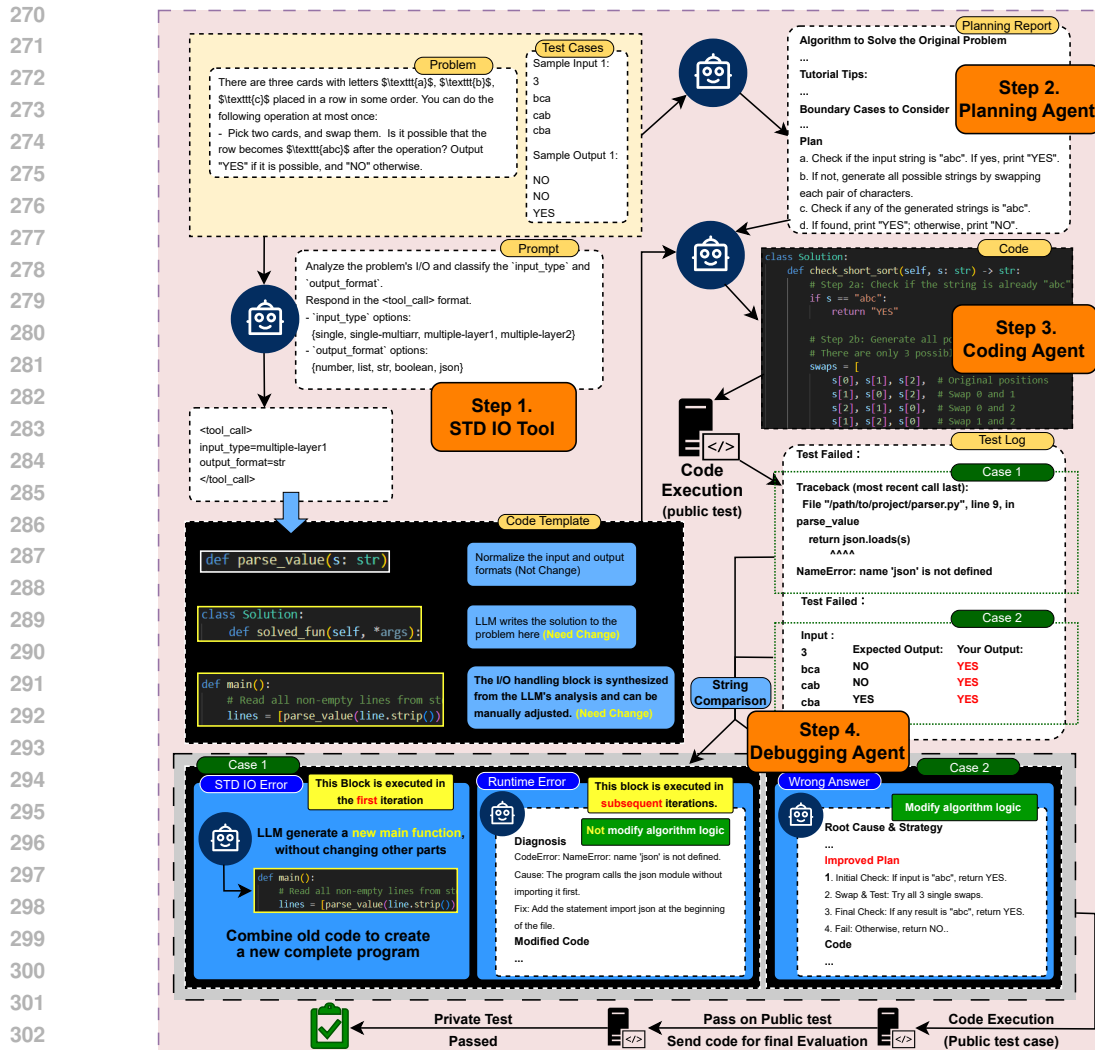


Figure 2: This figure illustrates the MACS-Coder pipeline in the Deep-Planning stage, featuring structured generation for stable, verifiable outputs and fine-grained debugging for precise error isolation, jointly improving efficiency and robustness in code generation.

range of models, including **Llama3.1-8B**, **Gemma2-9B**, **Mistral-8B**, and **Qwen2.5-Coder-7B**, in Appendix B.2. We contextualize our results against top proprietary models from the official **LiveCodeBench Leaderboard**.

Our MACS-Coder framework is governed by two key hyperparameters that control its nested-loop structure. These are  $p$ , the maximum number of outer *planning cycles*, and  $d$ , the maximum number of inner *debug attempts* per planning cycle. Key values for these hyperparameters and model-specific configurations used to suppress unwanted internal reasoning are detailed in Table 1.

Table 1: Hyperparameter and model-specific configurations.

Model Category	Parameters ( $p, d$ )	Configuration Detail
Primary Models	$p = 5, d = 5$	Qwen3: "Non-Thinking Mode" + specific prompt (Ma et al., 2025). gpt-oss: "Reasoning:Low" prompt (OpenAI, 2025).
Generalization Models	$p = 1, d = 5$	N/A

Contest-Level Problems (LiveCodeBench V5)											
LLM	Approach	Easy			Medium			Hard			Total ACC (%)
		ACC (%)	Time (s)	Tokens	ACC (%)	Time (s)	Tokens	ACC (%)	Time (s)	Tokens	
gpt-oss 20B	Direct	87.1	3.66	883	71.3	13.24	3,173	42.2	27.58	6,580	67.3
	CoT	87.8	3.07	721	68.3	10.48	2,465	38.9	23.59	5,459	65.4
	MapCoder	95.7	19.86	4,762	80.7	40.09	9,505	55.9	150.42	35,148	77.8
	CodeSIM	97.1	13.82	3,279	81.9	40.55	9,683	61.5	176.13	41,380	80.4
	<b>MACS-Coder</b>	<b>97.5</b>	<b>5.06</b>	<b>1,204</b>	<b>87.9</b>	<b>23.83</b>	<b>5,690</b>	<b>63.0</b>	<b>143.56</b>	<b>33,968</b>	<b>83.2</b>
Qwen3-14B	Direct	67.4	28.38	1,799	38.4	90.69	5,660	25.2	134.05	8313	43.5
	CoT	58.1	30.64	1,911	36.9	97.24	6,001	23.0	148.54	9,063	39.3
	MapCoder	87.8	82.82	5,144	65.0	<b>260.51</b>	<b>15,831</b>	32.6	<b>426.70</b>	<b>25,833</b>	62.2
	CodeSIM	91.8	142.54	8,622	61.6	467.18	28,107	28.5	726.95	44,037	61.0
	<b>MACS-Coder</b>	<b>95.3</b>	<b>64.69</b>	<b>3,971</b>	<b>72.8</b>	<b>308.64</b>	<b>18,586</b>	<b>35.9</b>	<b>768.47</b>	<b>46,220</b>	<b>68.6</b>
Qwen3-8B	Direct	61.3	32.46	2,072	39.0	51.57	3,231	20.7	60.68	3,748	40.4
	CoT	47.7	4.30	259	13.0	16.52	1,001	4.8	35.72	2,145	21.4
	MapCoder	72.0	<b>39.03</b>	<b>2,475</b>	31.1	<b>82.07</b>	<b>5,183</b>	10.7	<b>134.31</b>	<b>8,356</b>	52.0
	CodeSIM	87.1	89.34	5,688	53.2	262.34	16,592	22.2	678.91	41,591	54.4
	<b>MACS-Coder</b>	<b>94.6</b>	<b>49.71</b>	<b>3,130</b>	<b>66.5</b>	<b>150.07</b>	<b>9,331</b>	<b>26.7</b>	<b>319.72</b>	<b>19,716</b>	<b>63.1</b>

Table 2: Performance comparison of various methods on the LiveCodeBench V5 benchmark, utilizing state-of-the-art (SOTA) open-source SLMs including gpt-oss-20B, Qwen3-14B, and Qwen3-8B. We evaluate each approach based on three key metrics: ACC (Pass@1 Accuracy), Time (average seconds per problem), and Tokens (average generated response tokens per problem). Our proposed method, MACS-Coder, consistently achieves the highest accuracy across all tested models.

## 5 RESULTS

In Table 2, we evaluate the model’s performance on complex, competitive-level programming tasks. MACS-Coder demonstrates a significant superiority over all baseline methods in solving these complex problems. With the most powerful model, gpt-oss-20B, MACS-Coder achieves a total accuracy of 83.2%, setting a new state-of-the-art (SOTA) record and clearly surpassing the next-best methods, CodeSIM (80.4%) and MapCoder (77.8%).

Beyond top-tier accuracy, another core advantage of MACS-Coder lies in its exceptional computational efficiency. The data shows that MACS-Coder achieves higher accuracy while consuming far fewer computational resources (time and tokens) than its competitors. For instance, on the gpt-oss-20B model, MACS-Coder is not only the most accurate but also saves nearly 18% in token consumption on "Hard" level problems compared to CodeSIM (33,968 vs. 41,380). This efficiency advantage is due to our "fast-thinking" system, which quickly handles simple problems, strategically reserving computational resources for challenges that genuinely require deep thought.

The effectiveness of MACS-Coder remains consistent across models of different scales. Particularly noteworthy is its profound empowering effect on smaller models. The Qwen3-8B model equipped with MACS-Coder reaches a total accuracy of 63.1%, not only far exceeding other methods at the same scale but also outperforming the larger Qwen3-14B model when using MapCoder (62.2%). These results strongly indicate that MACS-Coder’s "Fast and Deep Planning" dual-system framework marks a significant advancement in handling the complexity of competitive-level problems.

**Comparison with State-of-the-Art Models** To further validate the absolute competitiveness of the MACS-Coder framework, we compared it with the current strongest open-source and proprietary models on the LiveCodeBench benchmark. The data for this comparison is sourced from the official LiveCodeBench leaderboard (Jain et al.), covering contest problems from **May 1, 2023, to February 1, 2025**. As shown in Table 3, the results indicate that MACS-Coder can elevate the capabilities of open-source small models to a level competitive with top-tier proprietary models.

The most striking result is that the gpt-oss-20B model equipped with MACS-Coder achieves a Pass@1 score of 83.2%. This performance not only surpasses powerful proprietary models like Grok-3-Mini (81.4%) and o3-Mini (80.5%) but is also within a close margin of Gemini-2.5 Pro (84.7%). This demonstrates that MACS-Coder is an extremely powerful performance enhancer, enabling a 20B-level open-source model to compete with the industry’s top models.

Furthermore, the value of MACS-Coder lies in its ability to enable smaller models to achieve SOTA-level performance, making top-tier capabilities more accessible and deployable. The Qwen3 14B with MACS-Coder (68.6%) easily surpasses the un-enhanced gpt-oss-20b (67.3%) and is on par with o1-Mini (68.4%). Even the smaller Qwen3 8B, empowered by MACS-Coder (63.1%), outperforms a range of powerful models like DeepSeek V3 (56.3%) and Claude-3.5-Sonnet (51.5%). These results strongly demonstrate that the MACS-Coder framework provides an efficient pathway for the open-source community to tackle complex programming challenges that were previously only manageable by top-tier proprietary models, using relatively smaller models.

Table 3: Effectiveness of **MACS-Coder** on LiveCodeBench V5 from **May 1, 2023, to February 1, 2025**. Models enhanced with our method show substantial Pass@1 score improvements. The best result in each column is in **bold** and the second best is underlined.

LLM	LiveCodeBench V5			
	Easy	Medium	Hard	Pass@1
<b>gpt-oss-120b (MACS-Coder)</b>	97.9	<u>90.6</u>	<b>70.7</b>	<b>86.8</b>
o4-Mini (High)	<u>98.2</u>	89.7	<u>67.4</u>	<u>85.6</u>
Gemini-2.5 Pro	<u>98.2</u>	<b>91.8</b>	61.9	84.7
<b>gpt-oss-20b (MACS-Coder)</b>	97.5	87.9	63.0	83.2
Grok-3-Mini (High)	98.6	88.8	54.4	81.4
gpt-oss-120b	93.2	85.5	62.2	80.8
o3-Mini (High)	<b>98.9</b>	88.5	51.5	80.5
o1 (Med)	<b>98.9</b>	84.6	49.3	78.3
DeepSeek-R1-Preview	<b>98.9</b>	84.8	47.7	77.9
<b>Qwen3 14B (MACS-Coder)</b>	95.3	72.8	35.9	68.6
O1-Mini	95.0	73.4	34.8	68.4
gpt-oss-20b	87.1	71.3	42.2	67.3
<b>Qwen3 8B (MACS-Coder)</b>	94.6	<u>66.5</u>	<u>26.7</u>	<u>63.1</u>
DeepSeek V3	90.8	59.4	17.0	56.3
o1-Preview	94.3	56.8	14.1	55.6
Claude-3.5-Sonnet	92.9	46.3	14.9	51.5
Qwen3 14B	67.4	38.4	25.2	43.5
GPT-4o	87.4	36.8	6.1	43.4
Gemini-Pro-1.5	87.2	31.1	8.9	42.1
Qwen3 8B	61.3	39.0	20.7	40.4

## 6 ANALYSES

### 6.1 PERFORMANCE-COST ANALYSIS

To visually demonstrate the efficiency advantage of MACS-Coder over the current SOTA method, CodeSIM, we plotted the relationship between accuracy and token consumption (as shown in Figure 3). The graph clearly indicates that MACS-Coder achieves both higher accuracy and lower token consumption across all benchmarks and models, establishing its significant superiority in performance efficiency.

The arrows in the figure consistently point from CodeSIM to MACS-Coder, always moving towards the top-left representing higher accuracy (y-axis) and lower token consumption (x-axis). This improvement is particularly pronounced in complex tasks. For example, on the LiveCodeBench benchmark, MACS-Coder not only increases the accuracy of Qwen3 8B by 8.7% but also saves over 10,000 tokens. Even with a more powerful model like gpt-oss 20B, MACS-Coder still improves accuracy by 2.8% while saving over 4,400 tokens. This comprehensive performance advantage validates the design of our "Fast and Deep Planning" framework. It intelligently allocates computational resources, avoiding unnecessary deep thought on simple problems to concentrate resources on tackling difficult ones. This makes MACS-Coder not only a more accurate solution but also a more economical and efficient framework for practical applications.

### 6.2 ANALYSIS OF PLANNING ITERATIONS

To analyze the relationship between planning-stage investment and final performance, we evaluated the Pass@1 accuracy of MACS-Coder against CodeSIM over varying planning iterations (Figure 4). The results reveal a significant efficiency advantage for MACS-Coder, which consistently outperforms CodeSIM across all tested models and iteration counts. This efficiency gain is so substantial that MACS-Coder with the Qwen3 8B model at a single iteration achieves 56.8% Pass@1, surpassing the larger Qwen3 14B model using CodeSIM (49.4%). This finding demonstrates that our framework offers significant "capability compression," enabling smaller models to achieve results previously

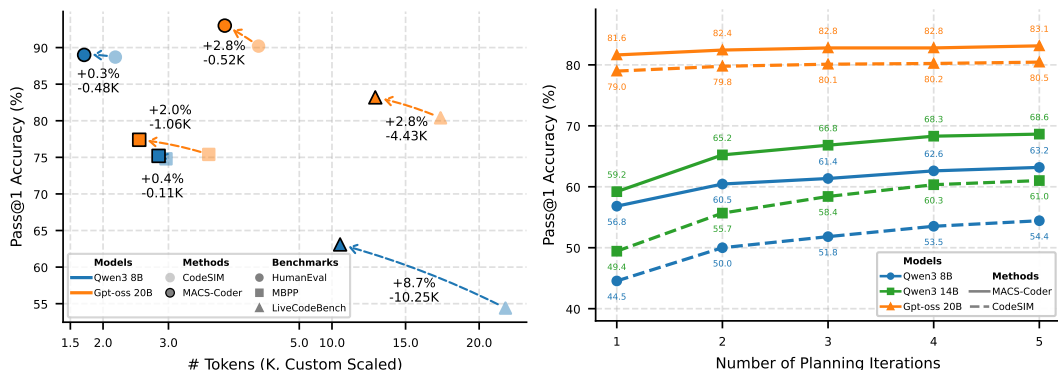


Figure 3: Accuracy-efficiency improvement achieved by our proposed methods. The plot shows that our approaches consistently push the models towards the ideal top-left corner (higher accuracy, fewer tokens).

requiring larger ones. This holds substantial implications for reducing computational costs while achieving state-of-the-art performance in practical applications.

## 7 CONCLUSION AND FUTURE WORK

In this paper, we introduced MACS-Coder, a novel framework designed to solve complex programming tasks. At its core, MACS-Coder features an innovative "Fast and Deep Planning" dual-system architecture that adaptively adjusts its strategy based on problem difficulty and integrates structured code template generation with a fine-grained debugging agent. Evaluation results across multiple mainstream code generation benchmarks show that MACS-Coder significantly surpasses existing SOTA methods in both accuracy and computational efficiency. Our research demonstrates that by empowering small open-source models, MACS-Coder successfully elevates their performance to a level competitive with top-tier proprietary models.

Future work will focus on extending this framework to other domains requiring complex reasoning, such as mathematical problem-solving and general question-answering, to broaden its scope and impact.

## 8 LIMITATIONS

Despite the strong performance and efficiency of MACS-Coder, some limitations exist. First, while our framework demonstrates strong generalization across multiple open-source SLMs, the magnitude of performance improvement is not entirely consistent. This variability is partly tied to the intrinsic capabilities of the base models themselves; an agentic framework can structure and guide an SLM’s reasoning, but its effectiveness is ultimately constrained by the model’s fundamental coding ability. If a model’s baseline performance is below a certain threshold, the framework’s capacity to elicit further gains is limited. Additionally, adapting our structured prompts for specific model architectures remains an area for optimization. Second, while MACS-Coder is more token-efficient than prior SOTA methods, its computational cost remains higher than Direct Prompting, marking an avenue for future work. Finally, the efficacy of our debugging agent is highly dependent on the quality and coverage of the provided test cases.

## REFERENCES

OpenAI Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, et al. Gpt-4 technical report. 2023. URL <https://api.semanticscholar.org/CorpusID:257532815>.

- 486 OpenAI Sandhini Agarwal, Lama Ahmad, Jason Ai, Sam Altman, and Others. gpt-oss-120b&gpt-  
487 oss-20b model card. 2025. URL <https://api.semanticscholar.org/CorpusID:280671456>.  
488
- 489 Anthropic. The claude 3 model family: Opus, sonnet, haiku. URL <https://api.semanticscholar.org/CorpusID:268232499>.  
490  
491
- 492 Anthropic. Introducing claude 4. <https://www.anthropic.com/news/claude-4>, May 2025. Ac-  
493 cessed 2025-09-01.  
494
- 495 Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan,  
496 Ellen Jiang, Carrie J. Cai, Michael Terry, Quoc V. Le, and Charles Sutton. Program synthesis with  
497 large language models. *ArXiv*, abs/2108.07732, 2021. URL <https://api.semanticscholar.org/CorpusID:237142385>.  
498
- 499 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Pondé, Jared Kaplan, Harrison  
500 Edwards, et al. Evaluating large language models trained on code. *ArXiv*, abs/2107.03374, 2021.  
501 URL <https://api.semanticscholar.org/CorpusID:235755472>.  
502
- 503 Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit  
504 Dhillon, et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long  
505 context, and next generation agentic capabilities. *ArXiv*, abs/2507.06261, 2025. URL <https://api.semanticscholar.org/CorpusID:280151524>.  
506
- 507 DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bing-Li Wang, Bochao Wu, Chengda Lu, Chenggang  
508 Zhao, et al. Deepseek-v3 technical report. *ArXiv*, abs/2412.19437, 2024. URL <https://api.semanticscholar.org/CorpusID:275118643>.  
509  
510
- 511 Yihong Dong, Ji Ding, Xue Jiang, Zhuo Li, Ge Li, and Zhi Jin. Codescore: Evaluating code generation  
512 by learning code execution. *ACM Transactions on Software Engineering and Methodology*, 34:1 –  
513 22, 2023. URL <https://api.semanticscholar.org/CorpusID:256105296>.  
514
- 515 Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad  
516 Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of  
517 models. *arXiv preprint arXiv:2407.21783*, 2024.  
518
- 519 Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao  
520 Bi, Yu Wu, Y. K. Li, Fuli Luo, Yingfei Xiong, and Wenfeng Liang. Deepseek-coder: When the  
521 large language model meets programming - the rise of code intelligence. *ArXiv*, abs/2401.14196,  
522 2024. URL <https://api.semanticscholar.org/CorpusID:267211867>.  
523
- 524 Binyuan Hui, Jian Yang, Zeyu Cui, Jiayi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun  
525 Zhang, Bowen Yu, Kai Dang, An Yang, Rui Men, Fei Huang, Shanghaoran Quan, Xingzhang  
526 Ren, Xuancheng Ren, Jingren Zhou, and Junyang Lin. Qwen2.5-coder technical report. *ArXiv*,  
527 abs/2409.12186, 2024. URL <https://api.semanticscholar.org/CorpusID:272707390>.  
528
- 529 Md. Ashraful Islam, Mohammed Eunos Ali, and Md. Rizwan Parvez. Mapcoder: Multi-agent code  
530 generation for competitive problem solving. In *Annual Meeting of the Association for Computa-  
531 tional Linguistics*, 2024a. URL <https://api.semanticscholar.org/CorpusID:269921148>.  
532
- 533 Md. Ashraful Islam, Mohammed Eunos Ali, and Md Rizwan Parvez. MapCoder: Multi-agent code  
534 generation for competitive problem solving. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar  
535 (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics  
536 (Volume 1: Long Papers)*, pp. 4912–4944, Bangkok, Thailand, August 2024b. Association for  
537 Computational Linguistics. doi: 10.18653/v1/2024.acl-long.269. URL <https://aclanthology.org/2024.acl-long.269/>.  
538
- 539 Md. Ashraful Islam, Mohammed Eunos Ali, and Md. Rizwan Parvez. Codesim: Multi-agent  
code generation and problem solving through simulation-driven planning and debugging. *ArXiv*,  
abs/2502.05664, 2025a. URL <https://api.semanticscholar.org/CorpusID:276249419>.

- 540 Md. Ashraful Islam, Mohammed Eunus Ali, and Md Rizwan Parvez. CodeSim: Multi-agent  
541 code generation and problem solving through simulation-driven planning and debugging. In  
542 Luis Chiruzzo, Alan Ritter, and Lu Wang (eds.), *Findings of the Association for Computational*  
543 *Linguistics: NAACL 2025*, pp. 5113–5139, Albuquerque, New Mexico, April 2025b. Association  
544 for Computational Linguistics. ISBN 979-8-89176-195-7. doi: 10.18653/v1/2025.findings-naacl.  
545 285. URL <https://aclanthology.org/2025.findings-naacl.285/>.
- 546 Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando  
547 Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free eval-  
548 uation of large language models for code. [https://livecodebench.github.io/leaderboard\\_](https://livecodebench.github.io/leaderboard_v5.html)  
549 [v5.html](https://livecodebench.github.io/leaderboard_v5.html). Accessed: 2025-09-01.
- 550 Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando  
551 Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free  
552 evaluation of large language models for code. *arXiv preprint arXiv:2403.07974*, 2024.
- 553 Chao Lei, Yanchuan Chang, Nir Lipovetzky, and Krista A. Ehinger. Planning-driven programming:  
554 A large language model programming workflow. In Wanxiang Che, Joyce Nabende, Ekaterina  
555 Shutova, and Mohammad Taher Pilehvar (eds.), *Proceedings of the 63rd Annual Meeting of the*  
556 *Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 12647–12684, Vienna,  
557 Austria, July 2025. Association for Computational Linguistics. ISBN 979-8-89176-251-0. doi:  
558 10.18653/v1/2025.acl-long.621. URL <https://aclanthology.org/2025.acl-long.621/>.
- 559 Wenjie Ma, Jingxuan He, Charlie Snell, Tyler Griggs, Sewon Min, and Matei Zaharia. Reasoning  
560 models can be effective without thinking. *ArXiv*, abs/2504.09858, 2025. URL [https://api.](https://api.semanticscholar.org/CorpusID:277781570)  
561 [semanticscholar.org/CorpusID:277781570](https://api.semanticscholar.org/CorpusID:277781570).
- 562 MetaAI. The llama 4 herd: The beginning of a new era of natively multimodal ai innovation.  
563 <https://ai.meta.com/blog/llama-4-multimodal-intelligence/>, April 2025. Accessed  
564 2025-09-01.
- 565 Mistral AI team. Codestral. <https://mistral.ai/news/codestral>, 2024. [Online; accessed  
566 2025-09-07].
- 567 MistralAI. Mistral large. <https://mistral.ai/news/mistral-large>, February 2024. Accessed  
568 2025-09-01.
- 569 OpenAI. Introducing openai o3 and o4-mini. [https://openai.com/index/](https://openai.com/index/introducing-o3-and-o4-mini/)  
570 [introducing-o3-and-o4-mini/](https://openai.com/index/introducing-o3-and-o4-mini/), 2025. [Online; accessed Sep 7, 2025].
- 571 OpenAI. Introducing gpt-5. <https://openai.com/gpt-5/>, August 2025. Accessed 2025-09-01.
- 572 OpenAI. GPT-OSS 20B Huggingface. <https://huggingface.co/openai/gpt-oss-20b>, 2025.  
573 Accessed: 2025-09-01.
- 574 Qwen Team. Qwen3-coder: Agentic coding in the world. [https://qwenlm.github.io/blog/](https://qwenlm.github.io/blog/qwen3-coder/)  
575 [qwen3-coder/](https://qwenlm.github.io/blog/qwen3-coder/), 2025. [Online; accessed 2025-09-07].
- 576 Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Tan, Yossi Adi,  
577 Jingyu Liu, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, I. Evtimov, Joanna Bitton, Manish P  
578 Bhatt, Cris tian Cantón Ferrer, Aaron Grattafiori, Wenhan Xiong, Alexandre D’efossez, Jade  
579 Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier, Thomas Scialom, and Gabriel  
580 Synnaeve. Code llama: Open foundation models for code. *ArXiv*, abs/2308.12950, 2023. URL  
581 <https://api.semanticscholar.org/CorpusID:261100919>.
- 582 Zhihong Shao, Damai Dai, Daya Guo, Bo Liu (Benjamin Liu), Zihan Wang, and Huajian Xin.  
583 Deepseek-v2: A strong, economical, and efficient mixture-of-experts language model. *ArXiv*,  
584 abs/2405.04434, 2024. URL <https://api.semanticscholar.org/CorpusID:269613809>.
- 585 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed H. Chi,  
586 Quoc V Le, and Denny Zhou. Chain of thought prompting elicits reasoning in large language  
587 models. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances*  
588 *in Neural Information Processing Systems*, 2022. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=_VjQ1MeSB_J)  
589 [\\_VjQ1MeSB\\_J](https://openreview.net/forum?id=_VjQ1MeSB_J).

594 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang  
595 Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng  
596 Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang,  
597 Jianxin Yang, Jiaxin Yang, Jingren Zhou, Jingren Zhou, Junyan Lin, Kai Dang, Keqin Bao, Ke-  
598 Pei Yang, Le Yu, Li-Chun Deng, Mei Li, Min Xue, Mingze Li, Pei Zhang, Peng Wang, Qin  
599 Zhu, Rui Men, Ruize Gao, Shi-Qiang Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao  
600 Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yi-  
601 Chao Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang,  
602 Zhipeng Zhou, and Zihan Qiu. Qwen3 technical report. *ArXiv*, abs/2505.09388, 2025. URL  
603 <https://api.semanticscholar.org/CorpusID:278602855>.

604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617

618 Qwen An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan  
619 Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu,  
620 Jianwei Zhang, Jianxin Yang, Jiaxin Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu,  
621 Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji  
622 Lin, Tianhao Li, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yi-Chao  
623 Zhang, Yunyang Wan, Yuqi Liu, Zeyu Cui, Zhenru Zhang, Zihan Qiu, Shanghaoran Quan, and  
624 Zekun Wang. Qwen2.5 technical report. *ArXiv*, abs/2412.15115, 2024. URL <https://api.semanticscholar.org/CorpusID:274859421>.

625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639

640 Li Zhong, Zilong Wang, and Jingbo Shang. Debug like a human: A large language model debugger via  
641 verifying runtime execution step by step. In *Annual Meeting of the Association for Computational*  
642 *Linguistics*, 2024. URL <https://api.semanticscholar.org/CorpusID:268032812>.

643  
644  
645  
646  
647

## A PSEUDOCODE OF THE MACS-CODER FRAMEWORK

---

### Algorithm 1 MACS-Coder Framework

---

```

648
649
650
651
652 1: Input: problem  $P$ , unit tests  $U_{tests}$ 
653 2: Output: solution code  $C_{sol}$ 
654 3:  $p \leftarrow$  max number of planning cycles
655 4:  $d \leftarrow$  max number of debugging steps per cycle
656
657                                     ▷ Stage 1: Fast Thinking System
658 5:  $C_{fast} \leftarrow$  FastThinkGenerate( $P$ )
659 6:  $passed, log \leftarrow$  Test( $C_{fast}, U_{tests}$ )
660 7: if  $passed$  then
661 8:   return  $C_{fast}$ 
662 9: end if
663
664                                     ▷ Stage 2: Deep Planning System
665 10: for  $i \leftarrow 1$  to  $p$  do
666                                     ▷ Start of Planning Cycle
667 11:    $T_{code} \leftarrow$  GenerateCodeTemplate( $P$ )
668                                     ▷ STD IO Tool
669 12:    $R_{plan} \leftarrow$  GeneratePlan( $P$ )
670                                     ▷ Planning Agent
671 13:    $C_{current} \leftarrow$  GenerateCode( $P, R_{plan}, T_{code}$ )
672                                     ▷ Coding Agent
673 14:   first_failure_in_cycle  $\leftarrow$  true
674
675                                     ▷ Start of Debugging Sub-Cycle
676 15:   for  $j \leftarrow 1$  to  $d$  do
677 16:      $passed, log \leftarrow$  Test( $C_{current}, U_{tests}$ )
678 17:     if  $passed$  then
679 18:       return  $C_{current}$ 
680 19:     end if
681 20:      $error\_type \leftarrow$  AnalyzeLog( $log$ )
682                                     ▷ Determine error type
683 21:     if  $error\_type =$  RUNTIME_ERROR then
684 22:       if first_failure_in_cycle then
685 23:          $C_{current} \leftarrow$  FixSTDIOError( $P, C_{current}$ )
686 24:         first_failure_in_cycle  $\leftarrow$  false
687 25:       else
688 26:          $C_{current} \leftarrow$  FixRuntimeError( $P, C_{current}, log$ )
689 27:       end if
690 28:     else if  $error\_type =$  WRONG_ANSWER then
691 29:        $C_{current} \leftarrow$  FixLogicalError( $P, R_{plan}, C_{current}, log$ )
692 30:     end if
693 31:   end for
694                                     ▷ End of Debugging Sub-Cycle
695 32: end for
696                                     ▷ End of Planning Cycle
697 33: return  $C_{current}$ 
698                                     ▷ Return the last generated code
699
700
701

```

---

## B DETAILED PERFORMANCE ON FOUNDATIONAL BENCHMARKS AND GENERALIZATION MODELS

### B.1 BASIC CODE GENERATION

In the basic code generation tasks shown in Table 4, the core advantage of MACS-Coder is its outstanding overall performance efficiency, achieving an optimal balance between top-tier accuracy and computational cost.

The efficient design of MACS-Coder allows it to reach SOTA-level accuracy with fewer resources. For example, on HumanEval, the Qwen3-8B model with MACS-Coder not only achieves a top average accuracy of 89.0% but also has significantly lower time (26.80s) and token (1,711) consumption compared to the similarly performing CodeSIM.

More importantly, this high efficiency enables smaller models to achieve performance beyond their scale. The Qwen3-8B with MACS-Coder on HumanEval (89.0%) performs on par with the larger Qwen3-14B using the CoT method (89.0%). This proves that our framework can achieve results with

Basic Programming Problems											
LLM	Approach	HumanEval					MBPP				
		ACC (%)	Time (s)	Tokens	ET ACC (%)	Avg ACC (%)	ACC (%)	Time (s)	Tokens	ET ACC (%)	Avg ACC (%)
gpt-oss 20B	Direct	62.2	4.07	981	55.5	58.8	47.6	3.51	848	31.2	39.4
	CoT	98.2	3.49	829	85.4	91.8	77.6	2.15	522	52.6	65.1
	LDB	98.8	<b>3.91</b>	<b>937</b>	86.0	92.4	90.7	<b>4.53</b>	<b>1096</b>	58.7	74.7
	MapCoder	97.0	22.31	5314	84.8	90.9	92.9	21.04	5047	61.2	77.0
	CodeSIM	97.0	18.43	4377	83.5	90.2	91.4	15.03	3616	59.4	75.4
	MACS-Coder	<b>99.4</b>	<b>16.23</b>	<b>3861</b>	<b>86.6</b>	<b>93.0</b>	<b>93.2</b>	<b>10.69</b>	<b>2553</b>	<b>61.7</b>	<b>77.4</b>
Qwen3-14B	Direct	48.2	14.16	910	43.3	45.7	40.1	14.16	911	30.0	35.0
	CoT	93.9	15.86	1000	<b>84.1</b>	89.0	76.6	14.25	912	51.6	64.1
	LDB	95.1	<b>23.21</b>	<b>1465</b>	<b>84.1</b>	89.6	90.9	<b>36.09</b>	<b>2263</b>	<b>62.5</b>	<b>76.7</b>
	MapCoder	94.5	46.16	2896	81.7	88.1	88.9	79.38	4851	59.2	74.0
	CodeSIM	95.1	89.62	5460	<b>84.1</b>	89.6	<b>92.4</b>	88.72	5413	59.7	76.0
	MACS-Coder	<b>95.7</b>	<b>44.87</b>	<b>2812</b>	<b>84.1</b>	<b>89.9</b>	<b>90.7</b>	<b>67.33</b>	<b>4181</b>	<b>59.9</b>	<b>75.3</b>
Qwen3-8B	Direct	48.8	25.23	1621	41.5	45.1	20.4	23.21	1500	14.6	17.5
	CoT	79.3	1.54	96	70.1	74.7	74.1	1.42	92	51.9	63.0
	LDB	88.4	<b>11.80</b>	<b>758</b>	78.0	83.2	84.1	<b>12.07</b>	<b>775</b>	55.2	69.6
	MapCoder	91.5	26.29	1685	78.0	84.7	86.1	29.57	1746	56.2	71.1
	CodeSIM	<b>94.5</b>	35.09	2188	82.9	88.7	89.2	47.68	2962	<b>60.5</b>	74.8
	MACS-Coder	<b>93.9</b>	<b>26.80</b>	<b>1711</b>	<b>84.1</b>	<b>89.0</b>	<b>89.9</b>	<b>48.24</b>	<b>2849</b>	<b>60.5</b>	<b>75.2</b>

Table 4: Performance comparison on the HumanEval and MBPP benchmarks using open-source models (gpt-oss-20B, Qwen3-14B, and Qwen3-8B). The evaluation metrics include: ACC (Pass@1 Accuracy), Time (average seconds per problem), Tokens (average generated response tokens per problem), and ET ACC (a more stringent evaluation with additional private test cases). The results show that MACS-Coder not only achieves the highest accuracy but also demonstrates superior efficiency by consuming fewer Tokens than other high-performing methods like CodeSIM and MapCoder.

more economical smaller models that previously required more expensive larger models, offering a solution for the code generation field that combines top-tier capability with practical utility.

## B.2 PERFORMANCE ACROSS OPEN-SOURCE SLMs

To further demonstrate the generalization capability of the MACS-Coder framework, we evaluated its performance across several mainstream open-source SLMs, including Llama3.1-8B, Gemma2-9B, Ministral-8B, and Qwen2.5-Coder-7B. As shown in Table 5, the results highlight MACS-Coder’s exceptional adaptability and its superior balance between accuracy and efficiency.

On the more complex **LiveCodeBench V5** dataset, MACS-Coder consistently and significantly outperforms CodeSIM in accuracy across all tested models. For instance, with **Ministral-8B**, MACS-Coder achieves an accuracy of 20.0%, a substantial improvement over CodeSIM’s 14.0%. Similar advantages are observed on **Llama3.1-8B** (15.1% vs. 13.9%) and **Gemma2-9B** (12.5% vs. 10.3%), proving its robust problem-solving capabilities.

On the **HumanEval** dataset, MACS-Coder demonstrates a more nuanced but equally compelling advantage. While its accuracy is highly competitive—and even superior on models like **Ministral-8B** (85.4% vs. 79.3%)—its primary strength lies in its remarkable efficiency. For example, on **Qwen2.5-Coder-7B**, although CodeSIM’s accuracy is slightly higher (87.8% vs. 86.6%), MACS-Coder completes the task **29% faster** while consuming **25% fewer tokens**. This trend of achieving comparable or higher accuracy with significantly less computational cost is a consistent theme across the models.

These results strongly indicate that MACS-Coder’s dual-system framework is not just optimized for a specific model but can be widely adapted to enhance different language model architectures.

Table 5: Performance comparison on HumanEval and LiveCodeBench V5 across different open-source LLMs.

LLM	Approach	HumanEval				LiveCodeBench V5		
		ACC (%)	Time (s)	Tokens	ET ACC (%)	ACC (%)	Time (s)	Tokens
Llama3.1 8B Instruct	Direct	50.0	5.58	395	43.3	10.7	6.51	451
	CodeSIM	<b>82.3</b>	69.12	4650	<b>71.3</b>	13.9	94.94	6304
	MACS-Coder	80.5	<b>56.92</b>	<b>3881</b>	67.7	<b>15.1</b>	<b>63.69</b>	<b>4271</b>
Gemma2 9B Instruct	Direct	66.5	8.98	477	58.5	12.0	10.54	537
	CodeSIM	80.5	87.35	<b>3981</b>	67.1	10.3	<b>84.09</b>	<b>3870</b>
	MACS-Coder	<b>81.1</b>	<b>83.00</b>	4029	<b>67.7</b>	<b>12.5</b>	110.01	5165
Ministral 8B Instruct	Direct	74.4	6.41	447	64.6	13.9	8.64	579
	CodeSIM	79.3	33.29	2161	67.7	14.0	85.40	5214
	MACS-Coder	<b>85.4</b>	<b>27.00</b>	<b>1760</b>	<b>71.3</b>	<b>20.0</b>	<b>77.78</b>	<b>4803</b>
Qwen2.5-Coder 7B Instruct	Direct	70.1	3.31	233	62.8	18.5	6.08	426
	CodeSIM	<b>87.8</b>	25.07	1661	<b>76.8</b>	23.0	74.85	5207
	MACS-Coder	86.6	<b>17.70</b>	<b>1246</b>	75.0	<b>23.3</b>	<b>64.09</b>	<b>4427</b>

Its ability to deliver high accuracy while drastically improving computational efficiency validates MACS-Coder as a universal and highly practical enhancement framework for code generation.

## C ABLATION STUDY OF MACS-CODER COMPONENTS

### C.1 EFFECTIVENESS OF THE DUAL-SYSTEM FRAMEWORK

To validate the value of our "Fast and Deep Planning" dual-system architecture, we conducted a critical ablation study with results shown in Table 6. The findings demonstrate a decisive victory for the complete framework: enabling the fast-thinking system not only saves **16.5%** in token consumption but, critically, boosts Pass@1 accuracy from 52.4% to 63.1%. This seemingly counter-intuitive accuracy gain powerfully validates our dual-system philosophy.

We attribute this phenomenon to the avoidance of "overthinking," where the sophisticated "Deep-Planning" agentic workflow, designed for complex tasks, can introduce unnecessary failure points or negative guidance when applied to simpler problems. The "fast-thinking" system acts as an intelligent triage mechanism that circumvents this risk by leveraging the model’s raw capabilities on straightforward tasks. This result strongly proves that our dual-system design is not merely a layering of processes but an intelligent resource allocation strategy that raises the problem-solving ceiling by applying the right cognitive tool to the right problem, optimizing for both accuracy and computational cost.

Table 6: Token consumption analysis with and without Fast Thinking. (on LiveCodeBench V5)

Fast Thinking	Pass@1	Tokens	Tokens Saved (%)
✗	52.4%	12,650	-
✓	63.1%	10,551	16.5%

### C.2 CONTRIBUTION OF EACH AGENT IN DEEP-PLANNING SYSTEM

To quantify the contribution of the core agents in the "Deep-Planning" system, we conducted further ablation studies, with results shown in Table 7. The data clearly reveals the critical impact of the STD IO Tool and the Debugging Agent on final performance. When we remove the Debugging Agent, the Pass@1 accuracy drops from 63.1% to 55.3%, a performance loss of 7.8%. Removing the STD

IO Tool leads to an even greater performance decline, with accuracy falling to 53.6%, a drop of 9.5%. If both are disabled, the accuracy plummets to 47.8%. This means that each agent contributes approximately 8-9% to the accuracy, which, in a large test set, is equivalent to successfully solving 70 to 80 additional problems. This result strongly demonstrates that our structured template generation and fine-grained debugging processes, designed to mimic human engineers, are the two pillars that enable MACS-Coder to stably solve complex problems.

Table 7: Ablation study on STD IO Tool and Debugging Agent. (on LiveCodeBench V5)

STD IO Tool	Debugging Agent	Pass@1	Performance Drop
X	X	47.8%	15.3%
✓	X	55.3%	7.8%
X	✓	53.6%	9.5%
✓	✓	63.1%	-

### C.3 ANALYSIS OF FINE-GRAINED DEBUGGING MODULES

To validate the value of each specialized module within our fine-grained Debugging Agent, we conducted an ablation study, with results shown in Table 8. The data indicates that every error-handling module makes an indispensable contribution to the final performance. The module for handling **Runtime Errors** has the most significant impact; its removal leads to a 5.0% performance drop, highlighting the importance of managing runtime exceptions in complex programming. In comparison, the modules for handling **STD IO Errors** and **Wrong Answer** have a smaller but still crucial impact (0.5% and 0.9%, respectively), and their presence collectively ensures the integrity and robustness of the debugging process. This result proves that our fine-grained design, which allows the agent to apply the most appropriate correction strategy for different failure reasons, is key to its efficient debugging capabilities.

Table 8: Ablation study of Debugging Agent components. (on LiveCodeBench V5)

Runtime Error	STD IO Error	Wrong Answer	Pass@1	Performance Drop (%)
✓	✓	✓	63.1%	-
X	✓	✓	58.1%	5.0%
✓	X	✓	62.6%	0.5%
✓	✓	X	62.2%	0.9%

## D THE MULTI-STAGE CHAT-TEMPLATE DESIGN OF MACS-CODER

### D.1 THE PROMPT OF THE STD IO TOOL IN THE DEEP PLANNING SYSTEM

```
The prompt of the STD IO Tool in the Deep Planning System (In Step 1)

System Prompt: You are an automated code analysis tool. Based on the following problem’s
input and output examples, determine which "input type" and "output format" should be used.
Use the <tool_call> tag and output only this:

<tool_call>
input_type = {single | single-multiarr | multiple-layer1 | multiple-layer2}
output_format = {number | list | str | boolean | json}
</tool_call>
```

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

**Classification rules:**

1. **single** – Only one test case provided per execution (with variants):

- *single-simple*: a single plain input (no extra parameters)
- *single-multiarr*: matrix or multiple-row array, for example:

```
4 2
1 2 3 4
4 3 1 2
```

2. **multiple** – Multiple test cases in one run (with T first, followed by per-case data):

- *multiple-layer1*: T lines, each a single-line test case, e.g.:

```
6
abc
acb
...
```

- *multiple-layer2*: multiple test cases, each spanning multiple lines, e.g.:

```
4
4
2 2 1 2
...
```

**Output format** must be one of:

- *number*: an integer or floating-point number
- *list*: space-separated items
- *str*: a single string
- *boolean*: true or false
- *json*: valid JSON structure (e.g. {"a":1,"b":2})

**Example 1**

Input:

```
{"input":"4 2\n1 2 3 4\n4 3 1 2\n", "output":"1 4\n"}
```

Expected response:

```
<tool_call>
input_type = single-multiarr
output_format = list
</tool_call>
```

**Example 2**

Input:

```
{"input":"4\n4\n2 2 1 2\n3\n0 1 2\n5\n4 3 2 3 4\n9\n9 9 9 9 9 9 9 9 9\n", "output":"16\n2\n432\n430467210\n"}
```

Expected response:

```
<tool_call>
input_type = multiple-layer2
output_format = number
</tool_call>
```

918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971

**Example 3**

Input:

```
{"input": "[1, 4, 3, 8, 5]", "output": "[1, 3]" }
```

Expected response:

```
<tool_call>
input_type = single
output_format = list
</tool_call>
```

**Example 4**

Input:

```
{"input": "\"51230100\"", "output": "\"512301\""} 
```

Expected response:

```
<tool_call>
input_type = single
output_format = str
</tool_call>
```

**Important:** Analyze the problem’s intended function signature, not just the test cases. Determine the correct `input_type` and `output_format` for the function you’ll implement. Below is the problem description:  
{Problem}

## D.2 THE PROMPT OF THE PLANNING AGENT IN THE DEEP PLANNING SYSTEM

## The prompt of the Planning Agent in the Deep Planning System. (In Step 2)

You are a programmer tasked with generating an appropriate plan to solve a given problem using the `{language}` programming language.

PROBLEM

{problem}

**Expected Output:**

Your response must be structured as follows:

PROBLEM UNDERSTANDING

- Think about the original problem. Develop an initial understanding of the problem.

PROBLEM SETTER PERSPECTIVE

Analyze the problem as if you are the problem setter:

- What core concept or algorithm is being tested?
- What kind of solution do you expect the solver to come up with?
- Why might this problem have been designed in this particular way?

POSSIBLE SOLUTION METHODS

- List **all possible algorithms or techniques** that could be applied to solve the problem.
- For each method, briefly describe:
  - Its general idea
  - Its complexity
  - Its pros and cons in the context of this problem

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

#### CHOSEN METHOD: SIMPLICITY AND RELIABILITY

- Based on the methods listed above, choose the one that is simplest, most robust, and most likely to succeed given the problem constraints.
- Justify why this method is chosen over the others.

#### RECALL EXAMPLE PROBLEM

Recall a relevant and distinct problem (different from the one mentioned above):

- Describe it
- Write the {language} code step-by-step to solve that problem
- Discuss the algorithm used in that example
- Generate a plan to solve that example problem

#### ALGORITHM TO SOLVE THE ORIGINAL PROBLEM

- Describe the chosen algorithm again, but now in the specific context of the original problem.
- Include tutorial-style tips:
  - How to approach this type of algorithm
  - Important pitfalls to avoid
  - Any tricks to make implementation easier

#### BOUNDARY CASES TO CONSIDER

- Based on the problem description and constraints, list all boundary cases that need special handling.
- Use clear natural language to describe these boundary cases.

#### PLAN

- Write down a detailed, step-by-step plan to solve the **original problem**.

---

#### Important Instruction:

- Strictly follow the instructions.
- Do not output any code.
- The plan must be written in natural language, not code.

### D.3 THE PROMPT OF THE CODE TEMPLATE IN THE DEEP PLANNING SYSTEM

#### The prompt of the Code Template in the Deep Planning System. (In Step 3)

##### CODE FORMAT TEMPLATE

```
import sys, json
from typing import List, Union

# Universal input parser: automatically parse JSON, int, float, space-separated
# list, etc.
def parse_value(s: str) -> Union[int, float, str, list, dict, bool]:
    s = s.strip()
    # 1. Try JSON parsing (handles list, dict, quoted string, boolean, number)
    try:
        return json.loads(s)
    except json.JSONDecodeError:
        pass
    # 2. Integer detection
    if s.isdigit() or (s.startswith('-') and s[1:].isdigit()):
```

```

1026     return int(s)
1027 # 3. Float detection
1028     try:
1029         return float(s)
1030     except ValueError:
1031         pass
1032 # 4. Split on spaces *only* if the value is not enclosed in quotes
1033     if '_' in s and not (s.startswith('"') and s.endswith('"')):
1034         return [parse_value(part) for part in s.split()]
1035 # 5. Strip surrounding quotes and keep inner spaces
1036     if s.startswith('"') and s.endswith('"'):
1037         return s[1:-1]
1038 # Fallback: return raw string
1039     return s
1040
1041 # Replace this class with your own logic implementation
1042 class Solution:
1043     def solved_fun(self, *args):
1044         pass # TODO: Replace with your actual implementation
1045
1046     def solve(self, *args) -> Union[int, str]:
1047         return self.solved_fun(*args)
1048
1049 # Main function: handles most common competitive input/output formats
1050 def main():
1051     # Read all non-empty lines from stdin
1052     lines = [parse_value(line.strip()) for line in sys.stdin if line.strip()]
1053
1054     # stdin template
1055     {stdin_template}
1056
1057 # Python script entry point
1058 if __name__ == "__main__":
1059     main()

```

This Python template is designed to handle a wide variety of standard input/output formats common in programming contests. You should edit this template in three specific places when solving a new problem.

### 1. import Section

If your solution requires additional Python libraries (e.g., math, collections, etc.), please add them to the import section at the top of the file.

```

import sys, json
from typing import List, Union
# Add additional imports here if needed
# e.g., from collections import defaultdict

```

### 2. Write Your Solution Inside the Solution Class

Implement your logic inside the solved\_fun() method. Rename the method according to the problem's context (e.g., def count\_pairs(self, ...)). Make sure to update the call inside the solve() method accordingly.

```

class Solution:
    def your_function_name(self, ...): # <-- Rename this function
        # Implement your solution here
        return ...

    def solve(self, *args) -> Union[int, str]:
        return self.your_function_name(*args) # <-- Make sure this matches

```

1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133

### 3. Customize the main() Function's Input Parser

You **must retain** the following line for reading input:

```
lines = [parse_value(line.strip()) for line in sys.stdin if line.strip()]
```

**Inspect the problem's input format** carefully:

- If the problem provides multiple lines of input, adjust how the input is grouped.
- If you need to handle multiple test cases, add a loop around the function calls.
- If input is structured (e.g., a matrix or multiple arrays), modify the logic to parse accordingly.

```
def main():
    lines = [parse_value(line.strip()) for line in sys.stdin if line.
              strip()]
    # stdin template # <-- Make sure the following template matches your
    # problem's input format
    ...
    # print result # <-- Make sure the following template matches your
    # problem's output format
```

## D.4 THE PROMPT OF THE CODING AGENT IN THE DEEP PLANNING SYSTEM

### The prompt of the Coding Agent in the Deep Planning System. (In Step 3)

A step-by-step plan has already been created to solve this problem. Please follow the instructions below carefully:

PROBLEM

{problem}

PLANNING REPORT

{plan}

INSTRUCTIONS

1. Implement the code strictly according to the provided plan.
2. Each part of the code must clearly indicate which step in the plan it corresponds to.
3. Add meaningful in-line comments in the code to explain the logic.
4. Please provide the final complete code; do not just return the class solution.

{prompt\_of\_code\_template}

OUTPUT FORMAT

CODE

```
# Your code here, following the plan step by step and using the format provided
```

### Important Instructions:

- Strictly follow the instructions.
- Please modify the code format according to the format of **Code Format Sample**.

1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187

## D.5 DEBUGGING AGENT PROMPT IN THE DEEP PLANNING SYSTEM

### Debugging Agent Prompt for STD I/O Errors in the Deep Planning System. (In Step 4)

You are given the following problem description and Python code.

**Problem:** {Problem}

**Code:**

{Code}

Please focus only on fixing the `main()` function.

Your task is to **rewrite only the `main()` function** so that the program correctly handles input and output according to the problem description and the design of the `Solution` class.

**Requirements:**

- Do **not** modify any part of the `Solution` class.
- You **must retain** the following line for reading input:

```
lines = [parse_value(line.strip()) for line in sys.stdin if line.strip()]
```

- Based on the `lines` list and the problem's input format, write the rest of the `main()` logic so that it:
  1. Correctly parses the input.
  2. Calls the appropriate method(s) from `Solution`.
  3. Prints the result in the correct format.

Below is a sample input/output format you should follow:

{sample\_io}

Return only the modified `main()` function code. Do not return any explanations or other parts of the script.

### Debugging Agent Prompt for Runtime Errors in the Deep Planning System. (In Step 4)

You are a programmer who has received a solution of a problem written in **{language}**, but it fails to compile or throws an immediate runtime error. Your task is to fix the code so that it compiles and runs correctly.

BUGGY CODE

{code}

{test\_log}

**Expected Output:**

Your response must be structured as follows:

DIAGNOSIS

- Analyze the error message and buggy code.
- Identify the cause of the failure (e.g., syntax error, missing variable, etc.).
- Explain how to fix it concisely.

MODIFIED CODE

- Please keep all program comments.

1188  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241

- If there are errors, please also update the comments.

```
# Corrected code with concise comments on the fix.
```

#### Important Instructions:

- Strictly follow the instructions.
- Do not change the algorithm unless necessary to fix the compile error.
- Focus on fixing syntax or execution-level issues only.
- Do not add testing code, for example, `assert` statements in your code.

#### Debugging Agent Prompt for Wrong Answer Cases in the Deep Planning System. (In Step 4)

You are a programmer who needs to fix a solution to a problem written in `{language}`. The code runs, but produces incorrect results on certain test cases.

PROBLEM

```
{problem_with_planning}
```

BUGGY CODE

```
{code}
```

TEST CASE FAILURE LOG

```
{test_log}
```

#### Expected Output:

ANALYZE THE PLAN AND TEST FAILURE

- Read the problem plan and failed test case log.
- Simulate the plan's logic using one failed test input.
- Determine what the correct output should be.
- Identify what kind of mistake likely occurred (e.g., wrong algorithm, rounding error, missing constraint).

ROOT CAUSE & STRATEGY

- Clearly describe the most likely root cause.
- Do not patch old logic—redesign a **better and more reliable solution** if needed.
- Mention if precision issues suggest using better tools (e.g., decimal instead of float).

IMPROVED PLAN

- Describe an improved plan or algorithm to solve the problem more accurately and reliably.

CODE

- Please keep all program comments.
- If there are errors, please also update the comments.

```
# Rewritten code based on improved plan, using robust logic and tools.
```

#### Important Instructions:

- Strictly follow the instructions.

1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295

- Do not copy old logic.
- Be open to changing tools (e.g., decimal) or redesigning the algorithm if needed.
- Do not add testing code, for example, `assert t` statements in your code.

## E ROLE OF LLMs IN THIS RESEARCH

In the preparation of this manuscript, we utilized several Large Language Models (LLMs), including OpenAI’s GPT-5, Google’s Gemini 2.5 Pro, and Anthropic’s Claude 4 Sonnet, for two primary purposes. First, we employed these models to assist with proofreading and enhancing the grammatical clarity and readability of the text. Second, they were used as coding assistants to generate boilerplate code and utility functions within our implementation. The core research concepts, experimental design, and scientific analysis presented in this paper were conceived and executed entirely by the human authors. Our open-source code, partially developed with LLM assistance, is available to ensure full transparency and reproducibility.