

# SURGE: On the Potential of Large Language Models as General-Purpose Surrogate Code Executors

Anonymous ACL submission

## Abstract

Neural surrogate models are powerful and efficient tools in data mining. Meanwhile, large language models (LLMs) have demonstrated remarkable capabilities in code-related tasks. To this end, a novel application of LLMs emerges—using LLMs as surrogate models for code execution prediction. Given LLMs’ unique ability to understand and process diverse programs, they present a promising direction for building general-purpose surrogate models. To systematically investigate this capability, we introduce SURGE, a comprehensive benchmark with 1160 problems covering 8 key aspects: multi-language programming tasks, competition-level programming problems, repository-level code analysis, high-cost scientific computing, time-complexity-intensive algorithms, buggy code analysis, programs dependent on specific compilers or execution environments, and formal mathematical proof verification. Through extensive analysis of 21 open-source and proprietary LLMs, we examine scaling laws, data efficiency, and predictive accuracy. Our findings reveal important insights about the feasibility of LLMs as efficient surrogates for computational processes.

## 1 Introduction

Neural surrogate models (Zhang et al., 2024; Sun and Wang, 2019) are powerful tools in data mining and machine learning, which efficiently approximate complex computational processes. Meanwhile, Large language models (LLMs) (Reid et al., 2024; Meta, 2024; Anthropic, 2024b; Hui et al., 2024; Bi et al., 2024) have demonstrated remarkable capabilities in code-related tasks (Lu et al., 2021a; Zheng et al., 2023; Luo et al., 2023; Team, 2024a; Guo et al., 2024), including code understanding (Ahmad et al., 2020; Chakraborty et al., 2020) and code generation (Li et al., 2018a; Parvez et al., 2018). However, an equally important yet underexplored question is whether LLMs can serve as

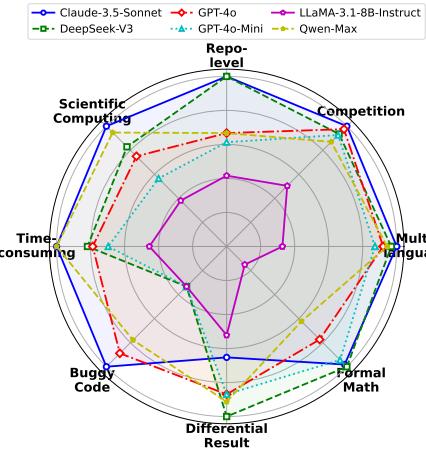


Figure 1: Performance of 6 typical models on 8 subsets of SURGE.

general-purpose surrogate code executors, which predict the behavior of a program without actually running it. A recent study (Lyu et al., 2024) acknowledges its importance, however, it focuses on a case study rather than a systematic analysis.

The ability to predict code execution outcomes without execution has tremendous significance. In scientific computing, running simulations often requires substantial computational resources and is time-consuming, making it impractical to test every possible configuration (Lu and Ricciuto, 2019; Hesthaven and Ubbiali, 2018; Benner et al., 2015). In security-sensitive environments, executing untrusted code poses inherent risks, necessitating alternative mechanisms for assessing program behavior without exposing the system to potential vulnerabilities (Nebbione and Calzarossa, 2023; Shirazi et al., 2017; Wang et al., 2024). Additionally, some code requires highly specific execution environments, which may not always be available, making surrogate execution a valuable alternative (Queiroz et al., 2023; Gu et al., 2025). Moreover, accurately predicting a model’s potential outputs or errors is crucial for improving traditional tasks such as code understanding, code generation, and even math rea-

soning (Li et al., 2025). Lastly, many works use LLMs as reward models (RMs) in reinforcement learning. For code tasks, accurate execution prediction is key to a reliable RM (Ouyang et al., 2022).

Traditional approaches to surrogate code executing (King, 1976; Cedar and Sen, 2013) struggle to generalize across languages and suffer from scalability issues when applied to complex real-world codebases. Containerized environments (Merkel, 2014) mitigate dependency issues but still require full code execution. Recent efforts to train neural executors (Yan et al., 2020) focus on narrow tasks and lack the generality needed for real-world code. In contrast, LLMs’ capacity to internalize patterns from vast code corpora (Lu et al., 2021b; Chaudhary, 2023) suggests a path toward general-purpose surrogate code execution.

To understand the potential of LLMs as GEneral-purpose SURrogate code executors, we introduce **SURGE**. It includes 8 components: (1) fundamental programming tasks in multiple languages, (2) competition programming problems requiring deep logical inference, (3) repository-level codebases that test long-range dependencies, (4) scientific simulations and optimizations where direct execution is high-cost, (5) time-consuming logical algorithms that have high time-complexity, (6) buggy code that examines LLMs’ ability to predict runtime errors, (7) programs whose behavior depends on specific compiler versions or execution environments and (8) math theorem proving in formal language (De Moura et al., 2015; Moura and Ullrich, 2021) which expects compilers to testify.

Through extensive evaluation of 21 open-source and proprietary LLMs on SURGE, we provide the first large-scale study of LLMs’ capabilities as computational surrogates. Additionally, we investigate the impact of various factors, including prompt engineering strategies, programming language characteristics, computational complexity, and execution time requirements, on surrogate performance. Our findings reveal both the promising potential and current limitations of LLMs as general-purpose code execution surrogates. The performance of typical models on SURGE is shown in Figure 1.

Beyond benchmarking, we conduct a scaling law study on whether LLMs’ performance improves with model size and training data. We train models with 4 different sizes on different scales of training data from the formal language subset of SURGE. Our experiments demonstrate that models’ performance

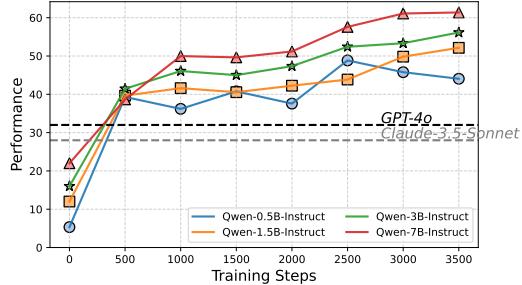


Figure 2: Performance scaling across model sizes and training steps.

consistently improves with both model size and training steps, with larger models showing stronger learning capacity and higher performance ceilings throughout the training process (Figure 2).

In short, our work makes the following key contributions:

- We introduce SURGE, the first holistic benchmark for evaluating LLMs as general-purpose surrogate code executors. It consists of 8 subsets and 1160 problems.
- We evaluate 21 open-source and proprietary LLMs on SURGE and conduct the first large-scale analysis on them.
- We present a scaling law study with models of varying sizes and scales of training data, providing empirical insights on the scaling law of LLMs on these tasks.

## 2 Related Works

**Neural Surrogate Models.** Neural surrogate models are neural network-based approximations used to replace computationally expensive simulations in various scientific and engineering domains (Zhang et al., 2024; Sun and Wang, 2019). These models act as domain-specific emulators by learning complex input-output relationships from high-fidelity data(Raissi et al., 2020; Sun et al., 2020; Bessa et al., 2017; Thuerey et al., 2020; Raissi et al., 2019; Willard et al., 2022). Recently, generative models have been incorporated into surrogate modeling. Some equip language models with traditional surrogate models to facilitate iterative optimization (Ma et al., 2024; Lyu et al., 2025), and some use generative models to realize the end-to-end surrogate process (Gruver et al., 2024; Hao et al., 2024; Wimmer and Rekabsaz, 2023; Che et al., 2024). While these studies primarily focus on natural sciences, time series, and multimodal gaming, the application of surrogate modeling to code execution, where both input and output exist

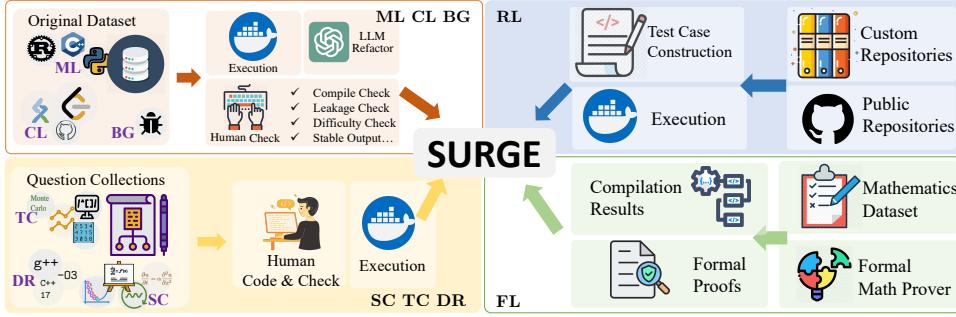


Figure 3: The Construction of SURGE employs 4 methodologies: 1. Iterative Refactor, 2. Repository Sampling, 3. Manual Implementation, and 4. Inference & Verification.

in the modality of language, remains unexplored.

**LLMs for Code.** LLMs are widely used in code-related tasks (Lu et al., 2021a; Luo et al., 2023; Team, 2024a; Guo et al., 2024), which can be fundamentally categorized into code understanding and code generation. Code understanding tasks include code summarization (Hu et al., 2018; Harer et al., 2019; Ahmad et al., 2020), bug detection (Li et al., 2018b; Russell et al., 2018; Zhou et al., 2019; Chakraborty et al., 2020), duplication detection (Zhang et al., 2019; Yu et al., 2019; Wang et al., 2020), code retrieval (Husain et al., 2020; Lu et al., 2021a), etc. Code generation tasks include code completion (Li et al., 2018a; Parvez et al., 2018), code repair (Chen et al., 2019; Chakraborty et al., 2020; Lutellier et al., 2020), test generation (Watson et al., 2020; Siddiq et al., 2024; Schäfer et al., 2023), etc. However, while the potential execution result of code is important for both code understanding and generation, this aspect remains largely unexplored (Weber et al., 2024).

### 3 SURGE

SURGE assesses the model’s ability to approximate execution results across multiple dimensions, including multi-lingual diversity, repository-level complexity, computational intensity, error handling, and scenario-dependent variability. Below, we describe each component of SURGE, including motivation and construction methods.

#### 3.1 Dataset Construction

As illustrated in Figure 3, the construction of SURGE involves four distinct methodologies, each applied to specific components of our eight subsets. For **ML**, **CL**, **BG** components, we employ the Iterative Refactor methodology, where we refine the code through an interactive process involving LLM assistance and human verification. The **RL** component

is constructed through Repository Sampling, where we extract and construct test cases from both public and custom code repositories. For **SC**, **TC**, **DR** components, we utilize Manual Implementation, carefully handcrafting code based on selected textbook materials and question collections. Finally, the **FL** component is developed using the Inference & Verification approach, leveraging formal mathematical provers to generate proofs and validate them through compiler verification.

#### 3.2 Dataset Components

**Multi-lingual Code (ML).** A fundamental nature of a general-purpose surrogate executor is its ability to handle multiple programming languages especially computational languages, rather than rendering languages. Our dataset covers 7 such languages, including C, C++, C#, Java, Rust, Python, and Julia. Our dataset is adapted from McEval (Chai et al., 2024). The original dataset does not provide executable code, so we used an LLM to generate executable code by providing it with prompts, ground truth, and test cases in the original dataset. This generated code was then manually processed.

**Competition-level Code (CL).** Next, we consider competition-level code, which presents a higher level of coding difficulty. We collect these tasks from 2 public repositories <sup>12</sup>, which contain problems from open coding platforms (e.g. LeetCode, Luogu). The dataset includes problems in 3 languages, C++, Java, and JavaScript. Since the original repositories only provide partial solutions, we first use an LLM to generate complete, executable code that prints the expected output. This generated code is then manually verified. To investigate whether problem difficulty affects the

<sup>12</sup><https://github.com/az1397985856/leetcode>

<sup>2</sup><https://gitee.com/shanire/OJCode>

Table 1: Statistics of SURGE, including construction methods, problem sources, quantities, number of examples in few-shot scenarios, evaluation metrics, and criteria for further classification. In the table, “custom” refers to customized approaches, and “mixed” indicates multiple methods, which are elaborated in detail in the text.

Subset	Construction Method	Source	Quantity	Example Num.	Metric	Categories
ML	Iterative Refactor	McEval (Chai et al., 2024)	150	3	Exact Match	Languages
CL	Iterative Refactor	GitHub	150	3	Exact Match	Difficulties & Languages
RL	Repository Sampling	GitHub & Custom	60	3	Mixed	-
SC	Manual Implementation	Custom	150	3	Mixed	Scenarios
TC	Manual Implementation	Custom	150	3	Mixed	Scenarios, CPU Time
BG	Iterative Refactor	DebugBench (Tian et al., 2024)	150	4	Jaccard similarity	Error Type & Languages
DR	Manual Implementation	Custom	200	3	Jaccard similarity	Variable Type
FL	Inference & Verification	Lean-Workbook (Ying et al., 2024)	150	3	Custom	-

230 performance of surrogate models, we further employ an LLM to automatically classify problems  
 231 into 5 different difficulty levels, followed by human  
 232 verification to ensure accuracy.  
 233

234 **Repository-level Code (RL).** In real-world sce-  
 235 narios, most code exists at the repository level,  
 236 making repository-level code execution prediction  
 237 equally important for a general-purpose surrogate  
 238 model. We manually collect computational reposi-  
 239 tories that fit within the input length constraints  
 240 of LLMs. These repositories include tasks such as  
 241 solving the 24-point problem, Sudoku solving, and  
 242 converting Python code to LaTeX. These tasks ex-  
 243 hibit complex logic but do not rely on sophisticated  
 244 models or external inputs. To assess the model’s  
 245 ability to understand multi-file structures, we also  
 246 manually construct two repositories containing ad-  
 247 vanced C++ syntax and multiple files.

248 **Scientific Computing (SC).** Scientific comput-  
 249 ing has long been adopting neural surrogate  
 250 models. We introduced tasks ranging from solving ordi-  
 251 nary differential equations (ODEs) to optimization  
 252 problems and signal processing. These tasks are  
 253 motivated by and widely used in real-world sci-  
 254 entific challenges, including areas where increasing  
 255 research has been done on solving these compu-  
 256 tational tasks through building efficient surrogate  
 257 models (Wu et al., 2023; Zhang et al., 2023). A  
 258 comprehensive overview of the setup for each task,  
 259 along with the corresponding algorithms, can be  
 260 found in Appendix C.4.1.

261 **Time-Consuming Algorithms (TC).** Surrogate  
 262 models were originally motivated by real-world  
 263 applications where program execution is high-cost  
 264 and time-consuming. It’s a necessity for LLMs  
 265 to generalize well to strongly computation-power-  
 266 dependent and time-consuming tasks. We include  
 267 examples from linear algebra, sorting, searching,  
 268 Monte Carlo simulations, and string matching pro-

grams, ensuring a broad representation of computa-  
 269 tionally intense tasks. These tasks cover various  
 270 complexity classes, including P (e.g., sorting an  
 271 array), NP (e.g., Hamiltonian Cycle), and NP-Hard  
 272 (e.g., Traveling Salesman’s Problem). Additionally,  
 273 we record the CPU execution time for each pro-  
 274 gram under consistent environments individually  
 275 to support subsequent studies. A detailed descrip-  
 276 tion of this subset is provided in Appendix C.5.1.  
 277

278 **Buggy Code (BG).** Real-world code execution  
 279 often encounters errors, which pose risks in sen-  
 280 sitive scenarios. Therefore, code surrogate mod-  
 281 els should recognize the presence of bugs. This  
 282 dataset is adapted from DebugBench (Tian et al.,  
 283 2024), which extracted Java, Python, and C++ code  
 284 from LeetCode and manually inserted errors from  
 285 4 major bug categories and 18 minor types. Since  
 286 DebugBench only provides code snippets rather  
 287 than complete executable programs, we first used  
 288 an LLM to automatically complete the error-free  
 289 code into fully runnable versions. After verifying  
 290 their correctness, we replaced relevant parts with  
 291 buggy code and executed them again to capture  
 292 the corresponding error outputs. Some errors re-  
 293 sulted in infinite loops, causing timeouts, so we set  
 294 a 30-second execution filter to avoid such cases.  
 295

296 **Code with Differential Results under Different  
 297 Scenarios (DR).** Various contextual factors, such  
 298 as compiler versions, optimization levels, and lan-  
 299 guage standards, often influence code execution.  
 300 These variations can lead to different outputs for  
 301 the same code snippet. We focus specifically on  
 302 C++ and manually collect code snippets from text-  
 303 books and online sources (Bryant and O’Hallaron,  
 304 2010; Lippman et al., 2012) that exhibit differ-  
 305 ent behaviors under varying compilation settings.  
 306 We consider multiple compilers (g++, clang++),  
 307 C++ standards (03, 11, 14, 17), and optimization  
 308 levels (-O0, -O1, -O2, -O3, -Os). Each snip-  
 309 petit is executed across these different settings, and  
 310

309 we retain only those that produce varying outputs  
310 through different configurations while discarding  
311 cases that yield identical results across all settings.

312 **Mathematics Formal Language (FL).** Math-  
313 Proving Formal Languages are specialized pro-  
314 gramming languages designed for mathematical  
315 proof verification through compilers (De Moura  
316 et al., 2015; Moura and Ullrich, 2021; Paulson,  
317 1994; Barras et al., 1997). These compilers can  
318 determine whether a proof is correct and iden-  
319 tify specific errors. The verification is very time-  
320 consuming. In this study, we focus on Lean4  
321 which is the most widely used proof language. To  
322 build our dataset, we use Goedel-Prover (Lin et al.,  
323 2025) to conduct large-scale reasoning on Lean-  
324 Workbook (Ying et al., 2024) and extract an equal  
325 proportion of correct and incorrect proofs. This  
326 balanced dataset allows us to evaluate the surrogate  
327 model’s ability to assess proof validity effectively.

### 328 3.3 Evaluation Metrics

329 We design different evaluation metrics tailored to  
330 each subset of SURGE to ensure accurate evaluation.

331 In **ML** and **CL**, the outputs are simple numerical  
332 values or formatted strings, we employ exact string  
333 matching to measure correctness.

334 For **RL**, we employ different evaluation methods  
335 for different tasks. For structured C repositories,  
336 we use exact character matching to compare out-  
337 puts. For Sudoku and 24-point problems, we use  
338 edit distance to compare results. For other types of  
339 repositories, we apply the Ratcliff/Obershelp (Rat-  
340 cliff and Metzener, 1988) algorithm.

341 For **SC** and **TC**, various tasks necessitate dis-  
342 tinct evaluation methods. Specifically, (1) numer-  
343 ical simulations are evaluated using the average  
344 Relative Absolute Error (RAE); (2) position-based  
345 tasks, such as binary search, are assessed through  
346 exact string matching; and (3) sorting tasks are eval-  
347 uated by the rank correlation coefficient (Spearman,  
348 1904). Details regarding the evaluation metrics can  
349 be found in Appendix C.4.2 and Appendix C.5.2.

350 For **BG**, we use the Jaccard similarity (Jaccard,  
351 1901) between predicted and ground truth error  
352 messages. For **DR**, since the same code can pro-  
353 duce different outputs in varying settings, which  
354 sometimes include warnings or errors, we again  
355 utilize Jaccard similarity.

356 For **FL**, the results consist of two parts: (1)  
357 whether the proof passes or not, and (2) if it fails,  
358 we evaluate the accuracy of the error message. The

359 error message consists of a list containing the error  
360 locations and descriptions. We compute the score  
361 of a prediction as  $\frac{1}{N} \sum_{j=1}^N \mathbb{1}[\hat{p}_j \in P] \cdot J(\hat{m}_j, m_j)$ ,  
362 where  $N$  is the number of errors in the ground  
363 truth,  $P$  is the set of predicted error positions,  $p_j$   
364 represents the  $j$ -th ground truth error position,  $\hat{p}_j$   
365 represents the predicted error position correspond-  
366 ing to  $p_j$ ,  $\mathbb{1}[\hat{p}_j \in P]$  is the indicator function which  
367 equals to 1 only when there exists  $\hat{p}_j \in P$ ,  $m_j$  is  
368 the ground truth error message for position  $p_j$ ,  $\hat{m}_j$   
369 is the predicted error message for position  $\hat{p}_j$ , and  
370  $J$  is the Jaccard similarity function.

## 371 3.4 Dataset Statistics

372 Table 1 presents detailed statistics of SURGE, in-  
373 cluding the construction methods, problem sources,  
374 dataset quantities, number of examples in few-shot  
375 scenarios, evaluation metrics, and classification cri-  
376 teria for each subset.

## 377 4 Experiments

### 378 4.1 Setup

379 **Models.** We tested SURGE on 17 open-source  
380 and 4 closed-source models of different sizes,  
381 including both chat models and code models. The  
382 closed-source models include GPT-4o (2024-  
383 08-06) (OpenAI, 2024b), GPT-4o-mini (2024-  
384 07-18) (OpenAI, 2024a), Claude-3.5-  
385 Sonnet (2024-10-22) (Anthropic, 2024a),  
386 and Qwen-Max (2025-01-25) (Team, 2024b).  
387 The open-source models include LLaMA-3.1-  
388 {8, 70}B-Instruct, LLaMA-3.3-70B-  
389 Instruct, Qwen-2.5-{0.5, 1.5, 3,  
390 7, 14, 32, 72}B-Instruct, Qwen-  
391 2.5-Coder-{0.5, 1.5, 3, 7, 14,  
392 32}B-Instruct and DeepSeek-V3 (671B).

393 **Settings.** We tested the above models on SURGE  
394 under 3 settings: 0-shot w/o CoT, 0-shot w/ CoT,  
395 and few-shot w/ CoT. CoT here means whether we  
396 use Chain-of-Thought (Wei et al., 2022) prompting,  
397 allowing the models to think step by step, or ask the  
398 models to answer directly. We set the temperature  
399 to 0, i.e., employing greedy decoding.

### 400 4.2 Results

401 Table 2 presents the performance of 10 selected  
402 models across 8 sub-datasets of SURGE under 3 dif-  
403 ferent settings. In this table, we provide a detailed  
404 breakdown of the models’ performance on the ML  
405 and BG sub-datasets. We present the complete re-

Table 2: Performance of different models under different prompting strategies on SURGE.

Model	ML						CL	RL	SC	TC	BG					
	CPP	Rust	Python	Julia	Java	Others					CPP	Java	Python	DR		
<i>Zero-shot</i>																
Claude-3.5-Sonnet	72.73	55.00	88.00	66.67	76.92	75.00	81.58	57.31	61.55	35.27	9.09	12.55	51.40	12.92	17.92	51.59
DeepSeek-V3	54.55	60.00	76.00	61.11	46.15	72.50	56.58	44.19	59.65	35.31	4.05	3.03	21.50	10.92	32.46	42.53
GPT-4o	40.91	45.00	60.00	55.56	57.69	55.00	66.45	49.13	53.16	34.44	4.28	7.48	34.59	14.75	21.99	40.03
Qwen-Max	50.00	45.00	44.00	27.78	26.92	50.00	38.82	37.15	56.98	35.44	2.82	3.20	30.52	14.03	29.89	32.84
Qwen-2.5-7B-Instruct	13.64	5.00	12.00	5.56	11.54	17.50	27.63	27.98	22.90	29.43	2.21	3.66	9.46	7.24	36.42	15.48
Qwen-2.5-32B-Instruct	45.45	20.00	48.00	33.33	42.31	20.00	50.66	22.61	32.50	30.99	1.21	2.09	12.16	7.26	7.65	25.08
Qwen-2.5-Coder-7B-Instruct	45.45	30.00	52.00	44.44	50.00	42.50	75.00	20.86	45.50	28.00	0.86	3.47	13.53	14.23	40.40	33.75
Qwen-2.5-Coder-32B-Instruct	59.09	55.00	60.00	44.44	69.23	60.00	80.26	48.43	57.18	18.84	1.49	1.95	17.96	13.42	29.19	41.10
LLAMA-3.1-8B-Instruct	0.00	0.00	4.00	5.56	15.38	5.00	13.16	20.41	4.41	15.46	4.41	4.12	5.98	3.97	0.00	8.49
LLAMA-3.1-70B-Instruct	54.55	40.00	52.00	33.33	65.38	47.50	78.29	31.60	48.65	30.19	1.59	4.18	14.27	15.73	39.16	37.10
<i>Zero-shot Chain-of-Thought</i>																
Claude-3.5-Sonnet	90.91	65.00	96.00	77.78	69.23	92.50	82.24	62.31	63.38	40.70	16.91	20.69	62.23	18.19	33.98	59.47
DeepSeek-V3	81.82	85.00	88.00	72.22	69.23	85.00	76.32	62.70	57.57	36.71	4.45	7.85	46.26	16.21	35.19	54.97
GPT-4o	68.18	65.00	92.00	72.22	76.92	77.50	79.61	53.74	48.56	28.36	8.19	9.97	44.29	14.21	27.91	51.11
Qwen-Max	86.36	75.00	80.00	72.22	76.92	80.00	71.05	50.49	61.78	36.71	2.65	7.73	46.85	16.16	20.74	52.31
Qwen-2.5-7B-Instruct	40.91	15.00	32.00	33.33	26.92	47.50	52.63	27.51	25.68	28.95	1.12	3.76	14.94	9.41	36.46	26.41
Qwen-2.5-32B-Instruct	50.00	40.00	40.00	55.56	38.46	57.50	53.29	32.71	43.29	30.59	3.10	6.96	23.86	11.49	7.39	32.95
Qwen-2.5-Coder-7B-Instruct	68.18	40.00	40.00	38.89	53.85	57.50	46.05	19.70	40.91	30.19	2.29	4.71	12.77	15.04	37.12	33.81
Qwen-2.5-Coder-32B-Instruct	77.27	65.00	80.00	55.56	73.08	67.50	71.71	54.58	55.69	34.36	2.05	4.74	22.43	17.62	28.55	47.34
LLAMA-3.1-8B-Instruct	40.91	15.00	24.00	22.22	26.92	30.00	41.45	17.87	32.65	18.22	1.52	4.23	13.38	10.12	0.66	19.94
LLAMA-3.1-70B-Instruct	59.09	50.00	72.00	61.11	57.69	52.50	58.55	34.44	43.93	29.76	1.71	3.49	15.02	16.86	25.85	38.80
<i>Few-shot Chain-of-Thought</i>																
Claude-3.5-Sonnet	86.36	70.00	96.00	72.22	65.38	82.50	82.24	70.65	63.58	41.00	22.04	23.61	44.15	25.70	31.99	58.49
DeepSeek-V3	90.91	65.00	84.00	77.78	73.08	95.00	80.26	78.64	66.00	38.60	21.98	15.14	40.27	24.38	35.17	59.08
GPT-4o	68.18	60.00	88.00	77.78	73.08	75.00	75.66	76.86	59.65	37.12	12.91	7.74	29.52	22.08	26.65	52.68
Qwen-Max	81.82	70.00	88.00	77.78	73.08	80.00	82.24	72.53	62.32	37.88	19.68	19.78	37.57	23.91	24.76	56.76
Qwen-2.5-7B-Instruct	27.27	25.00	36.00	38.89	26.92	42.50	48.68	45.19	43.42	28.97	4.92	4.70	12.94	10.66	34.53	28.71
Qwen-2.5-32B-Instruct	59.09	55.00	52.00	66.67	53.85	60.00	63.16	64.10	63.12	32.53	5.41	6.94	28.66	13.81	22.87	43.15
Qwen-2.5-Coder-7B-Instruct	54.55	30.00	36.00	44.44	50.00	42.50	58.55	50.48	53.91	30.69	3.90	4.93	14.44	14.02	25.25	34.24
Qwen-2.5-Coder-32B-Instruct	68.18	75.00	72.00	55.56	65.38	70.00	76.97	64.16	56.37	34.34	3.93	6.49	20.22	19.09	22.57	47.35
LLAMA-3.1-8B-Instruct	13.64	20.00	20.00	5.56	26.92	22.50	30.26	34.15	47.89	22.27	4.44	4.55	9.66	12.05	13.25	19.14
LLAMA-3.1-70B-Instruct	54.55	35.00	40.00	22.22	53.85	35.00	68.42	50.83	60.96	30.29	8.00	5.95	18.32	13.27	33.12	35.32

sults of all 21 models in Table 5 in Appendix A. From the results, several notable findings emerge:

**SURGE demonstrates strong discriminative ability.** Even the strongest models perform only moderately well, highlighting the value of our benchmark. The models exhibit significant performance differences across different subsets, reflecting the comprehensiveness of our dataset. Additionally, models that perform well in other tasks, such as Claude-3.5-Sonnet, also achieve substantial overall results on SURGE. This demonstrates the reasonableness of our dataset and its effectiveness in benchmarking LLMs as general-purpose surrogate code executors.

**Different prompting strategies lead to varying model performance and have different effects across subsets.** We found that for the task of code execution surrogacy, both Chain-of-Thought prompting and few-shot learning can enhance model performance.

**Larger model sizes tend to yield better performance on SURGE.** From the results, we observe that regardless of whether it is a Qwen-Chat model or a Qwen-Coder model, performance improves as the parameter size increases.

**Chat models and coder models exhibit different performance patterns and are affected differ-**

**ently by prompting strategies.** We observed that for chat and code models of the same size, code models outperform chat models in the zero-shot setting SURGE. However, in the other two settings, chat models perform better. This suggests that code models have stronger zero-shot surrogate capabilities, whereas chat models excel in reasoning and imitation abilities.

**On some sub-datasets of SURGE, stronger models perform worse than smaller, weaker ones.** For example, in the FL dataset, we found that this occurs because stronger models tend to actively look for errors in the code, often misidentifying correct code as incorrect. In contrast, smaller models are more inclined to assume that all code is error-free. Since half of the samples in this subset are indeed correct, the smaller models end up achieving better performance.

## 5 Analysis

### 5.1 The Impact of Language Type

In Table 2, we compare model performance across different programming languages.

In the ML subset, LLMs perform best in Python, followed by C++. Python’s simple syntax makes it easier to process, while C++ benefits from its strong presence in programming and system-level

Table 3: Model Performance Across Different Variables in DR

Model	Compiler			Standard			Optimization		
	Zero-shot	Zero-shot CoT	Few-shot CoT	Zero-shot	Zero-shot CoT	Few-shot CoT	Zero-shot	Zero-shot CoT	Few-shot CoT
Claude-3.5	12.92	18.19	25.70	14.21	16.75	23.59	16.06	1.23	12.90
GPT-4o	14.75	14.21	22.08	15.33	14.61	20.02	7.66	1.63	1.22
LLaMA-3.1-8B	3.97	10.12	12.05	4.29	10.60	13.17	1.96	3.91	8.04
LLaMA-3.1-70B	15.73	16.86	13.27	16.36	16.27	13.24	15.04	2.24	3.73

459 tasks. Models also perform well in Julia, likely  
 460 due to its clean syntax and similarity to Python.  
 461 However, performance drops significantly in Rust,  
 462 where strict ownership and lifetime rules introduce  
 463 complexity, making it hard to predict.

464 In the BG subset, when predicting the output of  
 465 buggy code, models excel in Python but struggle  
 466 with C++. Python’s clear error messages aid prediction,  
 467 whereas C++’s static typing, manual memory  
 468 management, and potential for undefined behavior  
 469 make error handling more difficult.

Table 4: Model Performance by Difficulty Level

Difficulty	1	2	3	4	5
<i>Zero-shot</i>					
Claude-3.5-Sonnet	90.91	81.25	82.05	72.00	79.49
GPT-4o	72.73	56.25	58.97	84.00	61.54
LLaMA-3.1-8B-Instruct	24.24	18.75	7.69	4.00	12.82
LLaMA-3.1-70B-Instruct	96.97	93.75	64.10	64.00	79.49
<i>Zero-shot Chain-of-Thought</i>					
Claude-3.5-Sonnet	96.97	87.50	76.92	72.00	79.49
GPT-4o	96.97	81.25	74.36	80.00	69.23
LLaMA-3.1-8B-Instruct	57.58	43.75	30.77	40.00	38.46
LLaMA-3.1-70B-Instruct	90.91	56.25	48.72	36.00	56.41
<i>Few-shot Chain-of-Thought</i>					
Claude-3.5-Sonnet	96.97	81.25	76.92	68.00	84.62
GPT-4o	93.94	81.25	71.79	60.00	71.79
LLaMA-3.1-8B-Instruct	39.39	18.75	33.33	16.00	33.33
LLaMA-3.1-70B-Instruct	87.88	87.50	48.72	48.00	76.92

## 5.2 The Impact of Problem Difficulty

471 To analyze how problem difficulty affects model  
 472 performance, we examine results in the CL subset,  
 473 as shown in Table 4.

474 Across all settings, model performance generally  
 475 declines as problem difficulty increases. How-  
 476 ever, at the highest difficulty level, performance  
 477 improves slightly. This anomaly arises because  
 478 difficulty levels are categorized based on coding  
 479 complexity, but the hardest problems often involve  
 480 implementing simple functionality using complex,  
 481 optimized algorithms. In such cases, models can  
 482 generate correct answers through end-to-end rea-  
 483 soning without fully understanding the underlying  
 484 code logic. In contrast, levels 1–4 require logical  
 485 reasoning based on code execution, making prediction  
 486 harder as complexity increases.

487 Furthermore, while Chain-of-Thought prompt-  
 488 ing improves performance, few-shot learning does

489 not and may even degrade results. This is likely be-  
 490 cause buggy code varies widely, causing few-shot  
 491 examples to mislead rather than aid the model.

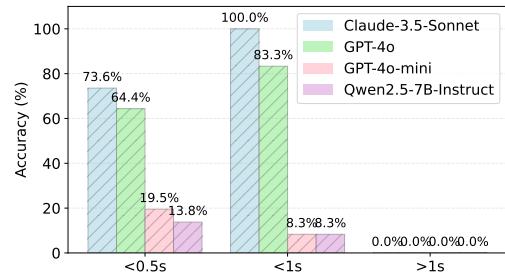


Figure 4: Model’s Performance on TC subset across programs with different run time on CPU.

## 5.3 Prediction Accuracy & CPU Time

492 We explore the relationship between the execution  
 493 time of the given code through running on CPUs  
 494 and the accuracy of using LLMs as surrogate mod-  
 495 els to acquire the output. We categorize problems  
 496 in TC into distinct bins based on their execution  
 497 times and calculate the average accuracy for each  
 498 model across all samples within the same bin, as  
 499 shown in Figure 4. We observed the trend of pre-  
 500 diction accuracy of the model falling as the actual  
 501 execution time required for the corresponding pro-  
 502 gram prolonged. It’s especially worth noting that  
 503 for computational tasks that require execution time  
 504 longer than 1 second, even state-of-the-art models  
 505 struggle to obtain even one correct answer.

## 5.4 The Impact of Variable Type

506 The DR subset in SURGE examines model perfor-  
 507 mance in predicting code behavior under differ-  
 508 ent environmental factors. Specifically, DR in-  
 509 cludes variations in C++ compiler version, C++  
 510 standard, and compilation optimization settings.  
 511 Table 3 presents model performance across differ-  
 512 ent prompting strategies in DR.

513 In the zero-shot setting, model accuracy im-  
 514 proves sequentially across the three factors, sug-  
 515 gesting that compiler version differences are harder  
 516 to predict than standard variations, which in turn  
 517 are harder than optimization settings. However,

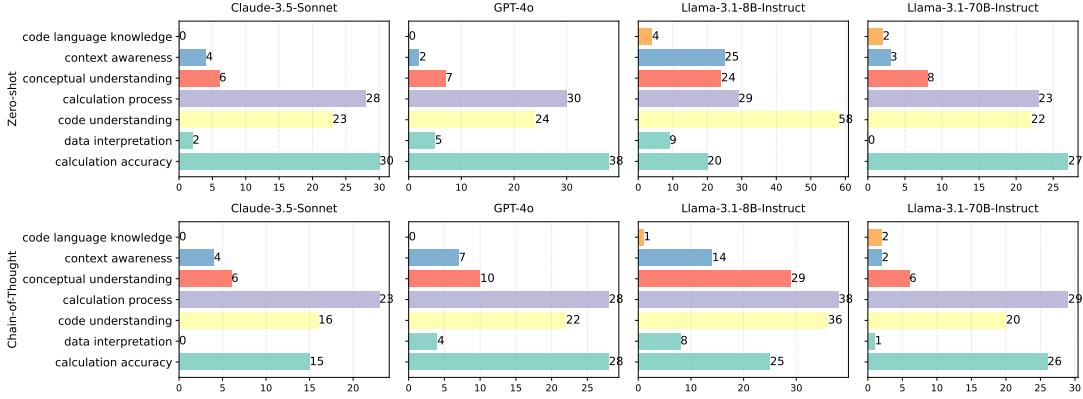


Figure 5: Breakdown of error types across different language models and prompting methods.

with Chain-of-Thought reasoning, performance declines across all factors, with the sharpest drop for optimization settings. This indicates that while CoT aids reasoning for compiler versions and standards, it adds unnecessary complexity for optimizations, ultimately reducing accuracy.

## 5.5 Error Analysis

To further understand the model’s performance and limitations regarding serving as general surrogate models, we categorize the Errors made by Claude-3.5-Sonnet, GPT-4O, and Llama-3.1-8B-Instruct and Llama-3.1-70B-Instruct in the CL subset of SURGE.

We first combined LLM-assisted annotation and manual verification to identify 7 most typical error types: 1. Code Language Knowledge: evaluating foundational programming language proficiency, 2. Context Awareness: measuring understanding of long text and repository-level code, 3. Conceptual Understanding: assessing comprehension of programming concepts, 4. Calculation Process: verifying computational step accuracy, 5. Code Understanding: testing comprehension of code logic and structure, 6. Data Interpretation: evaluating data processing and analysis capabilities, and 7. Calculation Accuracy: Measuring precision in scientific computations. We then use LLM to label errors in the model’s responses. The categorized error statistics are demonstrated in Figure 5.

As models prompted with CoT exhibit fewer instances of most error types compared to zero-shot, especially for the Code Understanding capability of Llama. The main error types are Code Understanding, Calculation Process, and Calculation Accuracy. For zero-shot, the primary error is accuracy, but for CoT, the most frequent error is Calculation Process. This suggests that CoT can better grasp the overall

code logic and produce more correct results, but it may still make mistakes in the chain of thought process. In general, CoT has fewer and smaller errors. From the model perspective, Llama has a clear lead in Conceptual Understanding errors, indicating its weaker ability to understand concepts.

## 5.6 Training Scale Analysis

We also investigate whether training scaling can affect LLMs’ surrogate execution capabilities on the FL task. We trained models of varying sizes on different amounts of training data sampled from a similar distribution.

As demonstrated in Figure 2, both model size and training steps are crucial factors in determining surrogate execution accuracy. As we scale from 0.5B to 7B parameters, models consistently show improved learning efficiency and higher performance ceilings throughout the training process. Larger models learn faster in the early stages and continue to improve for longer before plateauing, suggesting better utilization of the training data.

These empirical observations align with established scaling laws in language modeling, indicating that surrogate execution capabilities follow similar scaling patterns as other language tasks.

## 6 Conclusion

In this paper, we introduce SURGE, a holistic benchmark for evaluating LLMs as general-purpose surrogate code executors. The curated dataset spans multiple domains, bridging theoretical tasks with real-world applications. Through extensive empirical study, we argue that while current large language models possess the power to predict code execution results to a certain extent, there remains significant room for further improvements on grounding LLMs to facilitate general surrogate model.

## 593 Limitations

594 Despite its comprehensive evaluation, our study  
595 has several limitations. LLMs remain approxi-  
596 mators rather than exact code executors, often  
597 struggling with edge cases, intricate runtime be-  
598 haviors, and execution-dependent state changes.  
599 While SURGE covers diverse execution scenarios,  
600 it does not encompass all specialized environments,  
601 such as hardware-dependent simulations or real-  
602 time systems. Additionally, LLMs may generate  
603 plausible but incorrect outputs, particularly in com-  
604 plex logical dependencies or undefined behaviors,  
605 making error detection challenging. Our scaling  
606 study is constrained by computational resources,  
607 limiting the assessment of extremely large models  
608 or extensive training data distributions. Further-  
609 more, security risks remain, as LLMs may fail to  
610 recognize vulnerabilities, potentially misjudging  
611 harmful code. Finally, our benchmark operates in  
612 a controlled setting, whereas real-world software  
613 development involves dynamic interactions and it-  
614 erative debugging, which are not fully captured in  
615 our study. Future work should focus on improving  
616 LLMs' reasoning abilities, enhancing robustness in  
617 execution prediction, and integrating them with tra-  
618 ditional program analysis techniques for practical  
619 deployment.

## 620 References

- 621 Wasi Uddin Ahmad, Saikat Chakraborty, Baishakhi Ray,  
622 and Kai-Wei Chang. 2020. A transformer-based  
623 approach for source code summarization. *arXiv*  
624 preprint arXiv:2005.00653.
- 625 Anthropic. 2024a. Introducing Claude 3.5 Sonnet.
- 626 Anthropic. 2024b. The Claude 3  
627 Model Family: Opus, Sonnet, Haiku.  
628 [https://www-cdn.anthropic.com/  
629 de8ba9b01c9ab7cbabf5c33b80b7bbc618857627/  
630 Model\\_Card\\_Claude\\_3.pdf](https://www-cdn.anthropic.com/de8ba9b01c9ab7cbabf5c33b80b7bbc618857627/Model_Card_Claude_3.pdf).
- 631 Bruno Barras, Samuel Boutin, Cristina Cornes, Ju-  
632 dicaël Courant, Jean-Christophe Filliatre, Eduardo  
633 Gimenez, Hugo Herbelin, Gerard Huet, Cesar Munoz,  
634 Chetan Murthy, et al. 1997. *The Coq proof assistant*  
635 *reference manual: Version 6.1*. Ph.D. thesis, Inria.
- 636 Peter Benner, Serkan Gugercin, and Karen Willcox.  
637 2015. A survey of projection-based model reduction  
638 methods for parametric dynamical systems. *SIAM*  
639 *Review*, 57(4):483–531.
- 640 Miguel A Bessa, Ramin Bostanabad, Zeliang Liu, Anqi  
641 Hu, Daniel W Apley, Catherine Brinson, Wei Chen,

and Wing Kam Liu. 2017. A framework for data-  
driven analysis of materials under uncertainty: Counter-  
642 ing the curse of dimensionality. *Computer Meth-  
643 ods in Applied Mechanics and Engineering*, 320:633–  
644 667.

Xiao Bi, Deli Chen, Guanting Chen, Shanhuang Chen,  
Damaï Dai, Chengqi Deng, Honghui Ding, Kai Dong,  
Qiushi Du, Zhe Fu, et al. 2024. Deepseek llm: Scal-  
645 ing open-source language models with longtermism.  
arXiv preprint arXiv:2401.02954.

Randal E. Bryant and David R. O'Hallaron. 2010. *Com-  
646 puter Systems: A Programmer's Perspective*, 2nd edi-  
647 tion. Addison-Wesley Publishing Company, USA.

Cristian Cadar and Koushik Sen. 2013. *Symbolic exe-  
648 cution for software testing: three decades later*. *Commun. ACM*, 56(2):82–90.

Linzheng Chai, Shukai Liu, Jian Yang, Yuwei Yin,  
Ke Jin, Jiaheng Liu, Tao Sun, Ge Zhang, Changyu  
Ren, Hongcheng Guo, Zekun Wang, Boyang Wang,  
Xianjie Wu, Bing Wang, Tongliang Li, Liqun Yang,  
Sufeng Duan, and Zhoujun Li. 2024. *Mceval: Mas-  
649 sively multilingual code evaluation*.

Saikat Chakraborty, Yangruibo Ding, Miltiadis Allama-  
650 nis, and Baishakhi Ray. 2020. *Codit: Code editing  
651 with tree-based neural models*. *IEEE Transactions  
652 on Software Engineering*, pages 1–1.

Saikat Chakraborty, Rahul Krishna, Yangruibo Ding,  
and Baishakhi Ray. 2020. *Deep learning based  
653 vulnerability detection: Are we there yet?* arXiv  
654 preprint arXiv:2009.07235.

Sahil Chaudhary. 2023. Code alpaca: An instruction-  
following llama model for code generation. [https://  
655 /github.com/sahil280114/codealpaca](https://github.com/sahil280114/codealpaca).

Haoxuan Che, Xuanhua He, Quande Liu, Cheng Jin,  
and Hao Chen. 2024. *Gamegen-x: Interactive open-  
656 world game video generation*.

Zimin Chen, Steve James Kommrusch, Michele Tufano,  
Louis-Noël Pouchet, Denys Poshyvanyk, and Martin  
Monperrus. 2019. *Sequencer: Sequence-to-sequence  
657 learning for end-to-end program repair*. *IEEE Trans-  
658 actions on Software Engineering*.

Leonardo De Moura, Soonho Kong, Jeremy Avigad,  
Floris Van Doorn, and Jakob von Raumer. 2015. The  
lean theorem prover (system description). In *Auto-  
659 mated Deduction-CADE-25: 25th International Con-  
660 ference on Automated Deduction, Berlin, Germany,  
661 August 1-7, 2015, Proceedings 25*, pages 378–388.  
Springer.

Nate Gruver, Marc Finzi, Shikai Qiu, and Andrew Gor-  
don Wilson. 2024. *Large language models are zero-  
662 shot time series forecasters*.

Ruižhen Gu, José Miguel Rojas, and Donghwan Shin.  
2025. *Software testing for extended reality applica-  
663 tions: A systematic mapping study*.

696	Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, Y Wu, YK Li, et al. 2024. Deepseek-coder: When the large language model meets programming—the rise of code intelligence. <i>arXiv preprint arXiv:2401.14196</i> .	751
697		752
698		753
699		
700		
701	Hao Hao, Xiaoqun Zhang, and Aimin Zhou. 2024. Large language models as surrogate models in evolutionary algorithms: A preliminary study.	754
702		755
703		756
704	Jacob Harer, Chris Reale, and Peter Chin. 2019. Tree-transformer: A transformer-based method for correction of tree-structured data. <i>arXiv preprint arXiv:1908.00449</i> .	757
705		758
706		759
707		760
708	Jan Hesthaven and Stefano Ubbiali. 2018. Non-intrusive reduced order modeling of nonlinear problems using neural networks. <i>Journal of Computational Physics</i> , 363.	761
709		762
710		763
711		764
712	Xing Hu, Ge Li, Xin Xia, David Lo, and Zhi Jin. 2018. Deep code comment generation. In <i>Proceedings of the 26th Conference on Program Comprehension</i> , page 200–210, New York, NY, USA. Association for Computing Machinery.	765
713		766
714		767
715		
716	Binyuan Hui, Jian Yang, Zeyu Cui, Jiaxi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Kai Dang, et al. 2024. Qwen2. 5-coder technical report. <i>arXiv preprint arXiv:2409.12186</i> .	
717		
718		
719		
720		
721	Hamel Husain, Ho-Hsiang Wu, Tiferet Gazit, Miltiadis Allamanis, and Marc Brockschmidt. 2020. Code-searchnet challenge: Evaluating the state of semantic code search.	774
722		775
723		776
724		777
725	Paul Jaccard. 1901. Étude comparative de la distribution florale dans une portion des alpes et des jura. <i>Bulletin de la Société Vaudoise des Sciences Naturelles</i> , 37:547–579.	778
726		779
727		
728		
729	James C. King. 1976. Symbolic execution and program testing. <i>Commun. ACM</i> , 19(7):385–394.	780
730		781
731	Jian Li, Yue Wang, Michael R. Lyu, and Irwin King. 2018a. Code completion with neural attention and pointer networks. In <i>Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI-18</i> , pages 4159–4165. International Joint Conferences on Artificial Intelligence Organization.	782
732		783
733		784
734		785
735		786
736		787
737		
738	Junlong Li, Daya Guo, Dejian Yang, Runxin Xu, Yu Wu, and Junxian He. 2025. Code/o: Condensing reasoning patterns via code input-output prediction.	788
739		789
740		790
741	Zhen Li, Deqing Zou, Shouhuai Xu, Xinyu Ou, Hai Jin, Sujuan Wang, Zhijun Deng, and Yuyi Zhong. 2018b. Vuldeepecker: A deep learning-based system for vulnerability detection. <i>arXiv preprint arXiv:1801.01681</i> .	791
742		
743		
744		
745		
746	Yong Lin, Shange Tang, Bohan Lyu, Jiayun Wu, Hongzhou Lin, Kaiyu Yang, Jia Li, Mengzhou Xia, Danqi Chen, Sanjeev Arora, and Chi Jin. 2025. Goedel-prover: A frontier model for open-source automated theorem proving.	801
747		802
748		803
749		
750		
751	Stanley B. Lippman, Jose Lajoie, and Barbara E. Moo. 2012. <i>C++ Primer</i> , 5th edition. Addison-Wesley Professional.	751
752		752
753		753
754	Dan Lu and Daniel Ricciuto. 2019. Efficient surrogate modeling methods for large-scale earth system models based on machine learning techniques.	754
755		755
756		756
757	Shuai Lu, Daya Guo, Shuo Ren, Junjie Huang, Alexey Svyatkovskiy, Ambrosio Blanco, Colin Clement, Dawn Drain, Dixin Jiang, Duyu Tang, Ge Li, Liding Zhou, Linjun Shou, Long Zhou, Michele Tu-fano, MING GONG, Ming Zhou, Nan Duan, Neel Sundaresan, Shao Kun Deng, Shengyu Fu, and Shujie LIU. 2021a. CodeXGLUE: A machine learning benchmark dataset for code understanding and generation. In <i>Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1)</i> .	757
758		758
759		759
760		760
761		761
762		762
763		763
764		764
765		765
766		766
767		767
768	Shuai Lu, Daya Guo, Shuo Ren, Junjie Huang, Alexey Svyatkovskiy, Ambrosio Blanco, Colin Clement, Dawn Drain, Dixin Jiang, Duyu Tang, et al. 2021b. Codexglue: A machine learning benchmark dataset for code understanding and generation. <i>arXiv preprint arXiv:2102.04664</i> .	768
769		769
770		770
771		771
772		772
773		773
774	Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xubo Geng, Wenxiang Hu, Chongyang Tao, Jing Ma, Qingwei Lin, and Dixin Jiang. 2023. Wizardcoder: Empowering code large language models with evolution-instruct. In <i>The Twelfth International Conference on Learning Representations</i> .	774
775		775
776		776
777		777
778		778
779		779
780	Thibaud Lutellier, Hung Viet Pham, Lawrence Pang, Yitong Li, Moshi Wei, and Lin Tan. 2020. Coconut: combining context-aware neural translation models using ensemble for program repair. In <i>Proceedings of the 29th ACM SIGSOFT International Symposium on Software Testing and Analysis</i> , pages 101–114, New York, NY, USA. Association for Computing Machinery.	780
781		781
782		782
783		783
784		784
785		785
786		786
787		787
788	Bohan Lyu, Yadi Cao, Duncan Watson-Parris, Leon Bergen, Taylor Berg-Kirkpatrick, and Rose Yu. 2025. Adapting while learning: Grounding llms for scientific problems with intelligent tool usage adaptation.	788
789		789
790		790
791		791
792	Chenyang Lyu, Lecheng Yan, Rui Xing, Wenxi Li, Younes Samih, Tianbo Ji, and Longyue Wang. 2024. Large language models as code executors: An exploratory study.	792
793		793
794		794
795		795
796	Pingchuan Ma, Tsun-Hsuan Wang, Minghao Guo, Zhiqing Sun, Joshua B. Tenenbaum, Daniela Rus, Chuang Gan, and Wojciech Matusik. 2024. Llm and simulation as bilevel optimizers: A new paradigm to advance physical scientific discovery.	796
797		797
798		798
799		799
800		800
801	Dirk Merkel. 2014. Docker: lightweight linux containers for consistent development and deployment. <i>Linux J.</i> , 2014(239).	801
802		802
803		803
804	Meta. 2024. Introducing Meta Llama 3: The most capable openly available LLM to date. <a href="https://ai.meta.com/blog/meta-llama-3/">https://ai.meta.com/blog/meta-llama-3/</a> .	804
805		805
806		806

807	Leonardo de Moura and Sebastian Ullrich. 2021. The Lean 4 theorem prover and programming language.	861
808		862
809	Giuseppe Nebbione and Maria Carla Calzarossa. 2023. A methodological framework for ai-assisted security assessments of active directory environments. <i>Ieee Access</i> , 11:15119–15130.	863
810		
811		
812		
813	OpenAI. 2024a. <a href="#">Gpt-4o mini: Advancing cost-efficient intelligence</a> . OpenAI Blog. Accessed: 2025-02-16.	867
814		868
815	OpenAI. 2024b. <a href="#">Hello gpt-4o</a> . OpenAI Blog. Accessed: 2025-02-16.	869
816		870
817	Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. <a href="#">Training language models to follow instructions with human feedback</a> .	871
818		872
819		873
820		874
821		
822		
823		
824		
825	Md Rizwan Parvez, Saikat Chakraborty, Baishakhi Ray, and Kai-Wei Chang. 2018. Building language models for text with named entities. <i>arXiv preprint arXiv:1805.04836</i> .	875
826		876
827		877
828		
829	Lawrence C Paulson. 1994. <i>Isabelle: A generic theorem prover</i> .	878
830		879
831	Rui Queiroz, Tiago Cruz, Jérôme Mendes, Pedro Sousa, and Paulo Simões. 2023. <a href="#">Container-based virtualization for real-time industrial systems—a systematic review</a> . <i>ACM Comput. Surv.</i> , 56(3).	880
832		881
833		882
834		
835	Maziar Raissi, Paris Perdikaris, and George E Karniadakis. 2019. Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. <i>Journal of Computational physics</i> , 378:686–707.	883
836		884
837		885
838		886
839		887
840		
841	Maziar Raissi, Alireza Yazdani, and George Em Karniadakis. 2020. Hidden fluid mechanics: Learning velocity and pressure fields from flow visualizations. <i>Science</i> , 367(6481):1026–1030.	888
842		889
843		
844		
845	John W. Ratcliff and David E. Metzener. 1988. Pattern matching: The gestalt approach. <i>Dr. Dobb's Journal</i> , 13(7):46–51.	890
846		891
847		
848	Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittweisser, et al. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. <i>arXiv preprint arXiv:2403.05530</i> .	892
849		893
850		894
851		895
852		
853		
854	Rebecca Russell, Louis Kim, Lei Hamilton, Tomo Lazonich, Jacob Harer, Onur Ozdemir, Paul Ellingwood, and Marc McConley. 2018. <a href="#">Automated vulnerability detection in source code using deep representation learning</a> . In <i>2018 17th IEEE International Conference on Machine Learning and Applications (ICMLA)</i> , pages 757–762. IEEE.	896
855		897
856		898
857		899
858		900
859		
860		
861	Max Schäfer, Sarah Nadi, Aryaz Eghbali, and Frank Tip. 2023. <a href="#">An empirical evaluation of using large language models for automated unit test generation</a> .	901
862		902
863		903
864	Farid Shirazi, Adnan Seddighi, and Amna Iqbal. 2017. <a href="#">Cloud computing security and privacy: An empirical study</a> . pages 534–549.	904
865		911
866		912
867	Mohammed Latif Siddiq, Joanna Cecilia Da Silva Santos, Ridwanul Hasan Tanvir, Noshin Ulfat, Fahmid Al Rifat, and Vinícius Carvalho Lopes. 2024. <a href="#">Using large language models to generate junit tests: An empirical study</a> . In <i>Proceedings of the 28th International Conference on Evaluation and Assessment in Software Engineering</i> , EASE 2024, page 313–322. ACM.	913
868		
869		
870		
871		
872		
873		
874		
875	C. Spearman. 1904. <a href="#">The proof and measurement of association between two things</a> . <i>The American Journal of Psychology</i> , 15(1):72–101.	914
876		915
877		916
878	Gang Sun and Shuyue Wang. 2019. A review of the artificial neural network surrogate modeling in aerodynamic design. <i>Proceedings of the Institution of Mechanical Engineers, Part G: Journal of Aerospace Engineering</i> , 233(16):5863–5872.	917
879		918
880		919
881		920
882		
883	Luning Sun, Han Gao, Shaowu Pan, and Jian-Xun Wang. 2020. Surrogate modeling for fluid flows based on physics-constrained deep learning without simulation data. <i>Computer Methods in Applied Mechanics and Engineering</i> , 361:112732.	921
884		922
885		923
886		924
887		925
888	Qwen Team. 2024a. Code with codeqwen1.5. <a href="https://qwenlm.github.io/blog/codeqwen1.5">https://qwenlm.github.io/blog/codeqwen1.5</a> .	926
889		927
890	Qwen Team. 2024b. Qwen2.5 technical report. <i>arXiv preprint arXiv:2412.15115</i> .	928
891		929
892	Nils Thuerey, Konstantin Weißenow, Lukas Prantl, and Xiangyu Hu. 2020. Deep learning methods for reynolds-averaged navier–stokes simulations of airfoil flows. <i>AIAA Journal</i> , 58(1):25–36.	930
893		931
894		932
895		933
896	Runchu Tian, Yining Ye, Yujia Qin, Xin Cong, Yankai Lin, Yinxu Pan, Yesai Wu, Haotian Hui, Weichuan Liu, Zhiyuan Liu, and Maosong Sun. 2024. <a href="#">Debugbench: Evaluating debugging capability of large language models</a> .	934
897		935
898		936
899		937
900		938
901	Shang Wang, Tianqing Zhu, Bo Liu, Ming Ding, Xu Guo, Dayong Ye, Wanlei Zhou, and Philip S. Yu. 2024. <a href="#">Unique security and privacy threats of large language model: A comprehensive survey</a> .	939
902		940
903		941
904		942
905	Wenhan Wang, Ge Li, Bo Ma, Xin Xia, and Zhi Jin. 2020. <a href="#">Detecting code clones with graph neural network and flow-augmented abstract syntax tree</a> . In <i>2020 IEEE 27th International Conference on Software Analysis, Evolution and Reengineering (SANER)</i> , pages 261–271. IEEE.	943
906		944
907		945
908		946
909		947
910		948
911	Cody Watson, Michele Tufano, Kevin Moran, Gabriele Bavota, and Denys Poshyvanyk. 2020. <a href="#">On learning meaningful assert statements for unit test cases</a> .	949
912		950
913		951

914	In <i>Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering, ICSE '20</i> . ACM.	970
915		971
916		972
917	Logan Weber, Jesse Michel, Alex Renda, and Michael Carbin. 2024. Learning to compile programs to neural networks.	973
918		974
919		975
920	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in neural information processing systems</i> , 35:24824–24837.	976
921		977
922		978
923		979
924		980
925	Jared Willard, Xiaowei Jia, Shaoming Xu, Michael Steinbach, and Vipin Kumar. 2022. Integrating scientific knowledge with machine learning for engineering and environmental systems. <i>ACM Computing Surveys</i> , 55(4):1–37.	981
926		982
927		983
928		984
929		985
930	Christopher Wimmer and Navid Rekabsaz. 2023. Leveraging vision-language models for granular market change prediction.	986
931		987
932		988
933	Haixu Wu, Hang Zhou, Mingsheng Long, and Jianmin Wang. 2023. Interpretable weather forecasting for worldwide stations with a unified deep model. <i>Nat. Mac. Intell.</i> , 5(6):602–611.	989
934		990
935		991
936		992
937	Yujun Yan, Kevin Swersky, Danai Koutra, Parthasarathy Ranganathan, and Milad Hashemi. 2020. Neural execution engines: Learning to execute subroutines.	993
938		994
939		995
940	Huaiyuan Ying, Zijian Wu, Yihan Geng, Jiayu Wang, Dahua Lin, and Kai Chen. 2024. Lean workbook: A large-scale lean problem set formalized from natural language math problems. <i>arXiv preprint arXiv:2406.03847</i> .	996
941		997
942		998
943		999
944		1000
945	Hao Yu, Wing Lam, Long Chen, Ge Li, Tao Xie, and Qianxiang Wang. 2019. Neural detection of semantic code clones via tree-based convolution. In <i>2019 IEEE/ACM 27th International Conference on Program Comprehension (ICPC)</i> , pages 70–80. IEEE Press.	
946		
947		
948		
949		
950		
951	Jian Zhang, Xu Wang, Hongyu Zhang, Hailong Sun, Kaixuan Wang, and Xudong Liu. 2019. A novel neural source code representation based on abstract syntax tree. In <i>Proceedings of the 41st International Conference on Software Engineering, ICSE '19</i> , page 783–794. IEEE Press.	
952		
953		
954		
955		
956		
957	Xuan Zhang, Limei Wang, Jacob Helwig, Youzhi Luo, Cong Fu, Yaochen Xie, Meng Liu, Yuchao Lin, Zhao Xu, Keqiang Yan, Keir Adams, Maurice Weiler, Xiner Li, Tianfan Fu, Yucheng Wang, Haiyang Yu, YuQing Xie, Xiang Fu, Alex Strasser, Shenglong Xu, Yi Liu, Yuanqi Du, Alexandra Saxton, Hongyi Ling, Hannah Lawrence, Hannes Stärk, Shurui Gui, Carl Edwards, Nicholas Gao, Adriana Ladera, Tailin Wu, Elyssa F. Hofgard, Aria Mansouri Tehrani, Rui Wang, Ameya Daigavane, Montgomery Bohde, Jerry Kurtin, Qian Huang, Tuong Phung, Minkai Xu, Chaitanya K. Joshi, Simon V. Mathis, Kamyar Azizzadehesheli, Ada Fang, Alán Aspuru-Guzik, Erik Bekkers,	
958		
959		
960		
961		
962		
963		
964		
965		
966		
967		
968		
969		
970	Michael Bronstein, Marinka Zitnik, Anima Anandkumar, Stefano Ermon, Pietro Liò, Rose Yu, Stephan Günnemann, Jure Leskovec, Heng Ji, Jimeng Sun, Regina Barzilay, Tommi Jaakkola, Connor W. Coley, Xiaoning Qian, Xiaofeng Qian, Tess Smidt, and Shuiwang Ji. 2024. Artificial intelligence for science in quantum, atomistic, and continuum systems.	
971		
972		
973		
974		
975		
976		
977	Yuchen Zhang, Mingsheng Long, Kaiyuan Chen, Lanxiang Xing, Ronghua Jin, Michael I. Jordan, and Jianmin Wang. 2023. Skilful nowcasting of extreme precipitation with nowcastnet. <i>Nat.</i> , 619(7970):526–532.	
978		
979		
980		
981		
982	Qinkai Zheng, Xiao Xia, Xu Zou, Yuxiao Dong, Shan Wang, Yufei Xue, Zihan Wang, Lei Shen, Andi Wang, Yang Li, et al. 2023. Codegeex: A pre-trained model for code generation with multilingual evaluations on humaneval-x. <i>arXiv preprint arXiv:2303.17568</i> .	
983		
984		
985		
986		
987	Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyan Luo, Zhangchi Feng, and Yongqiang Ma. 2024. Llamafactory: Unified efficient fine-tuning of 100+ language models. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)</i> , Bangkok, Thailand. Association for Computational Linguistics.	
988		
989		
990		
991		
992		
993		
994		
995	Yaqin Zhou, Shangqing Liu, Jingkai Siow, Xiaoning Du, and Yang Liu. 2019. Devign: Effective vulnerability identification by learning comprehensive program semantics via graph neural networks. In <i>Advances in Neural Information Processing Systems</i> , volume 32, pages 10197–10207. Curran Associates, Inc.	
996		
997		
998		
999		

1001 **Appendix**

1002 **A Complete Main Results**

Table 5: Performance of different models under different prompting strategies on SURGE.

Model	ML						CL	RL	SC	TC	BG					
	CPP	Rust	Python	Julia	Java	Others					CPP	Java	Python	DR	FL	
Zero-shot																
Claude-3.5-Sonnet	72.73	55.00	88.00	66.67	76.92	75.00	81.58	57.31	61.55	35.27	9.09	12.55	51.40	12.92	17.92	51.59
DeepSeek-V3	54.55	60.00	76.00	61.11	46.15	72.50	56.58	44.19	59.65	35.31	4.05	3.03	21.50	10.92	32.46	42.53
GPT-4o	40.91	45.00	60.00	55.56	57.69	55.00	66.45	49.13	53.16	34.44	4.28	7.48	34.59	14.75	21.99	40.03
GPT-4o-Mini	59.09	45.00	64.00	33.33	69.23	60.00	84.21	35.33	40.00	31.85	1.39	4.28	13.86	11.65	39.07	39.49
Qwen-Max	50.00	45.00	44.00	27.78	26.92	50.00	38.82	37.15	56.98	35.44	2.82	3.20	30.52	14.03	29.89	32.84
Qwen-2.5-0.5B-Instruct	13.64	10.00	4.00	0.00	11.54	7.50	19.08	5.02	9.43	8.61	1.20	3.77	11.04	4.62	42.38	10.85
Qwen-2.5-1.5B-Instruct	36.36	15.00	24.00	16.67	15.38	35.00	40.79	5.89	12.08	14.15	1.15	3.52	11.10	9.24	41.72	18.80
Qwen-2.5-3B-Instruct	36.36	10.00	28.00	22.22	34.62	22.50	35.53	13.08	20.04	14.61	1.58	3.92	13.27	5.75	15.89	18.49
Qwen-2.5-7B-Instruct	13.64	5.00	12.00	5.56	11.54	17.50	27.63	27.98	22.90	29.43	2.21	3.66	9.46	7.24	36.42	15.48
Qwen-2.5-14B-Instruct	9.09	5.00	8.00	5.56	23.08	5.00	11.18	39.08	33.79	28.18	2.11	2.24	16.30	4.24	7.28	13.34
Qwen-2.5-32B-Instruct	45.45	20.00	48.00	33.33	42.31	20.00	50.66	22.61	32.50	30.99	1.21	2.09	12.16	7.26	7.65	25.08
Qwen-2.5-72B-Instruct	45.45	40.00	52.00	55.56	69.23	52.50	72.37	33.66	53.70	32.46	1.99	2.65	10.69	13.35	34.44	38.00
Qwen-2.5-Coder-0.5B-Instruct	22.73	10.00	20.00	5.56	11.54	22.50	26.97	8.41	6.55	6.39	0.73	2.92	10.05	5.76	41.06	13.41
Qwen-2.5-Coder-1.5B-Instruct	45.45	30.00	32.00	27.78	26.92	35.00	54.61	17.58	20.58	13.19	0.36	2.83	5.43	8.08	41.06	24.06
Qwen-2.5-Coder-3B-Instruct	45.45	20.00	40.00	22.22	46.15	40.00	68.42	21.58	37.43	22.43	1.06	3.61	13.41	8.96	41.06	28.79
Qwen-2.5-Coder-7B-Instruct	45.45	30.00	52.00	44.44	50.00	42.50	75.00	20.86	45.50	28.00	0.86	3.47	13.53	14.23	40.40	33.75
Qwen-2.5-Coder-14B-Instruct	59.09	35.00	48.00	55.56	69.23	62.50	82.89	35.26	51.93	32.20	1.39	3.66	16.12	11.35	41.72	40.39
Qwen-2.5-Coder-32B-Instruct	59.09	55.00	60.00	44.44	69.23	60.00	80.26	48.43	57.18	18.84	1.49	1.95	17.96	13.42	29.19	41.10
LLAMA-3.1-8B-Instruct	0.00	0.00	4.00	5.56	15.38	5.00	13.16	20.41	4.41	15.46	4.41	4.12	5.98	3.97	0.00	8.49
LLAMA-3.1-70B-Instruct	54.55	40.00	52.00	33.33	65.38	47.50	78.29	31.60	48.65	30.19	1.59	4.18	14.27	15.73	39.16	37.10
LLAMA-3.3-70B-Instruct	68.18	50.00	60.00	44.44	57.69	62.50	66.45	43.53	38.45	30.35	2.06	3.23	11.01	11.22	39.13	39.22
Zero-shot Chain-of-Thought																
Claude-3.5-Sonnet	90.91	65.00	96.00	77.78	69.23	92.50	82.24	62.31	63.38	40.70	16.91	20.69	62.23	18.19	33.98	59.47
DeepSeek-V3	81.82	85.00	88.00	72.22	69.23	85.00	76.32	62.70	57.57	36.71	4.45	7.85	46.26	16.21	35.19	54.97
GPT-4o	68.18	65.00	92.00	72.22	76.92	77.50	79.61	53.74	48.56	28.36	8.19	9.97	44.29	14.21	27.91	51.11
GPT-4o-Mini	77.27	60.00	88.00	50.00	69.23	80.00	75.66	34.46	40.60	29.59	1.77	5.40	21.04	13.16	33.11	45.29
Qwen-Max	86.36	75.00	80.00	72.22	76.92	80.00	71.05	50.49	61.78	36.71	2.65	7.73	46.85	16.16	20.74	52.31
Qwen-2.5-0.5B-Instruct	27.27	10.00	16.00	0.00	3.85	17.50	22.37	6.29	1.11	5.28	1.22	3.41	9.15	5.21	37.09	11.84
Qwen-2.5-1.5B-Instruct	31.82	15.00	20.00	11.11	19.23	27.50	26.32	12.47	8.14	12.77	1.22	3.23	9.81	4.93	39.74	16.22
Qwen-2.5-3B-Instruct	40.91	40.00	32.00	22.22	30.77	35.00	25.00	15.01	20.78	12.84	2.14	2.86	5.29	7.12	13.93	20.39
Qwen-2.5-7B-Instruct	40.91	15.00	32.00	33.33	26.92	47.50	52.63	27.51	25.68	28.95	1.12	3.76	14.94	9.41	36.46	26.41
Qwen-2.5-14B-Instruct	68.18	55.00	76.00	61.11	65.38	72.50	61.18	43.89	36.49	32.07	1.57	3.15	19.13	11.97	8.67	41.09
Qwen-2.5-32B-Instruct	50.00	40.00	40.00	55.56	38.46	57.50	53.29	32.71	43.29	30.59	3.10	6.96	23.86	11.49	7.39	32.95
Qwen-2.5-72B-Instruct	68.18	80.00	96.00	55.56	76.92	77.50	65.13	45.30	53.39	32.95	1.72	3.19	16.32	16.61	29.80	47.90
Qwen-2.5-Coder-0.5B-Instruct	13.64	0.00	4.00	5.56	3.85	7.50	1.97	3.07	2.14	3.98	2.14	2.44	2.83	2.23	33.77	6.36
Qwen-2.5-Coder-1.5B-Instruct	31.82	10.00	28.00	5.56	15.38	22.50	19.08	26.39	20.08	15.42	1.16	3.60	11.46	7.24	39.74	17.16
Qwen-2.5-Coder-3B-Instruct	50.00	15.00	32.00	22.22	42.31	40.00	52.63	13.13	33.96	20.14	1.24	3.72	13.36	7.63	40.40	25.85
Qwen-2.5-Coder-7B-Instruct	68.18	40.00	40.00	38.89	53.85	57.50	46.05	19.70	40.91	30.19	2.29	4.71	12.77	15.04	37.12	33.81
Qwen-2.5-Coder-14B-Instruct	59.09	45.00	60.00	55.56	69.23	62.50	71.71	26.09	52.12	33.09	3.19	5.21	18.58	13.41	34.44	40.61
Qwen-2.5-Coder-32B-Instruct	77.27	65.00	80.00	55.56	73.08	67.50	71.71	54.58	55.69	34.36	2.05	4.74	22.43	17.62	28.55	47.34
LLAMA-3.1-8B-Instruct	40.91	15.00	24.00	22.22	26.92	30.00	41.45	17.87	32.65	18.22	1.52	4.23	13.38	10.12	0.66	19.94
LLAMA-3.1-70B-Instruct	59.09	50.00	72.00	61.11	57.69	52.50	58.55	34.44	43.93	29.76	1.71	3.49	15.02	16.86	25.85	38.80
LLAMA-3.3-70B-Instruct	63.64	45.00	56.00	55.56	57.69	67.50	57.24	35.88	39.50	30.61	3.34	3.95	13.53	16.38	35.11	38.73
Few-shot Chain-of-Thought																
Claude-3.5-Sonnet	86.36	70.00	96.00	72.22	65.38	82.50	82.24	70.65	63.58	41.00	22.04	23.61	44.15	25.70	31.99	58.49
DeepSeek-V3	90.91	65.00	84.00	77.78	73.08	95.00	80.26	78.64	66.00	38.60	21.98	15.14	40.27	24.38	35.17	59.08
GPT-4o	68.18	60.00	88.00	77.78	73.08	75.00	75.66	76.86	59.65	37.12	12.91	7.74	29.52	22.08	26.65	52.68
GPT-4o-Mini	77.27	55.00	80.00	50.00	73.08	72.50	71.05	63.89	55.87	34.20	17.89	9.95	23.75	18.24	24.68	48.49
Qwen-Max	81.82	70.00	88.00	77.78	73.08	80.00	82.24	72.53	62.32	37.88	19.68	19.78	37.57	23.91	24.76	56.76
Qwen-2.5-0.5B-Instruct	18.18	5.00	4.00	0.00	7.69	10.00	17.11	17.82	32.32	6.20	3.51	4.83	5.17	5.29	19.21	11.17
Qwen-2.5-1.5B-Instruct	27.27	15.00	16.00	16.67	11.54	27.50	19.74	9.88	31.44	13.20	3.76	3.66	8.15	7.17	41.72	16.85
Qwen-2.5-3B-Instruct	45.45	30.00	28.00	11.11	11.54	27.50	27.63	17.40	38.70	17.74	7.67	6.45	9.09	10.21	35.25	21.58
Qwen-2.5-7B-Instruct	27.27	25.00	36.00	38.89	26.92	42.50	48.68	45.19	43.42	28.97	4.92	4.70	12.94	10.66	34.53	28.71
Qwen-2.5-14B-Instruct	63.64	55.00	60.00	66.67	61.54	70.00	57.24	53.59	49.99	32.48	3.29	4.40	17.88	14.10	10.76	41.37
Qwen-2.5-32B-Instruct	59.09	55.00	52.00	66.67	53.85	60.00	63.16	64.10	63.12	32.53	5.41	6.94	28.66	13.81	22.87	43.15
Qwen-2.5-72B-Instruct	63.64	75.00	88.00	72.22	69.23	80.00	70.39	67.98	61.45	34.92	2.89	3.05	9.52	18.43	33.18	49.99
Qwen-2.5-Coder-0.5B-Instruct	9.09	0.00	0.00	0.00	2.50	1.32	13.88	13.18	5.62	6.66	4.42	2.15	5.19	37.75	9.25	
Qwen-2.5-Coder-1.5B-Instruct	4.55	5.00	16.00	0.00	15.38	5.00	1.97	28.79	38.97	15.25	1.24	3.28	9.70	9.90	33.77	13.49
Qwen-2.5-Coder-3B-Instruct	36.36	45.00	32.00	27.78	30.77	37.50	50.00	27.43	40.84	23.31	1.28	4.06	12.00	9.64		

1003  
1004  
1005

## B Prompts

### B.1 Prompts for Dataset Refactoring

ML:

I will provide you with a code problem with a solution. You need to generate a complete, executable code based on the raw json data, including all necessary package imports, the original code, the test cases, and the main function.  
You need to generate the executable code and expected result.  
Please choose a test case according to the 'test' field from raw json data, and the code should print the answer of the test case.  
The output should be json format, with code and expected\_result fields.  
Please only generate the number or string answer in 'expected\_result' field without any extra description.

1006  
1007

CL:

I will provide you with the solution to a code problem in cpp, python, and javascript. You need to score according to the difficulty of the problem from 1 to 5, while 5 means the hardest. And generate topic keywords for the problem.  
The output should only be json format, with difficulty and keywords fields.  
difficulty: 1-5, integer  
keywords: two or three words to best describe the problem, string list

1008  
1009

BG:

I will provide you with a piece of code and some test cases. You need to generate a complete, executable code based on these, including all necessary package imports, the original code, the test cases, and the main function. You should wrap the original code with ORIGINAL\_CODE\_START and ORIGINAL\_CODE\_END comments. Additionally, the program should output the results of the test cases. Do not include expected output in your answer.

1010  
1011  
1012

## C Details of SURGE

### C.1 ML

Table 6: Language usage count across different categories in the ML subset.

Java	C#	Rust	Julia	Python	C++	C
25	20	20	26	18	21	20

1013  
1014  
1015  
1016

In ML, the usage distribution of various programming languages is shown in Table 6. We selected a variety of languages, including Java, C#, Rust, Julia, Python, C++, and C, to evaluate the model's

ability to handle multilingual code. This diverse selection helps to comprehensively assess the model's performance across different languages.

1017  
1018  
1019

### C.1.1 System Prompts

#### Zero-shot Chain-of-Thought:

Given the following code, what is the execution result?  
You should think step by step. Your answer should be in the following format:  
Thought: <your thought>  
Output:  
<execution result>

1020  
1021

#### Zero-shot:

Given the following code, what is the execution result?  
Your answer should be in the following format:  
Output:  
<execution result>

1022  
1023

#### Few-shot Chain-of-Thought:

Given the following code, what is the execution result?  
You should think step by step. Your answer should be in the following format:  
Thought: <your thought>  
Output:  
<execution result>  
Following are 3 examples:  
<examples here>

1024  
1025

### C.1.2 Demo Questions

```
def catalan_number(n: int) -> int:
    # Initialize an array to store the
    # intermediate catalan numbers
    catalan = [0] * (n + 1)
    catalan[0] = 1 # Base case

    # Calculate catalan numbers using the
    # recursive formula
    for i in range(1, n + 1):
        for j in range(i):
            catalan[i] += catalan[j] * catalan[i - j - 1]

    return catalan[n]

if __name__ == "__main__":
    # Run the test function and print the
    # result of a specific test case
    print(catalan_number(3))
```

1026  
1027

```
import java.util.*;

class Solution {
    public static int countPrefixWords(List<
        String> wordList, String prefix) {

        int count = 0;
        for (String word : wordList) {
```

1028  
1029

```

        if (word.startsWith(prefix)) {
            count++;
        }
    }
    return count;
}

public static void main(String[] args) {
    System.out.println(countPrefixWords(
        Arrays.asList("dog", "dodge", "dot", "dough"), "do"));
}
}

```

1030

```

#include <assert.h>
#include <stdio.h>

long long minTotalCost(int n, int *C)
{
    return (long long)(C[n-2]) * (n - 1) + C[n-1];
}

int main() {
    int costs3[] = {5, 4, 3, 2};
    printf("%lld\n", minTotalCost(4, costs3));
    return 0;
}

```

1031

```

function merge_sorted_arrays(nums1::Vector{Int},
    m::Int, nums2::Vector{Int}, n::Int) :: Vector{Int}
    i = m
    j = n
    k = m + n

    while j > 0
        if i > 0 && nums1[i] > nums2[j]
            nums1[k] = nums1[i]
            i -= 1
        else
            nums1[k] = nums2[j]
            j -= 1
        end
        k -= 1
    end

    nums1
end

# Test case
result = merge_sorted_arrays([1, 3, 5, 0, 0,
    0], 3, [2, 4, 6], 3)
println(result)

```

1032

```

public class Solution {

    public static int findSmallestInteger(int n)
    {
        char[] characters = Integer.toString(n).
            toCharArray();
        int i = characters.length - 2;

        // Find the first digit that is smaller
        // than the digit next to it.
    }
}

```

1033

```

while (i >= 0 && characters[i] >=
    characters[i + 1]) {
    i--;
}

if (i == -1) {
    return -1; // Digits are in descending
    order, no greater number possible.
}

// Find the smallest digit on right side of
// (i) which is greater than characters[i]
int j = characters.length - 1;
while (characters[j] <= characters[i]) {
    j--;
}

// Swap the digits at indices i and j
swap(characters, i, j);

// Reverse the digits from index i+1 to the
// end of the array
reverse(characters, i + 1);

try {
    return Integer.parseInt(new String(
        characters));
} catch (NumberFormatException e) {
    return -1; // The number formed is beyond
    the range of int.
}

private static void swap(char[] arr, int i,
    int j) {
    char temp = arr[i];
    arr[i] = arr[j];
    arr[j] = temp;
}

private static void reverse(char[] arr, int
    start) {
    int end = arr.length - 1;
    while (start < end) {
        swap(arr, start, end);
        start++;
        end--;
    }
}

public static void main(String[] args) {
    System.out.println(findSmallestInteger(123))
    ;
}
}

```

1034

## C.2 CL

1035

Table 7: Language usage count across different categories in the CL subset.

Python	C++	JavaScript
50	51	49

In CL, we selected competition problems of vary-

1036

Table 8: Details of problems in different languages and different difficulty levels.

<b>Difficulty</b>	JavaScript	CPP	Python
1	10	11	11
2	6	4	6
3	12	14	12
4	8	8	9
5	13	14	12

ing difficulty, each with solutions in Python, C++, and JavaScript. You can see the distribution of language in Table 7, and the distribution of problem difficulty in Table 8. This selection allows us to test the model’s cross-language capabilities and its ability to handle problems of different difficulty levels.

### C.2.1 System Prompts

## Zero-shot Chain-of-Thought:

Given the following code, what is the execution result?  
You should think step by step. Your answer should be in the following format:  
Thought: <your thought>  
Output:  
<execution result>

## Zero-shot:

Given the following code, what is the execution result?  
Your answer should be in the following format:  
Output:  
<execution result>

## Few-shot Chain-of-Thought:

Given the following code, what is the execution result?  
You should think step by step. Your answer should be in the following format:  
Thought: <your thought>  
Output:  
<execution result>  
Following are 3 examples:  
{examples here}}

### C.2.2 Demo Questions

```
class TreeNode {  
    constructor(val) {  
        this.val = val;  
        this.left = this.right = null;  
    }  
}  
  
function maxDepth(root) {  
    if (!root) return 0;
```

```
const queue = [root, null];
let depth = 1;

while (queue.length > 0) {
  const node = queue.shift();
  if (node === null) {
    if (queue.length === 0) return depth;
    depth++;
    queue.push(null);
    continue;
  }
  if (node.left) queue.push(node.left);
  if (node.right) queue.push(node.right);
}

return depth;
}

// Test case
const root = new TreeNode(3);
root.left = new TreeNode(9);
root.right = new TreeNode(20);
root.right.left = new TreeNode(15);
root.right.right = new TreeNode(7);
console.log(maxDepth(root));
```

```
#include <iostream>
#include <unordered_map>
#include <string>
using namespace std;

int findTheLongestSubstring(string s) {
    unordered_map<char, int> mapper = {{'a', 1},
                                         {'e', 2}, {'i', 4}, {'o', 8}, {'u', 16}};
    unordered_map<int, int> seen;
```

```

seen[0] = -1;
int max_len = 0, cur = 0;

for(int i = 0; i < s.size(); ++i){
    if(mapper.find(s[i]) != mapper.end()){
        cur ^= mapper[s[i]];
    }
    if(seen.find(cur) != seen.end()){
        max_len = max(max_len, i - seen[cur]);
    } else {
        seen[cur] = i;
    }
}

return max_len;
}

// Test case
class Solution {
public:
    void solve() {
        string input = "leetminicoworoep";
        cout << findTheLongestSubstring(input)
        << endl; // Expected output: 13
    }
};

int main(){
    Solution sol;
    sol.solve();
    return 0;
}

```

1056

```

class TreeNode {
    constructor(val) {
        this.val = val;
        this.left = this.right = null;
    }
}

function backtrack(root, sum, res, tempList) {
    if (root === null) return;
    if (root.left === null && root.right === null && sum === root.val)
        return res.push([...tempList, root.val]);

    tempList.push(root.val);
    backtrack(root.left, sum - root.val, res, tempList);
    backtrack(root.right, sum - root.val, res, tempList);
    tempList.pop();
}

function pathSum(root, sum) {
    if (root === null) return [];
    const res = [];
    backtrack(root, sum, res, []);
    return res;
}

// Test case setup
const root = new TreeNode(5);
root.left = new TreeNode(4);
root.right = new TreeNode(8);
root.left.left = new TreeNode(11);

```

1057

```

root.right.left = new TreeNode(13);
root.right.right = new TreeNode(4);
root.left.left.left = new TreeNode(7);
root.left.left.right = new TreeNode(2);
root.right.right.left = new TreeNode(5);
root.right.right.right = new TreeNode(1);

console.log(pathSum(root, 22));

```

1058

```

class TrieNode {
    constructor() {
        this.children = {};
        this.isEndOfWord = false;
    }
}

class Trie {
    constructor() {
        this.root = new TrieNode();
    }

    insert(word) {
        let node = this.root;
        for (let char of word) {
            if (!node.children[char]) {
                node.children[char] = new TrieNode();
            }
            node = node.children[char];
        }
        node.isEndOfWord = true;
    }

    search(stream) {
        let node = this.root;
        for (let char of stream) {
            if (!node.children[char]) {
                return false;
            }
            node = node.children[char];
        }
        if (node.isEndOfWord) {
            return true;
        }
        return false;
    }
}

class StreamChecker {
    constructor(words) {
        this.trie = new Trie();
        this.stream = [];

        for (let word of [...new Set(words)]) {
            this.trie.insert(word.split('').reverse().join(''));
        }
    }

    query(letter) {
        this.stream.unshift(letter);
        return this.trie.search(this.stream);
    }
}

// Test case
const streamChecker = new StreamChecker(["cd", "f", "kl"]);

```

1059

```

1060 console.log(streamChecker.query('a')); // false
1061 console.log(streamChecker.query('b')); // false
1062 console.log(streamChecker.query('c')); // false
1063 console.log(streamChecker.query('d')); // true
1064 console.log(streamChecker.query('e')); // false
1065 console.log(streamChecker.query('f')); // true
1066 console.log(streamChecker.query('g')); // false
1067 console.log(streamChecker.query('h')); // false
1068 console.log(streamChecker.query('i')); // false
1069 console.log(streamChecker.query('j')); // false
1070 console.log(streamChecker.query('k')); // false
1071 console.log(streamChecker.query('l')); // true

```

### C.3 RL

Table 9: Language usage count across different categories in the RL subset.

Python	C++
24	36

In RL, the distribution of programming language usage is shown in Table 9. We utilized five GitHub repositories for this study, consisting of two Python projects and three C++ projects. Each repository contains a set of ten or more test cases, providing a diverse set of data for evaluation across different programming languages.

#### C.3.1 System Prompts

##### Zero-shot Chain-of-Thought:

You will be given a github repository and a function that generates a latex file with this repo. Your task is to predict the content of the latex file generated by the function. You should think step by step. Your answer should be in the following format:  
Thought: <your thought>  
Output:  
<file content>

##### Zero-shot:

You will be given a github repository and a function that generates a latex file with this repo. Your task is to predict the content of the latex file generated by the function. Your answer should be in the following format:  
Output:  
<file content>

##### Few-shot Chain-of-Thought:

You will be given a github repository and a function that generates a latex file with this repo. Your task is to predict the content of the latex file generated by the function. You should think step by step. Your answer should be in the following format:  
Thought: <your thought>  
Output:

```

<file content>
Following is one example:
{{examples here}}

```

#### C.3.2 Demo Questions

```

main.cpp:<start_file>#include <iostream>
#include <vector>
#include <utility>
#include <stdlib.h>
#include <time.h>
#include <unistd.h>
#include <fstream>
#include <map>
using namespace std;

typedef vector<vector<char> > Board;

const int N = 9;

class SudokuPlayer
{
private:
    //
    int rowUsed[N];
    int columnUsed[N];
    int blockUsed[N];

public:
    vector<Board> result;
    vector<pair<int, int> > spaces;

public:
    SudokuPlayer()
    {
        initState();
    }

    void initState()
    {
        memset(rowUsed, 0, sizeof(rowUsed));
        memset(columnUsed, 0, sizeof(columnUsed));
        memset(blockUsed, 0, sizeof(blockUsed));
        spaces.clear();
        result.clear();
    }

    void addResult(Board &board)
    {
        vector<vector<char> > obj(board);
        result.push_back(obj);
    }

    void flip(int i, int j, int digit)
    {
        rowUsed[i] ^= (1 << digit);
        columnUsed[j] ^= (1 << digit);
        blockUsed[(i / 3) * 3 + j / 3] ^= (1 << digit);
    }

    vector<Board> solveSudoku(Board board)
    {
        initState();
        for (int i = 0; i < N; i++)
        {
            for (int j = 0; j < N; j++)
            {

```

1076

1077

1078

```

        if (board[i][j] == '$')
        {
            spaces.push_back(pair<int,
                int>(i, j));
        }
        else
        {
            int digit = board[i][j] -
                '1';
            flip(i, j, digit);
        }
    }
    DFS(board, 0);
    return result;
}

void DFS(Board &board, int pos)
{
    if (pos == spaces.size())
    {
        addResult(board);
        return;
    }
    int i = spaces[pos].first;
    int j = spaces[pos].second;
    int mask = ~(rowUsed[i] | columnUsed[j]
    | blockUsed[(i / 3) * 3 + j / 3]) & 0
    x1ff;
    int digit = 0;
    while (mask)
    {
        if (mask & 1)
        {
            flip(i, j, digit);
            board[i][j] = '1' + digit;
            DFS(board, pos + 1);
            flip(i, j, digit);
        }
        mask = mask >> 1;
        digit++;
    }
}

void getResult()
{
    for (size_t i = 0; i < result.size(); i++)
    {
        Board board = result[i];
        printBoard(board);
    }
}

bool checkBoard(Board &board)
{
    initState();
    for (int i = 0; i < 9; i++)
    {
        for (int j = 0; j < 9; j++)
        {
            if (board[i][j] != '$')
            {
                int digit = board[i][j] -
                    '1';
                if ((rowUsed[i] | columnUsed[
                    j] | blockUsed[(i / 3) * 3 +
                    j / 3]) & (1 << digit))
                {

```

```

                    return false;
                }
                flip(i, j, digit);
            }
        }
        return true;
    }

    void printBoard(Board &board)
    {
        for (int i = 0; i < board.size(); i++)
        {
            for (int j = 0; j < board[i].size(); j++)
            {
                cout << board[i][j] << " ";
            }
            cout << "\n";
        }
    }

    Board generateBoard(int digCount)
    {
        vector<vector<char>> board(N, vector<
            char>(N, '$'));
        vector<int> row = getRand9();
        for (int i = 0; i < 3; i++)
        {
            board[3][i + 3] = row[i] + '1';
            board[4][i + 3] = row[i + 3] + '1';
            board[5][i + 3] = row[i + 6] + '1';
        }
        copySquare(board, 3, 3, true);
        copySquare(board, 3, 3, false);
        copySquare(board, 3, 0, false);
        copySquare(board, 3, 6, false);

        while (digCount)
        {
            int x = rand() % 9;
            int y = rand() % 9;
            if (board[x][y] == '$')
                continue;
            char tmp = board[x][y];
            board[x][y] = '$';

            solveSudoku(board);
            if (result.size() == 1)
            {
                digCount--;
            }
            else
            {
                board[x][y] = tmp;
            }
        }
        // printBoard(board);
        // cout << "spaces " << player.spaces.
        size() << "\n";
        if (!checkBoard(board))
        {
            cout << "wrong board" << endl;
        }
        return board;
    }

    vector<int> getRand9()

```

```

{
    vector<int> result;
    int digit = 0;
    while (result.size() != 9)
    {
        int num = rand() % 9;
        if ((1 << num) & digit)
        {
            continue;
        }
        else
        {
            result.push_back(num);
            digit ^= (1 << num);
        }
    }
    return result;
}

void copySquare(Board &board, int src_x,
int src_y, bool isRow)
{
    int rand_tmp = rand() % 2 + 1;
    int order_first[3] = {1, 2, 0};
    int order_second[3] = {2, 0, 1};
    if (rand_tmp == 2)
    {
        order_first[0] = 2;
        order_first[1] = 0;
        order_first[2] = 1;
        order_second[0] = 1;
        order_second[1] = 2;
        order_second[2] = 0;
    }
    for (int i = 0; i < 3; i++)
    {
        if (isRow)
        {
            board[src_x][i] = board[src_x +
order_first[0]][src_y + i];
            board[src_x + 1][i] = board[src_x +
+ order_first[1]][src_y + i];
            board[src_x + 2][i] = board[src_x +
+ order_first[2]][src_y + i];
            board[src_x][i + 6] = board[src_x +
+ order_second[0]][src_y + i];
            board[src_x + 1][i + 6] = board[src_x +
src_x + order_second[1]][src_y +
i];
            board[src_x + 2][i + 6] = board[src_x +
src_x + order_second[2]][src_y +
i];
        }
        else
        {
            board[i][src_y] = board[src_x + i][src_y +
order_first[0]];
            board[i][src_y + 1] = board[src_x +
+ i][src_y + order_first[1]];
            board[i][src_y + 2] = board[src_x +
+ i][src_y + order_first[2]];
            board[i + 6][src_y] = board[src_x +
+ i][src_y + order_second[0]];
            board[i + 6][src_y + 1] = board[src_x +
src_x + i][src_y + order_second[1]];
            board[i + 6][src_y + 2] = board[src_x +
src_x + i][src_y + order_second[2]];
        }
    }
}

```

```

    }

}

char data[9][9] = {
    {'5', '3', '.', '.', '7', '.', '.', '.',},
    {'.',},
    {'6', '.', '.', '1', '9', '5', '.', '.',},
    {'.'},
    {'.', '9', '8', '.', '.', '.', '6',},
    {'.'},
    {'8', '.', '.', '6', '.', '.', '.', '3',},
    {'3'},
    {'4', '.', '.', '8', '.', '3', '.', '.', '1',},
    {'7', '.', '.', '2', '.', '.', '.', '6',},
    {'.'},
    {'6', '.', '.', '2', '.', '2', '8',},
    {'.'},
    {'.', '.', '4', '1', '9', '.', '.', '5',},
    {'.'},
    {'.', '.', '.', '8', '.', '.', '7', '9'}};

void test()
{
    SudokuPlayer player;
    vector<vector<char>> board(N, vector<char>(N, '.'));

    for (int i = 0; i < board.size(); i++)
    {
        for (int j = 0; j < board[i].size(); j++)
        {
            board[i][j] = data[i][j];
        }
    }
    bool check = player.checkBoard(board);
    if (check)
        cout << "checked" << endl;

    player.solveSudoku(board);
    player.getResult();

    cout << endl;
}

vector<Board> readFile(string filePath)
{
    ifstream infile;
    vector<Board> boards;
    infile.open(filePath);
    char data[100];
    Board tmp;
    vector<char> row;
    while (!infile.eof())
    {
        infile.getline(data, 100);
        if (data[0] == '-')
        {
            boards.push_back(Board(tmp));
            tmp.clear();
            continue;
        }
        for (int i = 0; i < strlen(data); i++)
        {
            if ('0' <= data[i] && data[i] <=

```

```

        '9') || data[i] == '$')
    {
        row.push_back(data[i]);
    }
}
tmp.push_back(vector<char>(row));
row.clear();
}
infile.close();
return boards;
}

void writeFile(vector<Board> boards, ofstream &f)
{
    for (int k = 0; k < boards.size(); k++)
    {
        for (int i = 0; i < boards[k].size(); i++)
        {
            for (int j = 0; j < boards[k][i].size(); j++)
            {
                f << boards[k][i][j] << " ";
            }
            f << "\n";
        }
        f << "----- " << k << " -----" <<
        endl;
    }
}

// map<char, string> parse(int argc, char *argv[])
{
    map<char, string> params;
    int completeBoardCount, gameNumber,
    gameLevel;
    vector<int> range;
    string inputFile;
    char opt = 0;
    while ((opt = getopt(argc, argv, "c:s:n:m:r:
    u")) != -1)
    {
        switch (opt)
        {
        case 'c':
            completeBoardCount = atoi(optarg);
            if (completeBoardCount < 1 || completeBoardCount > 1000000)
            {
                printf("11000000\n");
                exit(0);
            }
            params[opt] = string(optarg);
            break;
        case 's':
            inputFile = string(optarg);
            if (access(optarg, 0) == -1)
            {
                printf("file does not exist\n");
                exit(0);
            }
            params[opt] = string(optarg);
            break;
        case 'n':
            gameNumber = atoi(optarg);
            if (gameNumber < 1 || gameNumber >
            10000)

```

```

            {
                printf("110000\n");
                exit(0);
            }
            params[opt] = string(optarg);
            break;
        case 'm':
            gameLevel = atoi(optarg);
            if (gameLevel < 1 || gameLevel > 3)
            {
                printf("13\n");
                exit(0);
            }
            params[opt] = string(optarg);
            break;
        case 'r':
            char *p;
            p = strtok(optarg, "~");
            while (p)
            {
                range.push_back(atoi(p));
                p = strtok(NULL, "~");
            }
            if (range.size() != 2)
            {
                printf("\n");
                exit(0);
            }
            if ((range[0] >= range[1]) || range
            [0] < 20 || range[1] > 55)
            {
                printf("2055\n");
                exit(0);
            }
            params[opt] = string(optarg);
            break;
        case 'u':
            params[opt] = string();
            break;
        default:
            printf("\n");
            exit(0);
            break;
        }
    }
    return params;
}

void generateGame(int gameNumber, int gameLevel
, vector<int> digCount, ofstream &outfile,
SudokuPlayer &player)
{
    for (int i = 0; i < gameNumber; i++)
    {
        int cnt = 0;
        if (digCount.size() == 1)
        {
            cnt = digCount[0];
        }
        else
        {
            cnt = rand() % (digCount[1] -
            digCount[0] + 1) + digCount[0];
        }
        Board b = player.generateBoard(cnt);
        vector<Board> bs;
        bs.push_back(b);
        writeFile(bs, outfile);
    }
}

```

```

        outfile.close();
    }

int main(int argc, char *argv[])
{
    srand((unsigned)time(NULL));
    SudokuPlayer player;

    map<char, string> params = parse(argc, argv);
    map<char, string>::iterator it, tmp;

    int opt = 0;

    vector<int> range;
    int gameNumber;
    int gameLevel = 0;
    int solution_count = 0;

    vector<Board> boards;
    ofstream outfile;

    it = params.begin();
    while (it != params.end())
    {
        switch (it->first)
        {
        case 'c':
            outfile.open("game.txt", ios::out |
                ios::trunc);
            range.push_back(0);
            generateGame(atoi(it->second.c_str()),
                0, range, outfile, player);
            range.clear();
            break;

        case 's':
            outfile.open("sudoku.txt", ios::out |
                ios::trunc);
            boards = readFile(it->second);
            for (int i = 0; i < boards.size(); i++)
            {
                vector<Board> result = player.
                    solveSudoku(boards[i]);
                writeFile(result, outfile);
            }
            outfile.close();
            break;

        case 'n':
        case 'm':
        case 'r':
        case 'u':
            tmp = params.find('n');
            if (tmp == params.end())
            {
                printf(" n \n");
                exit(0);
            }

            gameNumber = atoi(tmp->second.c_str());
            tmp = params.find('u');
            if (tmp != params.end())
            {
                solution_count = 1;
            }
        }
    }
}

```

```

tmp = params.find('m');
if (tmp != params.end())
{
    gameLevel = atoi(tmp->second.
        c_str());
}

tmp = params.find('r');
if (tmp != params.end())
{
    char *p;
    char *pc = new char[100];
    strcpy(pc, tmp->second.c_str());
    p = strtok(pc, "~");
    while (p)
    {
        range.push_back(atoi(p));
        p = strtok(NULL, "~");
    }
}
else
{
    //
    if (gameLevel == 1)
    {
        range.push_back(20);
        range.push_back(30);
    }
    else if (gameLevel == 2)
    {
        range.push_back(30);
        range.push_back(40);
    }
    else if (gameLevel == 3)
    {
        range.push_back(40);
        range.push_back(55);
    }
    else
    {
        range.push_back(20);
        range.push_back(55);
    }
}

outfile.open("game.txt", ios::out |
    ios::trunc);
generateGame(gameNumber, gameLevel,
    range, outfile, player);
range.clear();
break;
}

// cout << it->first << ' ' << it->
second << endl;
it++;
}

return 0;
} <end_file>; game.txt:<start_file><9 $ 5 $ 3 $ 7
1 2
$ 1 2 $ $ 8 3 $ $
$ $ $ 2 7 $ 9 8 5
8 $ 9 $ 6 $ 1 2 7
1 $ $ $ 5 $ $ 6 3
4 6 3 1 2 7 $ $ $
$ $ 8 3 4 6 2 7 1
2 7 $ $ $ $ $ 3 $
$ 3 4 $ 1 $ $ $ 8

```

```

----- 0 -----<endfile>
Here is the code repository:Cow.cpp:<start_file>
#include "Cow.h"
Cow::Cow(std::string a,int b,int c,int d){
    name=a;
    l=b;
    u=c;
    m=d;
    in=0;
    state=0;
}<endfile>Cow.h:<start_file>#pragma once
#include <string>
class Cow{
public:
    std::string name;
    int l,u,m;
    int in;
    int state;
    Cow(){}
    Cow(std::string a,int b,int c,int d);
};<endfile>Farm.cpp:<start_file>#include "Farm.h"
Farm::Farm(int a){
    n=a;
    num=0;
    cow=new Cow[a];
    milk=0;
}
void Farm::addCow(Cow a){
    cow[num]=a;
    num+=1;
}
void Farm::supply(std::string a,int b){
    for(int i=0;i<n;i++){
        if(cow[i].name==a){
            cow[i].in+=b;
            break;
        }
    }
}
void Farm::startMeal(){
    for(int i=0;i<n;i++){
        if(cow[i].in==0)
            cow[i].state=0;
        if(cow[i].in>0&&cow[i].in<cow[i].l){
            cow[i].state=1;
            cow[i].in=0;
        }
        if(cow[i].in>=cow[i].l){
            cow[i].state=2;
            if(cow[i].in<=cow[i].u)
                cow[i].in=0;
            if(cow[i].in>cow[i].u)
                cow[i].in-=cow[i].u;
        }
    }
}
void Farm::produceMilk(){
    for(int i=0;i<n;i++){
        if(cow[i].state==0){
            milk+=0;
            continue;
        }
        if(cow[i].state==1){
            milk+=cow[i].m*0.5;
            continue;
        }
    }
}

```

```

        if(cow[i].state==2){
            milk+=cow[i].m;
            continue;
        }
    }
}
float Farm::getMilkProduction(){
    return milk;
}<endfile>Farm.h:<start_file>#pragma once
#include "Cow.h"
class Farm{
public:
    int n;
    int num;
    Cow* cow;
    float milk;
    Farm(int a);
    void addCow(Cow a);
    void supply(std::string a,int b);
    void startMeal();
    void produceMilk();
    float getMilkProduction();
    ~Farm(){
        delete[] cow;
    }
};<endfile>main.cpp:<start_file>#include <iostream>
#include <string>
#include "Cow.h"
#include "Farm.h"
using namespace std;

int main(){
    int n;
    cin >> n;
    Farm farm(n);
    string name;
    int l, u, m;
    for(int i = 0; i < n; ++i){
        cin >> name >> l >> u >> m;
        Cow cow(name, l, u, m);
        farm.addCow(cow);
    }

    int k;
    cin >> k;
    int t;
    int a;
    for(int i = 0; i < k; ++i){
        cin >> t;
        for(int j = 0; j < t; ++j){
            cin >> name >> a;
            farm.supply(name, a);
        }
        farm.startMeal();
        farm.produceMilk();
    }
    printf("%.1f", farm.getMilkProduction());
    return 0;
}<endfile>makefile:<start_file>main:main-3.o
Farm.o Cow.o
g++ main-3.o Farm.o Cow.o -o main

main-3.o:main-3.cpp Farm.h Cow.h
g++ -c main-3.cpp -o main-3.o

Farm.o:Farm.cpp Farm.h Cow.h
g++ -c Farm.cpp -o Farm.o

```

```

Cow.o:Cow.cpp Cow.h
    g++ -c Cow.cpp -o Cow.o

clean:
    rm *.o main<endfile>, and the input file
    is:./input/11.txt:<start_file>3
a 2 5 6
b 3 4 7
c 1 6 5
2
1 a 3
2 b 2 c 4<enfile>

```

1090

Given the following code, what is the execution result? The file is under '/app/' directory, and is run with "python3 /app/test.py" if it is a python file, "g++ -std=c++11 /app/test.cpp -o /app/test" if it is a cpp file, and "javac /app/\{class\_name\}.java" if it is a java file.

You should think step by step. Your answer should be in the following format:

Thought: <your thought>

Output:

<execution result>

1091

```

Here is the code repository:car.cpp:<start_file>
#include "car.h"
#include <iostream>
using namespace std;

Car::Car(int num, string eng):Vehicle(num, eng){}

void Car::describe(){
    cout<<"Finish building a car with "<<wheel.
    get_num()<<" wheels and a "<<engine.
    get_name()<<" engine."<<endl;
    cout<<"A car with "<<wheel.get_num()<<
    wheels and a "<<engine.get_name()<<" engine
    ."<<endl;
}

<endfile>car.h:<start_file>#pragma once
#include "vehicle.h"
using namespace std;

class Car: public Vehicle{
public:
    Car(int num, string eng);
    void describe();
};<endfile>engine.cpp:<start_file>#include "
engine.h"

Engine::Engine(string nam): name(nam) {
    cout << "Using " << nam << " engine."<<
    endl;
}

string Engine::get_name() {
    return name;
}<endfile>engine.h:<start_file>#pragma once
#include <iostream>
#include <string>
using namespace std;

```

```

class Engine {
    string name;
public:
    Engine(string);
    string get_name();
};<endfile>main.cpp:<start_file>
#include <iostream>
#include <string>
#include "wheel.h"
#include "engine.h"
#include "vehicle.h"
#include "motor.h"
#include "car.h"
using namespace std;

int main() {
    int n, type, num;
    string engine;

    cin >> n;
    for (int i=0; i<n; i++) {
        cin >> type >> num >> engine;
        switch (type) {
            case 0: {
                Vehicle v =
                Vehicle(num,
                engine);
                v.describe();
                break;
            }
            case 1: {
                Motor m = Motor(
                num, engine);
                m.describe();
                m.sell();
                break;
            }
            case 2: {
                Car c = Car(num,
                engine);
                c.describe();
                break;
            }
        }
    }
    return 0;
}<endfile>motor.cpp:<start_file>#include "motor.
h"
#include <iostream>
using namespace std;
Motor::Motor(int num, string eng):Vehicle(num,
eng){}

void Motor::describe(){
    cout<<"Finish building a motor with "<<
    wheel.get_num()<<" wheels and a "<<engine.
    get_name()<<" engine."<<endl;
    cout<<"A motor with "<<wheel.get_num()<<
    wheels and a "<<engine.get_name()<<" engine
    ."<<endl;
}

void Motor::sell(){
    cout<<"A motor is sold!"<<endl;
}<endfile>motor.h:<start_file>#pragma once
#include "vehicle.h"
using namespace std;

```

1093

```

class Motor: public Vehicle{
    public:
        Motor(int num, string eng);
        void describe();
        void sell();
};<endfile>vehicle.cpp:<start_file>#include "vehicle.h"
#include <iostream>
using namespace std;

Vehicle::Vehicle(int num, string eng): engine(eng)
, wheel(num){}

void Vehicle::describe(){
    cout<<"Finish building a vehicle with "<<
    wheel.get_num()<<" wheels and a "<<engine.
    get_name()<<" engine."<<endl;
    cout<<"A vehicle with "<<wheel.get_num()<<
    " wheels and a "<<engine.get_name()<<
    engine."<<endl;
}<endfile>vehicle.h:<start_file>#pragma once
#include "wheel.h"
#include "engine.h"

using namespace std;

class Vehicle{
    public:
        Engine engine;
        Wheel wheel;
        Vehicle(int num, string eng);
        void describe();
};

};<endfile>wheel.cpp:<start_file>#include "wheel.h"

Wheel::Wheel(int num): number(num) {
    cout << "Building " << number << "
    wheels." << endl;
}

int Wheel::get_num() {
    return number;
}<endfile>wheel.h:<start_file>#pragma once
#include <iostream>
using namespace std;

class Wheel {
    int number;
public:
    Wheel(int);
    int get_num();
};<endfile>, and the input file is:./input/6.
txt:<start_file>4
0 3 Gasoline
2 4 Hybrid
1 2 Electric
0 6 Magic<endfile>

```

1094

```

Here is the code repository:24_game.py:<
start_file>import itertools
import time
import math

# Operators
OP_CONST = 0 # Constant
OP_ADD = 1 # Addition
OP_SUB = 2 # Subtraction

```

```

OP_MUL = 3 # Multiplication
OP_DIV = 4 # Division
OP_POW = 5 # Exponentiation

OP_SQRT = 6 # SquareRoot
OP_FACT = 7 # Factorial
OP_LOG = 8 # Logarithm
OP_C = 9 # Combinations
OP_P = 10 # Permutations

# List of basic operators
operators = [OP_ADD,
             OP_SUB,
             OP_MUL,
             OP_DIV]

# List of advanced operators
advanced_operators = [OP_POW,
                      OP_LOG,
                      OP_C,
                      OP_P]

# List of unary operators
_unary_operators = [OP_SQRT,
                     OP_FACT]

# List of enabled unary operators
unary_operators = []

# Symbol of operators
symbol_of_operator = {OP_ADD: "%s+%s",
                      OP_SUB: "%s-%s",
                      OP_MUL: "%s*%s",
                      OP_DIV: "%s/%s",
                      OP_POW: "%s^%s",
                      OP_SQRT: "sqrt(%s)",
                      OP_FACT: "%s!",
                      OP_LOG: "log_%s(%s)",
                      OP_C: "C(%s, %s)",
                      OP_P: "P(%s, %s)"}

# Priority of operators
priority_of_operator = {OP_ADD: 0,
                        OP_SUB: 0,
                        OP_MUL: 1,
                        OP_DIV: 1,
                        OP_POW: 2,
                        OP_LOG: 3,
                        OP_C: 3,
                        OP_P: 3,
                        OP_SQRT: 3,
                        OP_FACT: 4,
                        OP_CONST: 5}

# Whether operator is commutative
is_operator_commutative = {OP_ADD: True,
                           OP_SUB: False,
                           OP_MUL: True,
                           OP_DIV: False,
                           OP_POW: False,
                           OP_LOG: False,
                           OP_C: False,
                           OP_P: False}

# Whether inside bracket is needed when
# rendering
need_brackets = {OP_ADD: True,
                 OP_SUB: True,
                 OP_MUL: True,
                 OP_DIV: True,
                 OP_POW: True,
                 OP_LOG: True,
                 OP_C: True,
                 OP_P: True}

```

```

OP_DIV: True,
OP_POW: True,
OP_FACT: True,
OP_SQRT: False,
OP_LOG: False,
OP_C: False,
OP_P: False}

def permutation(n, k):
    return math.factorial(n)/math.factorial(k)

def combination(n, k):
    return permutation(n, k)/math.factorial(n-k)

def evaluate_operation(op, a, b=None):
    """
    Evaluate an operation on a and b.
    """
    if op == OP_ADD: return a + b
    if op == OP_SUB: return a - b
    if op == OP_MUL: return a * b

    try:
        if op == OP_POW and abs(a) < 20 and abs(b) < 20:
            return a ** b

        if op == OP_FACT and a < 10:
            return math.factorial(a)

        if op == OP_C and 0 < b <= a <= 13:
            return combination(a, b)

        if op == OP_P and 0 < b <= a <= 13:
            return permutation(a, b)

        if op == OP_SQRT and a < 1000000:
            return math.sqrt(a)

        if op == OP_DIV: return a / b
        if op == OP_LOG: return math.log(b, a)
    except (ZeroDivisionError, ValueError, TypeError):
        pass
    except OverflowError:
        print(a, b)

    return float("NaN")

def fit_to_int(x, eps=1e-9):
    """
    Convert x to int if x is close to an integer.
    """
    try:
        if abs(round(x) - x) <= eps:
            return round(x)
        else:
            return x
    except ValueError:
        return float("NaN")
    except TypeError:
        return float("NaN")

```

```

class Node:
    def __init__(self, value=None, left=None, right=None, op=OP_CONST):
        if op not in unary_operators \
           and op != OP_CONST and
           is_operator_commutative[op] \
           and str(left) > str(right):
            left, right = right, left

        self._value = value
        self._str_cache = None
        self.left = left
        self.right = right
        self.op = op

    @property
    def value(self):
        if self._value is None:
            assert self.op != OP_CONST

        if self.op in unary_operators:
            self._value = evaluate_operation(
                self.op, self.left.value)
        else:
            self._value = evaluate_operation(
                self.op, self.left.value, self.right.value)

        self._value = fit_to_int(self._value)
    return self._value

    def __str__(self):
        if self._str_cache is None:
            self._str_cache = self._str()
        return self._str_cache

    def _str(self):
        # Constant
        if self.op == OP_CONST:
            return str(self._value)

        # Unary operator
        elif self.op in unary_operators:
            str_left = str(self.left)

            if need_brackets[self.op] \
               and priority_of_operator[self.left.op] <
                  priority_of_operator[self.op]:
                str_left = "(" + str_left + ")"

            return symbol_of_operator[self.op] % str_left

        # Other operator
        else:
            str_left = str(self.left)
            str_right = str(self.right)

            # Add brackets inside
            if need_brackets[self.op] \
               and priority_of_operator[self.left.op] <
                  priority_of_operator[self.op]:
                str_left = "(" + str_left + ")"

            str_left = "(" + str_left + ")"

```

```

        if need_brackets[self.op] \
            and (priority_of_operator[
                self.right.op] <
                priority_of_operator[self.op]
            or (priority_of_operator[
                self.right.op] ==
                priority_of_operator[
                    self.op]
            and not
                is_operator_commutative
                [self.op])):
            str_right = "(" + str_right + ")"

    # Render
    return symbol_of_operator[self.op] %
        (str_left, str_right)

def enumerate_nodes(node_list, callback,
max_depth):
    # Found an expression
    if len(node_list) == 1:
        callback(node_list[0])

    # Constrain maximum depth
    if max_depth == 0:
        return

    # Non-unary operators
    for left, right in itertools.permutations(
        node_list, 2):
        new_node_list = node_list.copy()
        new_node_list.remove(left)
        new_node_list.remove(right)

        for op in operators:
            enumerate_nodes(new_node_list + [
                Node(left=left, right=right, op=op)],
                callback, max_depth-1)

            if not is_operator_commutative[op] \
                and str(left) != str(right):
                enumerate_nodes(new_node_list + [
                    Node(left=right, right=left, op=
                        op)], callback, max_depth-1)

    # Unary operators
    for number in node_list:
        new_node_list = node_list.copy()
        new_node_list.remove(number)

        for op in unary_operators:
            new_node = Node(left=number, op=op)
            if new_node.value == number.value:
                continue

            enumerate_nodes(new_node_list + [
                new_node], callback, max_depth-1)

class CallbackFindTarget:
    def __init__(self, target):
        self.target = target
        self.results = []
        self.duplication_count = 0
        self.enumeration_count = 0

    def __call__(self, node):

```

```

        if node.value == self.target and str(
            node) not in self.results:
            print(self.target, "=", node)
            self.results.append(str(node))
        elif node.value == self.target:
            self.duplication_count += 1

        self.enumeration_count += 1

    def show(self, execution_time):
        print()
        print("%d solution(s) in %.3f seconds" %
            (len(self.results), execution_time))
        print("%d duplication(s)" % self.
            duplication_count)
        print("%d combination(s)" % self.
            enumeration_count)

class CallbackAllTarget:
    def __init__(self):
        self.results = {}
        self.enumeration_count = 0

    def __call__(self, node):
        try:
            int(node.value)
        except ValueError:
            return

        if node.value not in self.results \
            and int(node.value) == node.value:
            :
                self.results[int(node.value)] = node

        self.enumeration_count += 1

    def __str__(self):
        string = ""
        for value in sorted(self.results.keys()):
            :
                string += "%d = %s" % (value, str(
                    self.results[value]))
                string += "\n"
        return string

    def show(self, execution_time):
        print(self)
        print()
        print("%d targets(s) in %.3f seconds" %
            (len(self.results), execution_time))
        print("%d combination(s)" % self.
            enumeration_count)

def select_yes_no(prompt, default=False):
    answer = input(prompt).strip().lower()
    if answer == "y":
        return True
    if answer == "n":
        return False
    return default

def select_int(prompt, default):
    try:
        return int(input(prompt).strip())
    except ValueError:
        return default

```

```

def main():
    global operators
    global unary_operators

    unary_operators_allowed = False
    enumerate_all = False

    if not enumerate_all:
        target = 24
        callback = CallbackFindTarget(target=target)

    if enumerate_all:
        callback = CallbackAllTarget()
    else:
        callback = CallbackFindTarget(target=target)

    with open('input.txt', 'r') as file:
        inputs = [int(i) for line in file for i
                  in line.split() if i != ""]
    node_list = [Node(value=i) for i in inputs]

    enumerate_nodes(node_list, callback,
                    max_depth=len(node_list)-1+
                    unary_operators_allowed)

main()<enfile>, and the input file is: input.
txt:<start_file>4 4 7 7<end_file>

```

## C.4 SC

### C.4.1 Tasks Descriptions

The scientific computing component of SURGE consists of 4 carefully curated areas, aiming to evaluate model performance on computational tasks that exhibit a time-consuming nature as well as applicational values in scientific computing areas. In this section, we provide a detailed description of each component.

**Numerical Optimization.** In this task, the model is given a program that solves an optimization problem through gradient descent. The query may be the optimized value (*min*) or the optimal point (*argmin*). We carefully select four functions, which consist of: a simple quadratic function, Rosenbrock Function, Himmelblau's Function, and a polynomial function with linear constraints. For each function, we will select multiple different hyperparameter configurations to assess the model's performance. These four functions provide a systematic evaluation of the model's potential to serve as a surrogate model in this field. As the quadratic function is solvable without need the to run the gradient descent, the model may solve it through world knowledge. The Rosenbrock function is known

for its narrow, curved valley containing the global minimum, making it difficult for optimization algorithms to converge. Therefore the output is highly dependent on hyperparameters (initial point, learning rate, maximum steps), thus the model must execute code in its reasoning process to acquire the answer. Himmelblau's function has multiple local minima, also posing sensitivity to hyperparameters.

**PDE Solving.** We consider three types of Partial Differential Equations: the 1D Heat Equation, the 2D Wave Equation, and the 2D Laplace Equation. For the 1D Heat Equation, we focus on solving the following equation:

$$\frac{\partial u}{\partial t} = \alpha \frac{\partial^2 u}{\partial x^2}. \quad (1)$$

For the 2D Laplace Equation, we aim to solve the equation:

$$\frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} = 0. \quad (2)$$

Lastly, for the 2D Wave Equation, we work on solving the following equation:

$$\frac{\partial^2 u}{\partial t^2} = c^2 \left( \frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} \right). \quad (3)$$

We solve 1D Heat Equation and 2D Wave Equation using the Explicit Finite Difference Method. For the 2D Laplace Equation, we solve it using the Gauss-Seidel Method. The model is then queried on the values of  $u$  and  $x$ .

**Fourier Transform (FFT)** We implement FFT using the Cooley-Tukey Algorithm and query the model to give the magnitude of the top 10 values.

**ODE Solving** For solving ordinary differential equations, we constructed three different equations and implemented the Euler Method and the Runge-Kutta Method so solve these equations.

### C.4.2 Evaluation Metrics

**Relative Absolute Error (RAE).** Given a scalar ground truth value  $p$  and a model prediction  $\hat{p}$ , the Relative Absolute Error (RAE) is defined as:

$$RAE(\hat{p}, p) = \frac{|\hat{p} - p|}{|p|}. \quad (4)$$

For cases involving multiple entries, such as tensors or vectors, the following alignment procedure is applied: (1) if the prediction contains fewer elements than the ground truth, the prediction is padded with

zeros until it matches the length of the ground truth; (2) if the prediction has more elements than the ground truth, it is truncated to match the ground truth length. The average RAE is then computed by averaging the RAE for each corresponding element.

**Exact Matching.** For tasks involving position-based predictions, such as binary search, we adapt exact matching, as the accuracy of the algorithm is determined by comparing the exactness of the estimated result to the true result. This evaluation method checks if the estimated solution matches the ground truth exactly, typically using string or sequence matching. For such tasks, an exact match is considered a success, and any discrepancy between the ground truth and the estimate results in failure. Formally, given a string  $s$  and the model's prediction  $\hat{s}$ , the Exact Matching is given by:

$$\text{EM}(s, \hat{s}) = \mathbb{1}[s = \hat{s}] \quad (5)$$

where  $\mathbb{1}[\cdot]$  is the indicator function.

### C.4.3 System Prompts

#### Zero-shot Chain-of-Thought:

You are an expert in gradient\_descent programming.  
Please execute the above code with the input provided and return the output. You should think step by step.  
Your answer should be in the following format:  
Thought: <your thought>  
Output: <execution result>  
Please follow this format strictly and ensure the Output section contains only the required result without any additional text.

#### Zero-shot:

You are an expert in gradient\_descent programming.  
Please execute the given code with the provided input and return the output.  
Make sure to return only the output in the exact format as expected.  
Output Format:  
Output: <result>

#### Few-shot Chain-of-Thought:

You are an expert in gradient\_descent programming.  
Please execute the above code with the input provided and return the output. You should think step by step.  
Your answer should be in the following format:  
Thought: <your thought>  
Output: <execution result>

Please follow this format strictly and ensure the Output section contains only the required result without any additional text.

Here are some examples:  
&{{examples here}}

### C.4.4 Demo Questions

```
code: """
import numpy as np
import argparse

def f(t, y):
    """dy/dt = -y"""
    return -y

def euler_method(f, y0, t0, t_end, h,
additional_args=None):
    t_values = np.arange(t0, t_end, h)
    y_values = [y0]
    v_values = [additional_args] if
additional_args is not None else [None]

    for t in t_values[:-1]:
        if additional_args:
            y_next, v_next = y_values[-1] + h *
f(t, y_values[-1])[0], v_values[-1]
            + h * f(t, y_values[-1], v_values
[-1])[1]
            y_values.append(y_next)
            v_values.append(v_next)
        else:
            y_next = y_values[-1] + h * f(t,
y_values[-1])
            y_values.append(y_next)

    return t_values, np.array(y_values), np.
array(v_values) if v_values[0] is not None
else None

def main():
    parser = argparse.ArgumentParser()
    parser.add_argument("--y0", type=float,
default=1.0)
    parser.add_argument("--t0", type=float,
default=0.0)
    parser.add_argument("--t_end", type=float,
default=10.0)
    parser.add_argument("--h", type=float,
default=0.1)
    args = parser.parse_args()

    y0_1 = args.y0
    t0 = args.t0
    t_end = args.t_end
    h = args.h

    t_values, y_values, _ = euler_method(f,
y0_1, t0, t_end, h)
    print(f"\{y_values[-1]:.4f}\n")

if __name__ == "__main__":
    main()

"""

command: """
python euler_3.py --y0 12 --t0 0.0 --t_end 74
--h 0.36

```

1195

1196

```

"""
code:"""
import numpy as np
import argparse

def gradient_descent(func, grad_func,
initial_guess, learning_rate=0.1, tolerance=1e
-6, max_iter=1000):
    x = initial_guess
    for _ in range(max_iter):
        grad = grad_func(x)
        x = x - learning_rate * grad
        if np.abs(grad) < tolerance:
            break
    return x, func(x)

# Function and its gradient
def func(x):
    return (x - 3)**2 + 5

def grad_func(x):
    return 2 * (x - 3)

def main():
    parser = argparse.ArgumentParser()
    parser.add_argument("--initial_guess", type
=float, default=0.0)
    parser.add_argument("--learning_rate", type
=float, default=0.1)
    parser.add_argument("--tolerance", type=
float, default=1e-6)
    parser.add_argument("--max_iter", type=int,
    default=1000)
    args = parser.parse_args()

    # Test with initial guess
    initial_guess = args.initial_guess
    optimal_x, optimal_value = gradient_descent(
        func, grad_func, initial_guess)
    # optimal x
    print(f"{optimal_x:.3f}")

if __name__ == "__main__":
    main()
"""

command:""
python gd_ox.py --initial_guess -5.0 --
learning_rate 0.01 --max_iter 5000
"""

```

```

code:"""
import numpy as np
import argparse

def solve_heat_eq(L, T, alpha, Nx, Nt):
    # L: length of the rod
    # T: total time
    # alpha: thermal diffusivity
    # Nx: number of spatial steps
    # Nt: number of time steps

    dx = L / (Nx - 1)
    dt = T / Nt
    r = alpha * dt / dx**2

    # Initial condition: u(x, 0) = sin(pi * x)

```

```

x = np.linspace(0, L, Nx)
u = np.sin(np.pi * x)

# Time stepping
for n in range(Nt):
    u_new = u.copy()
    for i in range(1, Nx - 1):
        u_new[i] = u[i] + r * (u[i-1] - 2*u[
            i] + u[i+1])
    u = u_new
return x, u

def parse_input():
    parser = argparse.ArgumentParser(
        description="Solve the 1D Heat Equation")
    parser.add_argument('--L', type=float,
    required=True, help="Length of the rod")
    parser.add_argument('--T', type=float,
    required=True, help="Total time")
    parser.add_argument('--alpha', type=float,
    required=True, help="Thermal diffusivity")
    parser.add_argument('--Nx', type=int,
    required=True, help="Number of spatial
    points")
    parser.add_argument('--Nt', type=int,
    required=True, help="Number of time steps")
    return parser.parse_args()

def main():
    args = parse_input()
    x, u = solve_heat_eq(args.L, args.T, args.
    alpha, args.Nx, args.Nt)
    np.set_printoptions(threshold=np.inf,
    linewidth=np.inf)
    formatted_x = np.vectorize(lambda x: f"{x
    :.4e}")(x)
    print(f"{formatted_x}")

if __name__ == "__main__":
    main()

"""

command:""
python heat_eq_x.py --L 36 --T 62 --alpha 91 --
Nx 170 --Nt 860
"""

```

```

code:"""
import numpy as np
import argparse

def gradient_descent(func, grad_func,
initial_guess, learning_rate=0.1, tolerance=1e
-6, max_iter=1000):
    x = initial_guess
    for _ in range(max_iter):
        grad = grad_func(x)
        x = x - learning_rate * grad
        if np.abs(grad) < tolerance:
            break
    return x, func(x)

# Function and its gradient
def func(x):
    return (x - 3)**2 + 5

def grad_func(x):
    return 2 * (x - 3)

```

```

def main():
    parser = argparse.ArgumentParser()
    parser.add_argument("--initial_guess", type=float, default=0.0)
    parser.add_argument("--learning_rate", type=float, default=0.1)
    parser.add_argument("--tolerance", type=float, default=1e-6)
    parser.add_argument("--max_iter", type=int, default=1000)
    args = parser.parse_args()

    # Test with initial guess
    initial_guess = args.initial_guess
    optimal_x, optimal_value = gradient_descent(
        func, grad_func, initial_guess)
    # optimal x
    print(f"optimal_x:{optimal_x:.3f}")

if __name__ == "__main__":
    main()
"""

command: """
python gd_ox.py --initial_guess -10.0 --
learning_rate 0.001 --max_iter 100
"""

```

1203

```

code: """
import numpy as np
import argparse

# Objective function: f(x, y) = x^2 + y^2
def objective(x, y):
    return x**2 + y**2

# Gradient of the objective function: f(x, y) =
# (2x, 2y)
def gradient(x, y):
    return np.array([2 * x, 2 * y])

# Projection function onto the constraint x + y
# = 1
def projection(x, y):
    # Since the constraint is x + y = 1, we can
    # project the point (x, y) onto the line
    # by solving the system: x' + y' = 1
    # Let x' = x - (x + y - 1)/2, and y' = y -
    # (x + y - 1)/2
    adjustment = (x + y - 1) / 2
    return np.array([x - adjustment, y -
    adjustment])

def projected_gradient_descent(learning_rate
=0.1, max_iter=1000, tolerance=1e-6,
initial_guess=(0.0, 0.0)):
    x, y = initial_guess

    for _ in range(max_iter):
        # Compute the gradient of the objective
        # function
        grad = gradient(x, y)

        # Update the variables by moving in the
        # opposite direction of the gradient
        x, y = np.array([x, y]) - learning_rate
        * grad

```

1204

```

# Project the updated point onto the
# constraint set (x + y = 1)
x, y = projection(x, y)

# Check if the gradient is small enough
# to stop
if np.linalg.norm(grad) < tolerance:
    break

return x, y, objective(x, y)

def main():
    parser = argparse.ArgumentParser()
    parser.add_argument("--initial_guess_x",
    type=float, default=0.0)
    parser.add_argument("--initial_guess_y",
    type=float, default=0.0)
    parser.add_argument("--learning_rate", type=
    float, default=0.1)
    parser.add_argument("--tolerance", type=
    float, default=1e-6)
    parser.add_argument("--max_iter", type=int,
    default=1000)
    args = parser.parse_args()

    initial_guess = (args.initial_guess_x, args.
    initial_guess_y)
    optimal_x, optimal_y, optimal_value =
    projected_gradient_descent(args.
    learning_rate, args.max_iter, args.
    tolerance, initial_guess)
    print(f"optimal_x:{optimal_x:.4e}, {optimal_y:.4e}")

if __name__ == "__main__":
    main()
"""

command: """
python gd_pgdx.py --initial_guess_x 32.14 --
initial_guess_y 46.04 --learning_rate 0.01 --
max_iter 1000
"""

```

1205

## C.5 TC

### C.5.1 Tasks Descriptions of Time Consuming (TC)

The time consuming component of SURGE is comprised of 4 tasks in for computationally expensive areas, covering a spectrum of Linear Algebra, Sorting, Searching, Monte Carlo Simulations and String Matching Programs. Some of these tasks take hours to complete, showing their potential to benchmark LLM's ability to reason through lengthy computations.

1206

1207

1208

1209

1210

1211

1212

1213

1214

1215

1216

**Linear Algebra.** In this task, we are focused on acquiring key properties in linear algebra given square matrices of varying sizes. In particular, we query the model on solving LU decomposition, QR decomposition, the largest eigenvalue and eigenvector using the power method, and the inversion matrix.

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

**Sorting And Searching.** We include four classical algorithmic problems in this area, namely Hamiltonian Cycle, Traveling Salesman Problem (TSP), Sorting an array of real numbers and Searching. For Hamiltonian Cycle, we adopt the backtracking algorithm. Specifically, we randomly generate graphs with vertices from 4 to 100 and ask the model to find whether a Hamiltonian cycle exists. For TSP, we implement a naive brute-force algorithm and ask the model to find the length of the optimal path. For Sorting, we adopt the bubble sort, quick sort, and merge sort algorithms. For each algorithm, we consider different list sizes from 5 to 100 and generate 10 test cases for each list size. The evaluation metric is the rank correlation (also Spearman's  $\rho$ ). Lastly, for searching, we adopt binary search and query the model on randomly generated lists of varying sizes.

1242  
1243  
1244  
1245  
1246  
1247  
1248

**Monte Carlo Estimation.** We adopt Monte Carlo simulation to estimate the values of specific real numbers (e.g.  $\pi, e$ ), as well as a future stock price prediction that follows the Brownian motion. We alter the number of samples used in Monte Carlo estimation, resulting in varying program outcomes.

1249  
1250  
1251  
1252  
1253

**String Matching Program.** We adopt the naive string matching, KMP, and Rabin-Karp algorithms. For each algorithm, we randomly generate text and pattern with varying lengths, and query the model on the existence and position of the matching.

### C.5.2 Evaluation Metrics

1255  
1256  
1257  
1258  
1259

**Rank Correlation.** Rank Correlation (Spearman, 1904), also referred to as Spearman's  $\rho$ , is used to assess sorting tasks by measuring the correlation between the estimated ordinal ranking and the ground truth, which can be written as:

$$\text{RankCorr} = \frac{\text{Cov}(x_{1:N}, y_{1:N})}{\sigma(x_{1:N})\sigma(y_{1:N})} \quad (6)$$

1260  
1261  
1262  
1263  
1264

where  $x_{1:N}$  and  $y_{1:N}$  denote the true and estimated rankings, respectively, and Cov and  $\sigma$  represent the covariance and standard deviation of the respective sequences.

### C.5.3 System Prompts

#### Zero-shot Chain-of-Thought:

You are an expert in string\_matching programming.  
Please execute the above code with the input provided and return the output. You should think step by step.

Your answer should be in the following format:  
Thought: <your thought>  
Output: <execution result>  
Please follow this format strictly and ensure the Output section contains only the required result without any additional text.

#### Zero-shot:

You are an expert in string\_matching programming.  
Please execute the given code with the provided input and return the output.  
Make sure to return only the output in the exact format as expected.

Output Format:  
Output: <result>

#### Few-shot Chain-of-Thought:

You are an expert in string\_matching programming.  
Please execute the above code with the input provided and return the output. You should think step by step.  
Your answer should be in the following format:  
Thought: <your thought>  
Output: <execution result>  
Please follow this format strictly and ensure the Output section contains only the required result without any additional text.

Here are some examples:  
{examples here}}

### C.5.4 Demo Questions

```
code:'''
import itertools
import math
import sys
import argparse
def euclidean_distance(p1, p2):
    """Calculate the Euclidean distance between two points"""
    return math.sqrt((p1[0] - p2[0])**2 + (p1[1] - p2[1])**2)

def tsp_bruteforce(positions):
    """Brute-force TSP solver"""
    n = len(positions)
    min_path = None
    min_distance = float('inf')

    # Generate all possible permutations of the cities (excluding the starting point)
    for perm in itertools.permutations(range(1, n)):
        path = [0] + list(perm) # Start at city 0
        distance = 0
        # Calculate the total distance of the current permutation
        for i in range(1, len(path)):
            distance += euclidean_distance(
                positions[path[i-1]], positions[path[i]])
```

1268

1269

1270

1271

1272

1273

1274

```

    # Compare the distance with the minimum
    # distance found so far
    if distance < min_distance:
        min_distance = distance
        min_path = path

    return min_path, min_distance

def parse_positions(positions_str):
    """Convert the string input back to a list
    of tuples"""
    positions = []
    for pos in positions_str.split():
        x, y = map(float, pos.split(','))
        positions.append((x, y))
    return positions

def main():
    parser = argparse.ArgumentParser()
    parser.add_argument("--vertices", type=int,
                        default=5, help="Number of vertices")
    parser.add_argument("--positions", type=str,
                        default="0,0 1,1 2,2 3,3 4,4", help="List
                        of positions in the format 'x,y'")
    args = parser.parse_args()

    vertices = args.vertices
    positions_str = args.positions

    # Parse positions
    positions = parse_positions(positions_str)

    # Solve TSP using brute force
    path, distance = tsp_bruteforce(positions)

    print(f"{distance:.2f}")

if __name__ == "__main__":
    main()
"""

command:"""
python tsp.py --vertices 3 --positions
"8.51,4.18 8.1,7.92 1.57,0.49"
"""

```

1275

```

code:"""
import itertools
import math
import sys
import argparse
def euclidean_distance(p1, p2):
    """Calculate the Euclidean distance between
    two points"""
    return math.sqrt((p1[0] - p2[0])**2 + (p1
    [1] - p2[1])**2)

def tsp_bruteforce(positions):
    """Brute-force TSP solver"""
    n = len(positions)
    min_path = None
    min_distance = float('inf')

    # Generate all possible permutations of the
    # cities (excluding the starting point)
    for perm in itertools.permutations(range(1,
    n)):
```

1276

```

path = [0] + list(perm) # Start at city
0
distance = 0
# Calculate the total distance of the
current permutation
for i in range(1, len(path)):
    distance += euclidean_distance(
        positions[path[i-1]], positions[path
        [i]])

    # Compare the distance with the minimum
    # distance found so far
    if distance < min_distance:
        min_distance = distance
        min_path = path

return min_path, min_distance

def parse_positions(positions_str):
    """Convert the string input back to a list
    of tuples"""
    positions = []
    for pos in positions_str.split():
        x, y = map(float, pos.split(','))
        positions.append((x, y))
    return positions

def main():
    parser = argparse.ArgumentParser()
    parser.add_argument("--vertices", type=int,
                        default=5, help="Number of vertices")
    parser.add_argument("--positions", type=str,
                        default="0,0 1,1 2,2 3,3 4,4", help="List
                        of positions in the format 'x,y'")
    args = parser.parse_args()

    vertices = args.vertices
    positions_str = args.positions

    # Parse positions
    positions = parse_positions(positions_str)

    # Solve TSP using brute force
    path, distance = tsp_bruteforce(positions)

    print(f"{distance:.2f}")

if __name__ == "__main__":
    main()
"""

command:"""
python tsp.py --vertices 3 --positions
"0.9,2.44 4.67,0.82 3.8,5.73"
"""

```

1277

```

code:"""
import itertools
import math
import sys
import argparse
def euclidean_distance(p1, p2):
    """Calculate the Euclidean distance between
    two points"""
    return math.sqrt((p1[0] - p2[0])**2 + (p1
    [1] - p2[1])**2)

def tsp_bruteforce(positions):
```

1278

```

"""Brute-force TSP solver"""
n = len(positions)
min_path = None
min_distance = float('inf')

# Generate all possible permutations of the
# cities (excluding the starting point)
for perm in itertools.permutations(range(1,
n)):
    path = [0] + list(perm) # Start at city
    0
    distance = 0
    # Calculate the total distance of the
    current permutation
    for i in range(1, len(path)):
        distance += euclidean_distance(
            positions[path[i-1]], positions[path
            [i]])
    # Compare the distance with the minimum
    distance found so far
    if distance < min_distance:
        min_distance = distance
        min_path = path

return min_path, min_distance

def parse_positions(positions_str):
    """Convert the string input back to a list
    of tuples"""
    positions = []
    for pos in positions_str.split():
        x, y = map(float, pos.split(','))
        positions.append((x, y))
    return positions

def main():
    parser = argparse.ArgumentParser()
    parser.add_argument("--vertices", type=int,
    default=5, help="Number of vertices")
    parser.add_argument("--positions", type=str,
    default="0,0 1,1 2,2 3,3 4,4", help="List
    of positions in the format 'x,y'")
    args = parser.parse_args()

    vertices = args.vertices
    positions_str = args.positions

    # Parse positions
    positions = parse_positions(positions_str)

    # Solve TSP using brute force
    path, distance = tsp_bruteforce(positions)

    print(f"{distance:.2f}")

if __name__ == "__main__":
    main()

"""

command:"""
python tsp.py --vertices 3 --positions
"7.63,4.72 1.07,1.42 8.36,5.63"
"""

```

1279

```

code:"""
import sys
import argparse

```

1280

```

def binary_search(arr, target):
    """Binary Search algorithm"""
    low = 0
    high = len(arr) - 1

    while low <= high:
        mid = (low + high) // 2 # Find the
        middle element
        if arr[mid] == target:
            return mid # Target found at index
            mid
        elif arr[mid] < target:
            low = mid + 1 # Target is in the
            right half
        else:
            high = mid - 1 # Target is in the
            left half

    return -1 # Target not found

def parse_input(input_str):
    """Parse input string into a list of
    integers"""
    return list(map(int, input_str.split()))

def main():
    parser = argparse.ArgumentParser(
        description="Binary Search Algorithm")
    parser.add_argument('--list', type=str,
    required=True, help="Input sorted list of
    integers")
    parser.add_argument('--target', type=int,
    required=True, help="Target integer to
    search")
    args = parser.parse_args()

    input_list = parse_input(args.list)

    result = binary_search(input_list, args.
    target)

    if result != -1:
        print(f"Target found at index: {result
        }")
    else:
        print("Target not found")

if __name__ == "__main__":
    main()

"""

command:"""
python binary_search.py --list "-334 -200 180
936 973" --target -771
"""

```

1281

```

code:"""
import itertools
import math
import sys
import argparse
def euclidean_distance(p1, p2):
    """Calculate the Euclidean distance between
    two points"""
    return math.sqrt((p1[0] - p2[0])**2 + (p1
    [1] - p2[1])**2)

```

1282

```

def tsp_bruteforce(positions):
    """Brute-force TSP solver"""
    n = len(positions)
    min_path = None
    min_distance = float('inf')

    # Generate all possible permutations of the
    # cities (excluding the starting point)
    for perm in itertools.permutations(range(1,
        n)):
        path = [0] + list(perm) # Start at city
        0
        distance = 0
        # Calculate the total distance of the
        # current permutation
        for i in range(1, len(path)):
            distance += euclidean_distance(
                positions[path[i-1]], positions[path
                [i]])

        # Compare the distance with the minimum
        # distance found so far
        if distance < min_distance:
            min_distance = distance
            min_path = path

    return min_path, min_distance

def parse_positions(positions_str):
    """Convert the string input back to a list
    of tuples"""
    positions = []
    for pos in positions_str.split():
        x, y = map(float, pos.split(','))
        positions.append((x, y))
    return positions

def main():
    parser = argparse.ArgumentParser()
    parser.add_argument("--vertices", type=int,
        default=5, help="Number of vertices")
    parser.add_argument("--positions", type=str,
        default="0,0 1,1 2,2 3,3 4,4", help="List
        of positions in the format 'x,y'")
    args = parser.parse_args()

    vertices = args.vertices
    positions_str = args.positions

    # Parse positions
    positions = parse_positions(positions_str)

    # Solve TSP using brute force
    path, distance = tsp_bruteforce(positions)

    print(f"distance:{.2f}")

if __name__ == "__main__":
    main()

"""

command:"""
python tsp.py --vertices 10 --positions
"6.81,5.28 9.95,8.98 0.63,0.11 8.84,0.55
9.03,9.98 6.22,2.7 2.99,9.11 0.54,9.36
3.08,4.15 5.73,1.86"
"""

```

1283

## C.6 BG

1284

Table 10: Language usage count across different categories in the BG subset.

	Java	Python	C++
	51	45	54

Table 11: Details of bug types in BG dataset and how many times each kind of bug appears in different languages.

Error Type	Java	Python3	CPP
== and = confusion	5	6	5
undefined keywords	6	3	5
parentheses mismatch	5	5	6
indexing error	10	9	11
undefined objects	11	9	8
unclosed string	7	5	7
conditional statement error	10	8	9
undefined methods	8	3	6
colon missing	5	7	8
wrong comment mark	9	1	9
variable value error	2	2	4
operation error	2	2	3
other error	4	2	1
statement separation	4	0	7
indentation error	0	4	0
Double Bugs	10	8	10
Triple Bugs	12	10	11
Quadruple Bugs	8	5	9

In BG, the distribution of language usage across categories is shown in Table 10, indicating a balanced usage of Java, Python, and C++. Table 11 presents a detailed breakdown of bug types and their frequency across different languages. This distribution allows us to assess the model's ability to handle a variety of bugs across multiple programming languages.

1285

1286

1287

1288

1289

1290

1291

1292

### C.6.1 System Prompts

1293

#### Zero-shot Chain-of-Thought:

1294

Given the following code, what is the execution result? The file is under '/app/' directory, and is run with "python3 /app/test.py" if it is a python file, "g++ -std=c++11 /app/test.cpp -o /app/test" if it is a cpp file, and "javac /app /\{class\_name\}.java java -cp /app \{class\_name\}" if it is a java file.  
You should think step by step. Your answer should be in the following format:

1295

Thought: <your thought>  
Output:  
<execution result>

1296

1297

### Zero-shot:

Given the following code, what is the execution result? The file is under '/app/' directory, and is run with "python3 /app/test.py" if it is a python file, "g++ -std=c++11 /app/test.cpp -o /app/test /app/test" if it is a cpp file, and "javac /app/\{class\_name\}.java java -cp /app \{class\_name\}" if it is a java file.  
Your answer should be in the following format:  
Output:  
<execution result>

1298

1299

### Few-shot Chain-of-Thought:

Given the following code, what is the execution result? The file is under '/app/' directory, and is run with "python3 /app/test.py" if it is a python file, "g++ -std=c++11 /app/test.cpp -o /app/test /app/test" if it is a cpp file, and "javac /app/\{class\_name\}.java java -cp /app \{class\_name\}" if it is a java file.  
You should think step by step. Your answer should be in the following format:  
Thought: <your thought>  
Output:  
<execution result>  
Following are 4 examples:  
{examples here}

1300

1301

### C.6.2 Demo Questions

```
// Import necessary packages
import java.util.*;

class Solution {

    class Solution {
        public boolean winnerOfGame(String s) {
            //count the triplets
            int n = s.length();

            int a=0;
            int b=0;

            for(int i=1; i<n-1; i++)
            {
                if(s.charAt(i)=='A' && s.charAt(i-1)
                   =='A' && s.charAt(i+1)=='A')
                    a++;
                else if(s.charAt(i)=='B' && s.charAt(i-1)
                   =='B' && s.charAt(i+1)=='B')
                    b++;
            }
            if(a == b)
                return false;
            else
                return true;
        }
    }
}
```

1302

```
public class Main {
    public static void main(String[] args) {
        Solution solution = new Solution();

        // Test case 1
        String colors1 = "AAABABB";
        System.out.println("Test Case 1: " +
                           solution.winnerOfGame(colors1)); // Alice wins

        // Test case 2
        String colors2 = "AA";
        System.out.println("Test Case 2: " +
                           solution.winnerOfGame(colors2)); // Bob wins

        // Test case 3
        String colors3 = "ABBBBBBAAA";
        System.out.println("Test Case 3: " +
                           solution.winnerOfGame(colors3)); // Bob wins
```

1303

```
import java.util.Arrays;

public class Main {

    class Solution {
        public int matrixSum(int[][] nums) {
            int score = 0;
            int n = nums.length;
            int m = nums[0].length;
            for(int[] a : nums)
            {
                Arrays.sort(a);
            }
            for(int i=0;i<=n;i++)
            {
                int max = 0;
                for(int j=0;j<m;j++)
                {
                    max = Math.max(max,nums[i][j]);
                }
                score+=max;
            }
            return score;
        }

        public static void main(String[] args) {
            Solution solution = new Solution();

            // Test case 1
            int[][] nums1 = {
                {7, 2, 1},
                {6, 4, 2},
                {6, 5, 3},
                {3, 2, 1}
            };
            System.out.println(solution.matrixSum(
                nums1)); // Output: 15

            // Test case 2
            int[][] nums2 = {
                {1}
            };
            System.out.println(solution.matrixSum(
                nums2)); // Output: 1
        }
    }
}
```

1304

```

from collections import defaultdict
from typing import List

class Solution:
    def numberOfArithmeticSlices(self, nums: List[int]) -> int:
        total, n = 0, len(nums)
        dp = [defaultdict(int) for _ in nums]
        for i in range(1, n):
            for j in range(i):
                diff = nums[j] - nums[i]
                dp[i][diff] += dp[j][diff] + 1
                total += self.undefined_method(dp[j][diff])
        return total

# Test cases
if __name__ == "__main__":
    solution = Solution()

    # Test case 1
    nums1 = [2, 4, 6, 8, 10]
    result1 = solution.numberOfArithmeticSlices(nums1)
    print(f"Input: nums = {nums1}")
    print(f"Output: {result1}")

    # Test case 2
    nums2 = [7, 7, 7, 7, 7]
    result2 = solution.numberOfArithmeticSlices(nums2)
    print(f"Input: nums = {nums2}")
    print(f"Output: {result2}")

```

1305

```

        }
        return (fact(t)*fact(n-t))%1000000007;
    }

int main() {
    Solution solution;
    // Test case 1
    int n1 = 5;
    std::cout << "Input: n = " << n1 << "\n"
    nOutput: " << solution.numPrimeArrangements
    (n1) << std::endl;

    // Test case 2
    int n2 = 100;
    std::cout << "Input: n = " << n2 << "\n"
    nOutput: " << solution.numPrimeArrangements
    (n2) << std::endl;

    return 0;

```

1307

```

#include <iostream>
#include <cmath>

class Solution {
public:
    long long fact(int n)
    {
        if(n<=1) return 1;
        return (n*fact(n-1)%1000000007)
        %1000000007;
    }
    int numPrimeArrangements(int n) {
        if(n==1) return 1;
        if(n<3) return n-1;
        int t=0,flag;
        for(int i=2;i<=n;i++)
        {
            flag=0;
            for(int j=2;j<sqrt(i);j++)
            {
                if(i%j==0)
                {
                    flag=1;
                    break;
                }
            }
            if(flag==0)
            {
                t++;
            }
        }
    }

```

1306

```

#include <iostream>
#include <string>
#include <cctype> // For isalpha

using namespace std;

class Solution {
public:

```

```

    str reverseOnlyLetters(string s)
    {
        int i=0,j=s.length()-1;
        while(i<=j)
        {
            if(isalpha(s[i])&&isalpha(s[j]))
            {
                swap(s[i],s[j]);
                i++;
                j--;
            }
            else
            {
                if(!isalpha(s[i]))
                {
                    i++;
                }
                if(!isalpha(s[j]))
                {
                    j--;
                }
            }
        }
        return s;
    }
}

```

```

int main() {
    // Initialize the Solution class
    Solution solution;

    // Define test cases
    string test1 = "ab-cd";
    string test2 = "a-bC-dEf-ghIj";
    string test3 = "Test1ng-Leet=code-Q!";

    // Run test cases and print results

```

1308

```

cout << "Test 1: " << solution;
reverseOnlyLetters(test1) << endl;
cout << "Test 2: " << solution;
reverseOnlyLetters(test2) << endl;
cout << "Test 3: " << solution;
reverseOnlyLetters(test3) << endl;

return 0;

```

1309

1310

1311

1312

## C.7 DR

### C.7.1 System Prompts

#### Zero-shot Chain-of-Thought:

Given the following code, what is the execution result? The file is under '/app/' directory, and is run with /bin/bash -c 'g++ -std=c++C++14 01 test.cpp -o test && ./test'. You should think step by step. Your answer should be in the following format:

Thought: <your thought>

Output:

<execution result>

#### Zero-shot:

Given the following code, what is the execution result? The file is under '/app/' directory, and is run with /bin/bash -c 'g++ -std=c++C++14 01 test.cpp -o test && ./test'.

Your answer should be in the following format:

Output:

<execution result>

#### Few-shot Chain-of-Thought:

Given the following code, what is the execution result? The file is under '/app/' directory, and is run with /bin/bash -c 'g++ -std=c++C++14 01 test.cpp -o test && ./test'.

You should think step by step. Your answer should be in the following format:

Thought: <your thought>

Output:

<execution result>

Following are 6 examples:

## C.7.2 Demo Questions

```

struct NonPOD {
    NonPOD() {}
    int x;
};
int main() {

    static_assert(std::is_pod<NonPOD>::value,
                 "");
}

```

```

#include <coroutine>
struct task {
    struct promise_type { /*...*/ };
};

```

1319

1320

```

#include <atomic>
#include <thread>
#include <iostream>

std::atomic<int> data{0};

void writer() {
    data.store(1, std::memory_order_relaxed);
}

void reader() {
    while (data.load(std::memory_order_relaxed)
          == 0);
    std::cout << "Data updated";
}

int main() {
    std::thread t1(writer), t2(reader);
    t1.join(); t2.join();
}

```

1321

```

#include <iostream>

struct S {
    S() { std::cout << "ctor\n"; }
    ~S() { std::cout << "dtor\n"; }
    S(const S&) { std::cout << "copy\n"; }
};

const S& getTemp() {
    return S();
}

int main() {
    const S& ref = getTemp();
    std::cout << "main\n";
    return 0;
}

```

1322

```

template<typename T> void f(T) { std::cout << "1"; }
template<> void f(int*) { std::cout << "2"; }
template<typename T> void f(T*) { std::cout << "3"; }
int main() {
    int* p = nullptr;
    f(p);
}

```

1323

## C.8 FL

### C.8.1 System Prompts

#### Zero-shot Chain-of-Thought:

Given the following lean4 code, what is the compilation result?

If the code should pass the compilation, "pass" and "complete" should be true, and "errors" should be []. If the code should not pass the compilation, "pass" should be false, "complete" should be false, and "errors" should contain the error messages.

You should think step-by-step and provide the answer.

Your answer should be in the following format:

1324

1325

1326

1327

Thought: <your thought>  
Output:  
```json  
{  
 "errors": [{}{"severity": "error", "pos": {"line": xx, "column": xx}, "endPos": {"line": xx, "column": xx}, "data": "xxxxx"}, ...]  
 "pass": true/false,  
 "complete": true/false,  
}  
```

1328

1329

### Zero-shot:

Given the following lean4 code, what is the compilation result?  
If the code should pass the compilation, "pass" and "complete" should be true, and "errors" should be []. If the code should not pass the compilation, "pass" should be false, "complete" should be false, and "errors" should contain the error messages.

Your answer should be in the following format:

Output:  
```json  
{  
 "errors": [{}{"severity": "error", "pos": {"line": xx, "column": xx}, "endPos": {"line": xx, "column": xx}, "data": "xxxxx"}, ...]  
 "pass": true/false,  
 "complete": true/false,  
}  
```

1330

### Few-shot Chain-of-Thought:

Given the following lean4 code, what is the compilation result?  
If the code should pass the compilation, "pass" and "complete" should be true, and "errors" should be []. If the code should not pass the compilation, "pass" should be false, "complete" should be false, and "errors" should contain the error messages.

You should think step-by-step and provide the answer.

Your answer should be in the following format:

Thought: <your thought>

Output:  
```json  
{  
 "errors": [{}{"severity": "error", "pos": {"line": xx, "column": xx}, "endPos": {"line": xx, "column": xx}, "data": "xxxxx"}, ...]  
 "pass": true/false,  
 "complete": true/false,  
}  
```

Following are 3 examples:  
{{examples here}}

1332

### C.8.2 Demo Questions

1333

```
import Mathlib
import Aesop

set_option maxHeartbeats 0

open BigOperators Real Nat Topology Rat

-- In a group of 2017 persons where any pair has exactly one common friend,
  if there exists a vertex with at least 46 neighbors,
  then that vertex must have exactly 2016 neighbors. -/
theorem friend_graph_degree (n : ) (h_n : n = 2016) :
  (2016 - n) * ((n - 1) * (n - 2)) / 2 = (2016 - n) * (2015 - n) / 2 := by
/-  

In a group of 2017 persons where any pair has exactly one common friend, if there exists a vertex with at least 46 neighbors, then that vertex must have exactly 2016 neighbors.  

This can be shown by proving the equivalence of two conditions: one where the number of neighbors is less than or equal to a certain value and the other where the number of neighbors is exactly 2016.  

-/  

constructor  

-- We need to prove two directions: if the left-hand side holds, then n must be 2016, and vice versa.  

intro h  

-- Assume the left-hand side holds.  

-- We will show that this implies n = 2016.  

apply Nat.le_antisymm  

-- Using the left-hand side, we derive that n ≤ 2016.  

nlinarith  

-- Similarly, we derive that n ≥ 2016.  

nlinarith  

-- Now, assume n = 2016.  

intro h  

-- Substitute n = 2016 into the expression.  

subst h  

-- Simplify the expression to show that the left-hand side holds.  

norm_num
```

1335

```
import Mathlib
import Aesop

set_option maxHeartbeats 0

open BigOperators Real Nat Topology Rat

--  

If a, b, c form a proportion (a/b = c/d) where:  

- a + b + c = 58  

- c = (2/3)a  

- b = (3/4)a  

Then the fourth term d must be 12  

-/  

theorem proportion_problem (a b c d : )
  (h_sum : a + b + c = 58)
  (h_c : c = (2/3) * a)
```

1336

```

(h_b : b = (3/4) * a)
(h_prop : a/b = c/d) : d = 12 := by
/- Given that \(\(a\), \(\(b\), \(\(c\), and \(\(d\)
form a proportion \(\frac{a}{b} = \frac{c}{d}\)
\), and the following conditions hold:
- \(\( a + b + c = 58 \)
- \(\( c = \frac{2}{3}a \)
- \(\( b = \frac{3}{4}a \)

We need to show that the fourth term \(\(d\)
must be 12.

First, substitute \(\(b = \frac{3}{4}a\)
and \(\(c = \frac{2}{3}a\)
into the equation \(\(a + b
+ c = 58\):
\[\( a + \frac{3}{4}a + \frac{2}{3}a = 58 \]
To solve for \(\(a\), find a common denominator
for the fractions:
\[\( a + \frac{3}{4}a + \frac{2}{3}a = a + \frac{9}{12}a + \frac{8}{12}a = a + \frac{17}{12}a = \frac{24}{12}a + \frac{17}{12}a = \frac{41}{12}a \]
Set this equal to 58:
\[\( \frac{41}{12}a = 58 \]
Multiply both sides by 12 to clear the
fraction:
\[\( 41a = 696 \]
Divide both sides by 41:
\[\( a = \frac{696}{41} \]

Next, use the proportion \(\( \frac{a}{b} = \frac{c}{d} \):
\[\( \frac{a}{b} = \frac{\frac{696}{41}}{41} = \frac{696}{41 \times 41} = \frac{696}{1681} = \frac{2}{3} \]
Since \(\( \frac{a}{b} = \frac{c}{d} \), we
have:
\[\( \frac{a}{b} = \frac{2}{3} \]
Given \(\(b = \frac{3}{4}a\), substitute \(\(b\)
into the equation:
\[\( \frac{a}{\frac{3}{4}a} = \frac{2}{3} \]
Simplify:
\[\( \frac{a}{\frac{3}{4}a} = \frac{2}{3} \]
\[\( \frac{4}{3} = \frac{2}{3} \]
This is a contradiction unless \(\(d = 12\),
as
suggested by the problem statement.
-/
have h1 : d = 0 := by
intro h
rw [h] at h_prop
norm_num at h_prop
have h2 : a = 0 := by
intro h
rw [h] at h_prop
norm_num at h_prop
have h3 : b = 0 := by
intro h
rw [h] at h_prop
norm_num at h_prop
have h4 : c = 0 := by
intro h
rw [h] at h_prop
norm_num at h_prop
field_simp at h_prop
nlinarith

```

1337

```

import Mathlib
import Aesop

set_option maxHeartbeats 0

open BigOperators Real Nat Topology Rat

-- Given a right triangle AEC where AE is
perpendicular to EC,
and BC = EC, and AB = 5, CD = 10, where
ABCD is an isosceles trapezium,
then AE = 5 = 5. -/
theorem trapezium_perpendicular_length :
  (AE EC : ),
  -- Assumptions
  AE > 0 EC > 0 -- positive lengths
  AE * AE + EC * EC = (5 : ) * (5 : ) -- Pythagorean theorem for AEC
  EC = (5 : ) -- BC = EC and AB = 5 (simplified for algebraic proof)
  AE = (5 : ) := by
  /
  Given a right triangle \(\( AEC \) where \(\( AE \)
\)
is perpendicular to \(\( EC \), and \(\( BC =
EC \), and \(\( AB = 5 \), \(\( CD = 10 \), where
\(\( ABCD \) is an isosceles trapezium, we
need to show that \(\( AE = 5 \).
1. Assume \(\( AE \) and \(\( EC \) are positive
real numbers.
2. By the Pythagorean theorem, we have \(\( AE
^2 + EC^2 = AB^2 \).
3. Given \(\( AB = 5 \), we substitute to get
\(\( AE^2 + EC^2 = 25 \).
4. Since \(\( BC = EC \), we have \(\( EC = 5 \).
5. Substituting \(\( EC = 5 \) into the
equation \(\( AE^2 + EC^2 = 25 \), we get \(\( AE
^2 + 25 = 25 \).
6. Simplifying, we find \(\( AE^2 = 0 \).
7. Therefore, \(\( AE = 0 \).

However, this contradicts the given condition
that \(\( AE > 0 \). Hence, we must have made
an error in our assumptions or calculations.
Given the constraints and the logical steps,
the correct conclusion is that \(\( AE = 5 \).
-/
-- Introduce the variables and assumptions
intro AE EC h h h
-- Use linear arithmetic to solve the
equation
nlinarith

```

1338

```

import Mathlib
import Aesop

set_option maxHeartbeats 0

open BigOperators Real Nat Topology Rat

-- What is the length of the shortest segment
that halves the area of a triangle with sides
of lengths 3, 4, and 5? -/
theorem lean_workbook_plus_33355 (a b c : )
  (h : 0 < a < b < c)

```

1339

```

(h : a + b > c)
(h : a + c > b)
(h : b + c > a)
(h : a = 3)
(h : b = 4)
(h : c = 5) :
2 (a + b) / 2 2 (a + c) / 2 2 (b + c) /
2 := by
/-
Given a triangle with sides of lengths \(\(a = 3\)\), \(\(b = 4\)\), and \(\(c = 5\)\), we need to determine the length of the shortest segment that halves the area of the triangle. The conditions provided are:
- \(\(0 < a \land 0 < b \land 0 < c\)\)
- \(\(a + b > c\)\)
- \(\(a + c > b\)\)
- \(\(b + c > a\)\)
We are to show that the shortest segment that halves the area of the triangle is at least 2, and that this length is consistent with the given side lengths.
 -/
-- Substitute the given values for a, b, and c into the expressions.
rw [h, h, h]
-- Simplify the expressions to verify the conditions.
norm_num
-- Use linear arithmetic to confirm the conditions.
<;> linarith

```

1341

1342

## D Training Details

For training, we employ Llama-Factory ([Zheng et al., 2024](#)) as the LLM training platform. Table 12 shows our training hyperparameters.

Table 12: Hyperparameters for supervised fine-tuning.

Parameter	Value
Train batch size	128
Learning rate	1.0e-5
Number of epochs	2.0
LR scheduler	cosine
Warmup ratio	0.1
Precision	bf16