

LIVECODEBENCH: HOLISTIC AND CONTAMINATION FREE EVALUATION OF LARGE LANGUAGE MODELS FOR CODE

Anonymous authors

Paper under double-blind review

ABSTRACT

Large Language Models (LLMs) applied to code-related applications have emerged as a prominent field, attracting significant interest from academia and industry. However, as new and improved LLMs are developed, existing evaluation benchmarks (e.g., HUMANEVAL, MBPP) are no longer sufficient for assessing their capabilities suffering from data contamination, overfitting, saturation, and focus on merely code generation. In this work, we propose LIVECODEBENCH, a comprehensive and contamination-free evaluation of LLMs for code, which collects *new* problems over time from contests across three competition platforms, LEETCODE, ATCODER, and CODEFORCES. Notably, our benchmark also focuses on a broader range of code-related capabilities, such as self-repair, code execution, and test output prediction, beyond just code generation. Currently, LIVECODEBENCH hosts over six hundred coding problems that were published between May 2023 and Aug 2024. We evaluate over 50 LLMs on LIVECODEBENCH (LCB for brevity) presenting the largest evaluation study of code LLMs on competition problems. Based on the study, we present novel empirical findings on contamination, overfitting, and holistic evaluations. We demonstrate that time-segmented evaluations serve as a robust approach to evade contamination; they are successful at detecting contamination across a wide range of open and closed models including GPT-4-O, CLAUDE, DEEPSEEK, and CODESTRAL. Next, we highlight overfitting and saturation of traditional coding benchmarks like HUMANEVAL and demonstrate LCB allows more reliable evaluations. Finally, our holistic evaluation scenarios allow for measuring the different capabilities of programming agents in isolation.

1 INTRODUCTION

Code has emerged as an important application area for LLMs, with a proliferation of code-specific models (Chen et al., 2021; Austin et al., 2021; Nijkamp et al., 2022; Li et al., 2023b; Roziere et al., 2023; Guo et al., 2024; Ridnik et al., 2024; Lozhkov et al., 2024) and their applications across various tasks such as program repair (Zheng et al., 2024; Olausson et al., 2023), optimization (Madaan et al., 2023a), test generation (Steenhoek et al., 2023), documentation (Luo et al., 2024), SQL (Sun et al., 2023). In contrast with these rapid advancements, evaluations have remained relatively stagnant, and current benchmarks like HUMANEVAL and MBPP paint a misleading picture. Firstly, while coding is a multi-faceted skill, these benchmarks only focus on code generation, thus overlooking broader code-related capabilities. Moreover, these benchmarks suffer from contamination or overfitting, as benchmark samples are present in the training datasets.

Motivated by these shortcomings, we introduce LIVECODEBENCH, built on the following principles:

1. **Live updates to prevent contamination.** LLMs are trained on massive inscrutable corpora and current benchmarks suffer from the risk of data contamination as they could be included in those training datasets. While previous works have attempted decontamination using both exact and fuzzy matches (Li et al., 2023b;d), it can be a non-trivial task (Team,

¹Note that for model comparisons – performances are averaged across multiple months and platforms achieving a larger sample size.

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

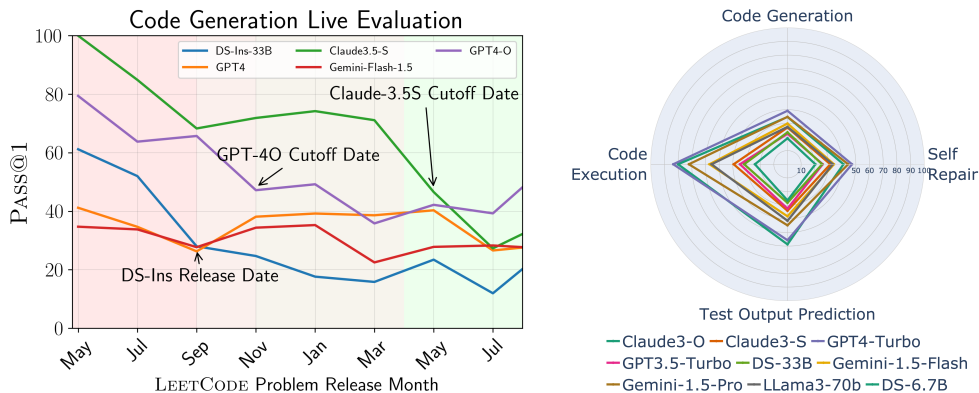


Figure 1: **Left.** LIVECODEBENCH comprises problems marked with release dates, allowing evaluations over different time windows. The figure depicts performances of DEEPSEEK, GPT-4-O, and CLAUDE-3s models over bimonthly time windows (~ 40 LEETCODE problems) showcasing a stark drop after their cutoff dates, highlighting contamination. Thus, we can detect and avoid contamination by evaluating models only on time-windows after the model’s cutoff date (green region)¹. **Right.** We evaluate LLMs across four scenarios that capture key coding capabilities necessary for building programming agents: code generation, repair, testing, and comprehension. The figure depicts various model performances across the scenarios available in LIVECODEBENCH in a radial plot – highlighting relative differences changing across the scenarios.

2024; Yang et al., 2023). Existing competition programming benchmarks (like APPS and CODESCOPE) are already contaminated and may fail to provide reliable evaluation of code LLM capabilities. In LIVECODEBENCH, to prevent the risk of problem contamination, we evaluate models on *new* problems tagged with a *release date*. Next, for newer models, we only consider problems released after the model’s cutoff date to ensure that the model has not been trained on it as demonstrated in Figure 1 left.

- Holistic Evaluation.** Current code evaluations primarily focus on natural language to code generation. However, programming is a multi-faceted task that requires capabilities beyond those measured by code generation. These broader capabilities are even more relevant for constructing programming agents that can interact with the execution environment. Therefore, we evaluate the execution-feedback-based multi-turn coding using the self-repair scenario (Olausson et al., 2023), assess code comprehension capabilities using the code execution scenario (Gu et al., 2024), and introduce the test output prediction scenario to evaluate the models’ test generation capabilities.
- High-quality problems and tests.** High-quality problems and tests are crucial for reliable evaluation of LLMs. However, existing benchmarks suffer from multiple deficiencies. First, existing competitive programming benchmarks (APPS, CODE-CONTESTS, XCODEEVAL, CODESCOPE) contain problems not amenable to input-output-based auto-grading. For example, CODE-CONTESTS (Li et al., 2022) (page 39, second-to-last paragraph) reports that about twenty-five percent of the problems in the benchmark accept multiple correct outputs for a single input. This *incorrectly penalizes correct solutions*, adding noise to a considerable fraction of the benchmark. Next, prior benchmarks like HUMANEVAL and APPS contain insufficient tests, further exacerbating noise in evaluations (Liu et al. (2023a) reports 8% drop in model performance). In LIVECODEBENCH, we source high-quality problems from reputable platforms and implement heuristics to detect and remove problems not amenable to input-output-based auto-grading. Finally, for every problem, we provide a substantial number of tests (over 18 on average) for reliable and efficient evaluations. In contrast to prior works, we also include several large tests designed for stress-testing solutions ensuring weak or worse complexity solutions do not pass test harnesses.
- Difficulty Guided Problem Curation.** Competitive programming is a challenging domain even for strong LLMs. As a result, these problems can be unsuitable for meaningfully comparing today’s open LLMs, because the variance in performance is low, often relying on less than 1% performance differences (within the margin of error). Therefore, we use problem difficulty ratings (sourced from the competition websites) to curate our problems, avoiding

Table 1: Comparing LCB with existing coding and competition programming benchmarks.

| Benchmark | Contamination Prevention | Problem Curation | Robust Test Based Eval. | Varied Difficulties | Not Saturated | Broader Eval. Scenarios | Comp. Analysis across Models |
|---------------------------------|--------------------------|------------------|-------------------------|---------------------|---------------|-------------------------|------------------------------|
| HUMANEVAL (Chen et al., 2021) | ✗ | ✓ | ✗ | ✓ | ✗ | ✗ | ✓ |
| HUMANEVAL+ (Liu et al., 2023b) | ✗ | ✓ | ✓ | ✓ | ✗ | ✗ | ✓ |
| MBPP (Austin et al., 2021) | ✗ | ✓ | ✗ | ✓ | ✗ | ✗ | ✓ |
| AAPP (Hendrycks et al., 2021) | ✗ | ✗ | ✗ | ✓ | ✓ | ✗ | ✗ |
| CODE-CONTESTS (Li et al., 2022) | ✗ | ✗ | ✓ | ✗ | ✓ | ✗ | ✗ |
| xCODEEVAL (Khan et al., 2023) | ✗ | ✗ | - | ✗ | ✓ | ✓ ² | ✗ |
| CODESCOPE (Yan et al., 2023) | ✗ | ✗ | - | ✗ | ✓ | ✓ ² | ✗ |
| TACO (Li et al., 2023c) | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ |
| USCAOBENCH (Shi et al., 2024) | ✗ | ✓ | ✓ | ✓ | ✓ | ✗ | ✗ |
| LIVECODEBENCH (Ours) | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |

those that are too difficult for current models. In particular, we avoid collecting CODEFORCES problems used by prior works (CODE-CONTESTS, CODESCOPE, xCODEEVAL) since we find they do not sufficiently distinguish models due to model performances (PASS@1) tending to zero. Indeed, LCB easy and medium problems can separate 7B models that are indistinguishable on the hard subset (see DS-7B vs SC2-7B).

With these principles in mind, we build LIVECODEBENCH, a continuously updated benchmark that avoids data contamination. Particularly, we collect 612 problems from contests occurring between May 2023 and Aug 2024 and use them to construct the different scenarios.

Empirical Findings. We have evaluated over 50 (18 base models and 34 instruction-tuned) models across different LIVECODEBENCH scenarios. Based on this study and following analysis, we present novel empirical findings which have not been revealed in prior benchmarks.

- Contamination.** We observe a stark drop in the performance of DEEPSEEK, GPT-4-O, CODESTRAL and CLAUDE-3S on LEETCODE problems released after Aug 2023, Oct 2023, Jan 2024, and April 2024 (Figure 1 left). These results highlight likely contamination in older problems and time-segmented evaluations prove effective for performing fair comparisons across varied set of models from both open and closed domains.
- Holistic Evaluation.** Our evaluations reveal that model performances are correlated across tasks, but the relative differences do vary. For example, in Figure 1 right, the gap between open and closed models further increases on tasks like self-repair or test output prediction. Similarly, CLAUDE-3-OPUS surpasses GPT-4 on the test output prediction scenario.
- HUMANEVAL Overfitting.** Upon comparing LIVECODEBENCH with HUMANEVAL, we find that models cluster into two groups, ones that perform well on both benchmarks, and others that perform well on HUMANEVAL but not on LIVECODEBENCH (see Figure 4). The latter group primarily comprises fine-tuned open-access models while the former comprises base and closed models. This indicates some level of overfitting to HUMANEVAL.
- Model Comparisons** (Figure 3). We provide model comparisons across different groups like base, instruct, open, and closed models and across groups like open vs closed models.

Concurrent Work. (Huang et al., 2023) also evaluate LLMs in a time-segmented manner. However, they only focus on CODEFORCES problems, while we combine problems across platforms and additionally propose a holistic evaluation across multiple code-related scenarios. Liu et al. (2024) evaluate the code comprehension capabilities of LLMs using execution. (Singhal et al., 2024) also evaluates LLMs on more tasks but focus on *non-functional-correctness aspects* of programming. Shi et al. (2024) evaluate LLMs on USACO problems while we focus on simple competition scenarios.

2 HOLISTIC EVALUATION

Code capabilities of LLMs are evaluated and compared using natural language to code generation tasks. However, this only captures one dimension of code-related capabilities. AlphaCodium (Ridnik et al., 2024) developed an intricate LLM pipeline for solving competition coding problems. By combining natural language reasoning, test case generation, code generation, and self-repair, they

²LIVECODEBENCH considers different evaluation scenarios compared to xCODEEVAL and CODESCOPE focusing on settings where reliable evaluation is possible.

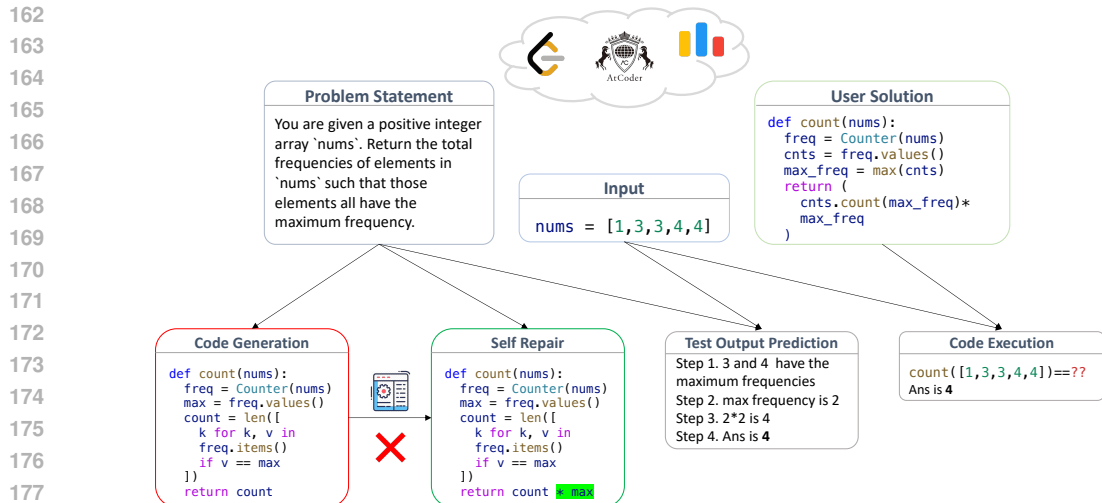


Figure 2: Overview of the four scenarios present in LIVECODEBENCH.

achieve significant improvements over a naive direct code generation baseline, showcasing the importance of these broader capabilities. Motivated by this, we propose a holistic evaluation of LLMs using a suite of evaluation setups that capture a broader range of code-related capabilities.

Specifically, we evaluate code LLMs in four scenarios, namely code generation, self-repair, code execution, and test output prediction. Our selection criterion was to pick settings that are useful components in code LLM workflows and in addition, have clear and automated evaluation metrics.

Following we describe each of these scenarios in detail.

Code Generation. The code generation scenario follows the standard setup for generating code from natural language. The model is given a problem statement, which includes a natural language description and example tests (input-output pairs), and is tasked with generating a correct solution. The evaluation is performed based on functional correctness, using a set of *unseen* test cases. We use the PASS@1 metric measured as the fraction of the problems for which the model was able to generate a program passing all tests. Figure 2 (left) provides an example of this scenario.

Self Repair. The self-repair scenario is based on previous works that tested the self-repair capabilities of LLMs (Olausson et al., 2023; Shinn et al., 2023; Chen et al., 2023). Here, the model is given a problem statement from which it generates a candidate program (similar to the single-step code generation scenario above). However, in case of a mistake, the model is additionally provided with error feedback (either the exception message or a failing test case) and is tasked with generating the fixed solution. Similar to the code generation scenario, the evaluation is performed via functional correctness on the final program, i.e. either the single-step correct generation or the attempted repair. We use the PASS@1 metric to measure the combined performance after the repair step. Figure 2 (mid-left) provides an example of this scenario.

Code Execution. The code execution scenario is based on the output prediction setup used in CRUXEVAL (Gu et al., 2024). The model is provided a program snippet consisting of a function (\mathcal{F}) along with a test input to the program and is tasked with predicting the output of the program on the input test case. The evaluation is performed via an execution based on an exact match correctness metric. Figure 2 (right) provides an example of the code execution scenario.

Test Case Output Prediction. Finally, we introduce a new task that is designed to study natural language reasoning and test generation. In this task, the model is given the problem statement along with a test case input, and it is tasked with generating the expected output for that input. This task follows a setup similar to the one used in CODET (Chen et al., 2022), where tests are generated solely from problem statements, without the need for the function’s implementation. A key difference is that we provide a fixed set of test inputs for each problem in our dataset, and the models are then prompted to only predict the expected output for those specific inputs. This approach allows for a straightforward evaluation of the test generation capabilities by avoiding test input prediction, a hard-to-evaluate task. Figure 2 (mid-right) provides an example of this scenario.

Finally, we would like to point out that LIVECODEBENCH also offers an extensible framework to add new scenarios in the future. So other relevant settings like input generation, program summarization, optimization, etc. can be integrated with our setup.

3 BENCHMARK CURATION

We curate our problems from three coding competition websites: LEETCODE, ATCODER, and CODEFORCES. These websites periodically host contests containing problems that assess the coding and problem-solving skills of participants. The problems consist of a natural language problem statement along with example input-output examples, and the goal is to write a program that passes a set of hidden tests. Further, thousands of participants participate, solving these problems thus ensuring that the problems are vetted for clarity.

3.1 DATA COLLECTION

We have written automated HTML scrapers for each of the above websites to collect problems and the corresponding metadata. To ensure quality and consistency, we parse mathematical formulas and exclude problems with images. We also exclude problems that are not suitable for grading by input-output examples, such as those that accept multiple correct answers or require the construction of data structures. Specifically, we use keyword-based heuristic filters to filter interactive problems and problems accepting multiple correct solutions. Besides parsing the problem descriptions, we also collect associated ground truth solutions and test cases whenever directly available. Thus for each problem, we collect tuples of natural language problem statement P , test cases T , and ground truth solution S . Finally, we associate the contest date D to mark the release date of each problem and use the collected attributes to construct problems for our four scenarios (detailed in Section C.2 ahead). Note that this process is completely automated and human involvement is only involved in modifying high-level design decisions such as updating problem difficulty settings and improving the keyword-based heuristics for filtering problems over different “live updates”.

Scrolling through time. As noted, we associate the contest date D for each problem. The release date allows us the measure the performance of LLMs over different time windows by filtering problems based on whether the problem release date falls within a time window (referred to as “scrolling” through time). This is crucial for evaluating and comparing models trained at different times. Specifically, for a new model and the corresponding cutoff date (normalized to the release date if the training cutoff date is not published), we can measure the performance of the model on benchmark problems released after the cutoff date. We have developed a UI that allows comparing models on problems released during different time windows (Figure 9).

Test collection. Tests are crucial for assessing the correctness of the generated outputs and are used in all four scenarios. We collect tests available on platform websites whenever possible and use them for the benchmark. Otherwise, following (Liu et al., 2023b), we use a LLM (here GPT-4-TURBO) to generate test *inputs* for the problems. A **key difference** between our test generation approach is that instead of generating inputs directly using the LLM, we construct generators that sample inputs based on the problem specifications using in-context learning. Details and examples of such input generators can be found in Section A.2. Finally, we collect a small fraction of failing tests from the platforms for adversarial tests.

Problem difficulty. Competition programming has remained a challenge for LLMs, with GPT-4 achieving an average CODEFORCES rating (ELO) of 392, placing it in the bottom 5 percentile (OpenAI, 2023). This makes it difficult to compare LLMs, as the variation in performance across models is low. In LIVECODEBENCH, we collect problems of diverse difficulties as labeled in competition platforms, excluding problems that are rated above a certain threshold that are likely too difficult for even the best models³. Further, we use these ratings to classify problems as EASY, MEDIUM, and HARD for more granular model comparisons.

We defer the platform and scenario specific curation details to the Appendix (Section C)

³From our early explorations, we find CODEFORCES problems being considerably more difficult than ATCODER and LEETCODE problems and thus focus primarily on the latter platforms.

4 EXPERIMENT SETUP

We describe the experimental setup in this section. First, we provide the common setup across the scenarios, followed by the scenario-specific setups in Section 4.1.

Models. We evaluate 52 models across various sizes, ranging from 1.3B to 70B, including base models, instruction models, and both open and closed models. Our experiments include models from different classes, such as GPTs, CLAUDES, GEMINIS, MISTRAL, LLAMA-3S, DEEPSEEKS, CODELLA-MAS, STARCODER2, CODEQWEN and their variants. Appendix D.1 provides the list of models.

Evaluation Metrics. We use the PASS@1 (Kulal et al., 2019; Chen et al., 2021) metric for our evaluations. Specifically, we generate 10 candidate answers for each problem from model providers or using vLLM (Kwon et al., 2023). We use nucleus sampling with temperature 0.2 and top_p 0.95 and calculate the fraction of correct solutions.

4.1 SCENARIO-SPECIFIC SETUP

The setup for each scenario is presented below. Note that the base models are only used in the code generation scenario since they do not easily follow the format for the other scenarios.

Code Generation. For the instruction-tuned models, we use a zero-shot prompt and follow the approach of Hendrycks et al. (2021) by adding appropriate instructions to generate solutions in either functional or stdin format (one-shot for base-models). Section D.2 depicts the prompt.

Self Repair. Similar to prior work Olausson et al. (2023), we use the programs generated during the code generation scenario along with the corresponding error feedback to build the zero-shot prompt for the self-repair scenario. The error feedback can be syntax errors, runtime errors, wrong answers, and time-limit errors. Section D.3 provides error feedback and the corresponding prompt.

Code Execution. We use few-shot prompts for the code execution scenario, both with and without chain-of-thought prompting (COT). Particularly, we use a 2-shot prompt without COT and a 1-shot prompt with COT with manually detailed steps. The prompts are detailed in Section D.4.

Test Output Prediction. We use a zero-shot prompt that queries the model to complete assertions, given the problem, function signature, and test input. We provide the prompt in Section D.5.

5 RESULTS

We first describe how LIVECODEBENCH helps detect and avoid benchmark contamination in Section 5.1. Next, we present the findings from our evaluations on LIVECODEBENCH in Section 5.2.

5.1 AVOIDING CONTAMINATION

Contamination in DEEPSEEK and GPT-4-O. LIVECODEBENCH curates problems released since May 2023. However, DEEPSEEK was released Sep 2023 and GPT-4-O’s official cutoff date is Nov 2023. We can measure the performance of these models on the benchmark problems released after these dates, thereby estimating the performance of the model on previously unseen problems. Figure 1 shows the performance of these models on LEETCODE problems released in different months from May 2023 and Aug 2024. First, we notice stark drops in model performance – DS-INS-33B model after Aug. 2023, GPT-4-O model after Nov. 2023, CLAUDE-3S model after Apr. 2023. These drops align with their release or cutoff dates which suggests that the earlier problems might indeed be contaminated. This trend is consistent across other LIVECODEBENCH scenarios like repair and code execution, as depicted in Figure 10. Concurrently, (Guo et al., 2024) (Section 4.1, last paragraph) also acknowledge the possibility of LEETCODE contamination, noting that “*models achieved higher scores in the LeetCode Contest held in July and August*”.

Interestingly, we find that this drop in performance primarily occurs for the LEETCODE problems only and that the model performance is relatively smooth across the months for problems from other platforms. Figure 11 shows a relatively stable performance for all models on ATCODER problems.

Performances of other models. We study performance variations in other models released more recently. Particularly, GPT-4-TURBO, MISTRAL-L, and CLAUDE-3S models were released in Nov’2023,

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

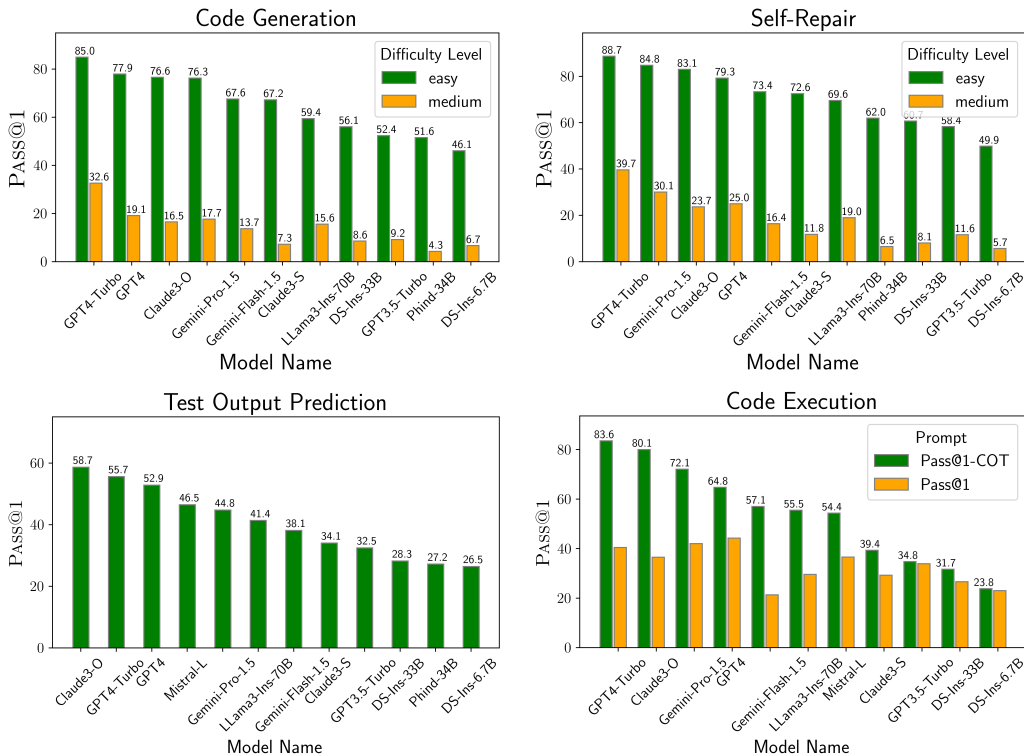


Figure 3: Model performances across the four scenarios available in LIVECODEBENCH .

Feb’2024, and Mar’2024 respectively. Empirically, we do not observe significant performance variations across the months, as shown in Figure 12. Interestingly, we find that even the DS-BASE-33B model also suffers from contamination dropping from PASS@1 ~ 60 in May problems to PASS@1 ~ 0 in September LEETCODE problems. Finally, CODESTRAL achieves PASS@1 36.5 on problems released between May’23 and Jan’24 and PASS@1 28.3 on problems post Jan’24.

5.2 PERFORMANCE AND MODEL COMPARISONS

We provide our model comparison findings here (and in Appendix E). To overcome contamination issues in DEEPSEEK models, we only consider problems released since Sep 2023 for all evaluations below. Figure 3 shows the performance of a subset of models across the four scenarios.

Holistic Evaluations. We have evaluated the models across the four scenarios currently available in LIVECODEBENCH. Figure 1 displays the performance of models on all scenarios along the axes of the polar chart. First, we observe that the relative order of models remains mostly consistent across the scenarios. This is also supported by high correlations between PASS@1 metric across the scenarios – over 0.88 across all pairs as shown in Figure 13. However, despite the strong correlation, the relative differences in performance do vary across the scenarios. For example, GPT-4-TURBO further gains performance gap over GPT-4 in the self-repair scenario after already leading in the code generation scenario. Similarly, CLAUDE-3-OPUS and MISTRAL-L perform well in tasks involving COT, particularly in the code execution and test output prediction scenarios. For instance, CLAUDE-3-OPUS even outperforms GPT-4-TURBO in the test output prediction scenario. These differences highlight the need for holistic evaluations beyond measuring code generation capabilities.

Comparison to HUMANEVAL. Next, we compare how code generation performance metrics translate between LIVECODEBENCH and HUMANEVAL, the primary benchmark used for evaluating coding capabilities. Note that we use HUMANEVAL+ providing more accurate evaluations. Figure 4 shows a scatter plot of PASS@1 on HUMANEVAL+ versus LCB-Easy code generation scenario. We find only a moderate correlation of 0.72, with much larger performance variations on LCB-Easy.

Additionally, we observe that the models cluster into two groups, shaded in red and green. The models in the green-shaded region lie close to the $x = y$ line, indicating that they perform similarly on both

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

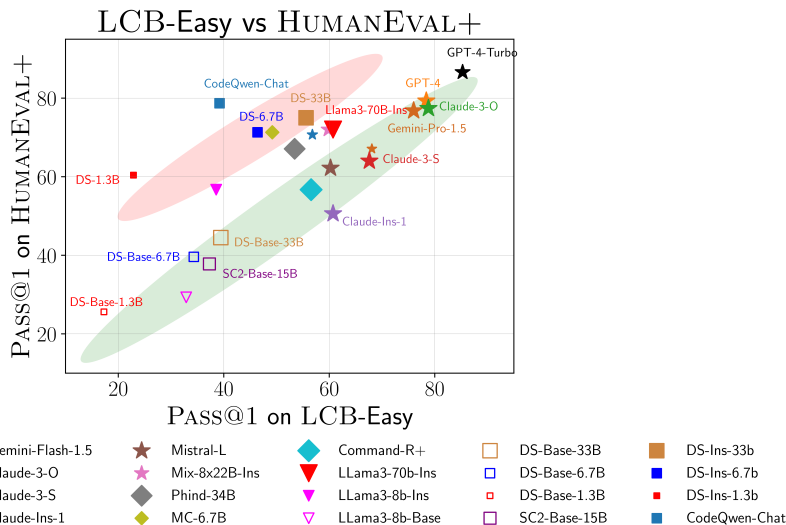


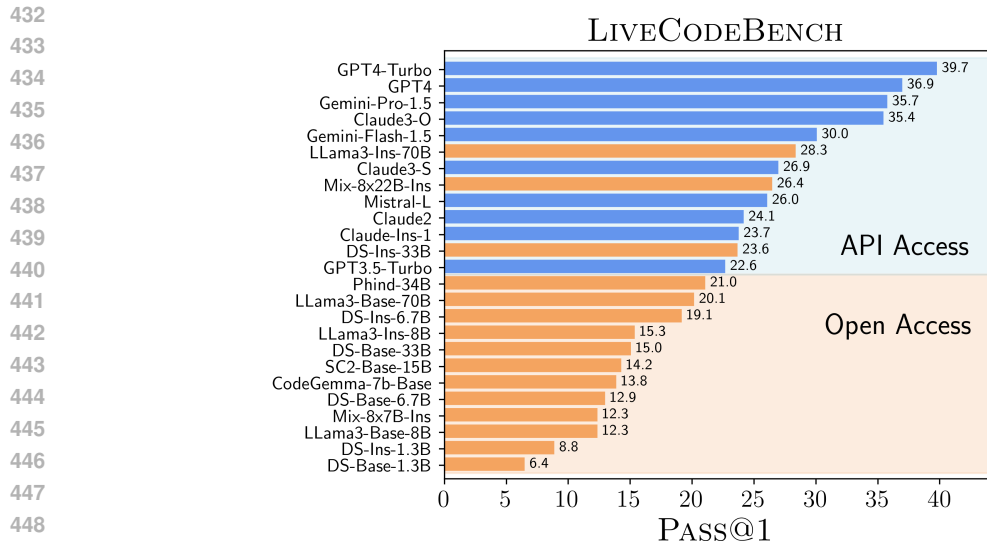
Figure 4: Scatter plot comparing PASS@1 of models on HUMAN EVAL+ versus LCB-Easy (time-window Sep 2023 to May 2024). We find that the models are separated into two groups – the green-shaded region where performances on the two datasets are *aligned* and the red-shaded region where models perform well on HUMAN EVAL+ but perform poorly on LIVECODEBENCH. This indicates potential overfitting on HUMAN EVAL+ and primarily occurs in the fine-tuned variants of open-access models. For example, DS-INS-1.3B achieves PASS@1 of 60 and 26 on HUMAN EVAL+ and LCB-Easy subset. Thus, while it ranks above CMD-R+ (~ 100x larger!) on HUMAN EVAL+, it performs significantly worse on the LCB. Similarly, DS-INS-6.7B and CODEQWEN outperform SONNET on HUMAN EVAL+ but are > 20 points behind on LCB-Easy.

benchmarks. On the other hand, models in the red shaded region perform well on HUMAN EVAL+ but not as well on LIVECODEBENCH. Interestingly, the green-shaded cluster contains base models or closed-access models, while the red-shaded cluster primarily comprises fine-tuned variants of open-access models. The well-separated clusters suggest that many models that perform well on HUMAN EVAL might be overfitting on the benchmark, and their performances do not translate well to problems from other domains or difficulty levels like those present in LIVECODEBENCH. Indeed, HUMAN EVAL has small and isolated programming problems and thus easier to overfit. In contrast, LIVECODEBENCH problems are sourced from reputable coding platforms offering more challenging problems with higher diversity and difficulty levels. We detail instances of this overfitting (DS-INS-1.3B, DS-INS-6.7B, and CODEQWEN) in Figure 4 caption.

Comparing Base Models. We use four families of base models and compare them on the code generation scenario. A one-shot prompt is used for all models to avoid any formatting and answer extraction issues. We find L3-BASE and DS-B models are significantly better than both CODELLAMA and STARCORDER2 base models with a DS-BASE-6.7B model even outperforming both CL-BASE-34B and SC2-BASE-15B models. Note that some LCB specific differences can potentially be attributed to data curation approaches. For instance, SC models (and potentially DS as discussed in Section 5.1) use competition problems during pre-training.

Role of Post Training. Post-training improves performance on both HUMAN EVAL+ and LIVECODEBENCH for the code generation scenario. Particularly, on LCB L3-INS-70B, DS-INS-33B and PHIND-34B improve PASS@1 over their base models by 8.2, 7.3 and 9.5 points respectively. This highlights the importance of good post-training datasets for building strong LLMs. At the same time, we note that the base models have *aligned* performances on LCB code generation and HUMAN EVAL+ benchmarks and lie within or close to the green shaded region in Figure 4. However, the fine-tuned open models exhibit a larger performance gap, with much better performances on HUMAN EVAL+. On the other hand, the closed-access models are still aligned across both benchmarks. This suggests the necessity of more diverse post-training data mixtures.

Open-Access vs Closed-Access Models. In general, closed (API) access model families generally outperform the open access models. The gap is only closed by three models, namely L3-INS-70B, MIXTRAL, and DS-INS-33B which reach the performance levels of the closed models. For instance, in



450 Figure 5: Comparison of open access and (closed) API access models on LIVECODEBENCH code
451 generation scenario. We find that closed-access models consistently outperform the open models
452 with only strong instruction-tuned variants of > 30B models (specifically L3-INS-70B, MIXTRAL and
453 DS-INS-33B models) crossing the performance gap.

455 the code generation scenario (Figure 5), these models reach close to or even outperform closed access
456 models like GEMINI-PRO, GPT-3.5-TURBO, and CLAUDE-3-SONNET. The performances vary across
457 scenarios with the closed-access models faring better in test output and code execution scenarios.

458 **Highlighting gap between SoTA** One distinct observation from our evaluations is the large gap
459 between SoTA models and open models across all scenarios. Particularly, GPT-4-TURBO, GPT-4,
460 GEMINI-PRO-1.5 and CLAUDE-3-OPUS lead across the benchmarks with wide performance margins
461 over other models. This distinguishes LIVECODEBENCH from prior benchmarks (like HUMANEVAL)
462 where various open models have achieved similar or better performance. For example, DS-INS-33B
463 is merely 4.3 point behind GPT-4-TURBO on HUMANEVAL+ but 16.2 points (69%) on LCB code
464 generation scenario. This gap either holds or sometimes even amplifies across other scenarios.

465 **Comparing open-access instruction-tuned models.** Here, we compare various fine-tuned variants
466 of the L3-BASE, DEEPSEEK and CODELLAMA base models across different model sizes. We find
467 that fine-tuned L3-BASE and DEEPSEEK models lead in performance, followed by PHIND-34B and
468 CODELLAMA models across most scenarios. Broadly, we find that model performances correlate with
469 model sizes. For example, PHIND-34B model outperforms the 6.7B models across all scenarios.

471 6 RELATED WORK

472
473 **Language Models for Code Generation.** Starting with Codex (Chen et al., 2021), there are
474 over a dozen code LLMs. These include CodeT5 (Wang et al., 2021; 2023), CodeGen (Nijkamp
475 et al., 2022), SantaCoder (Allal et al., 2023), StarCoder (Li et al., 2023b; Lozhkov et al., 2024),
476 InCoder (Fried et al., 2022), CodeGeeX (Zheng et al., 2023), L3-BASE, DEEPSEEK (Bi et al., 2024)
477 and CODELLAMA (Roziere et al., 2023).

478 **Code Generation Benchmarks.** Many benchmarks have been proposed to compare and evaluate
479 these models. These primarily focus on natural language to Python code generation: HUMANEVAL
480 (Chen et al., 2021), HUMANEVAL+ (Liu et al., 2023b), APPS (Hendrycks et al., 2021), CODE-CONTESTS
481 (Li et al., 2022), MBPP (Austin et al., 2021), L2CEval (Ni et al., 2023). Their variants have been
482 proposed to cover more languages, (Wang et al., 2022a; Zheng et al., 2023; Cassano et al., 2022;
483 Athiwaratkun et al., 2022). Many benchmarks have focused on code generation in APIs. Benchmarks
484 like DS-1000 (Lai et al., 2023), ARCADE (Yin et al., 2022), NumpyEval (Zhang et al., 2023b), and
485 PandasEval (Jain et al., 2022) focus on data science APIs. Other benchmarks measure using broader
APIs or general software engineering tasks, such as JuICe (Agashe et al., 2019), APiBench (Patil

et al., 2023), RepoBench (Liu et al., 2023c), ODEX (Wang et al., 2022b), SWE-Bench (Jimenez et al., 2023), GoogleCodeRepo (Shrivastava et al., 2023), RepoEval (Zhang et al., 2023a), ClasEval (Du et al., 2023) and Cocomic-Data (Ding et al., 2022).

A few benchmarks specifically measure competitive programming, such as APPS (Hendrycks et al., 2021), CodeContests (Li et al., 2022), CodeScope (Yan et al., 2023), xCodeEval (Khan et al., 2023), and LeetCode-Hard (Shinn et al., 2023), and TACO (Li et al., 2023c). Table 1 highlights differences between them. Methods such as AlphaCode (Li et al., 2022), AlphaCode 2 (Gemini Team et al., 2023), ALGO (Zhang et al., 2023d), Parsel (Zelikman et al., 2022), code cleaning (Jain et al., 2023), code explanations (Li et al., 2023a), analogical reasoning (Yasunaga et al., 2023), and AlphaCodium (Ridnik et al., 2024) have focused on improving LLMs.

6.1 HOLISTIC TASKS AND CONTAMINATION

Code Repair. (Chen et al., 2023; Olausson et al., 2023; Madaan et al., 2023b; Peng et al., 2023; Zhang et al., 2023c) have investigated self-repair for existing code LLM benchmarks. These methods use error feedback for models to improve inspiring our code repair scenario.

Code Execution. Code execution was first studied in (Austin et al., 2021; Nye et al., 2021) LIVE-CODEBENCH’s execution scenario is particularly inspired by CRUXEval (Gu et al., 2024), a recent benchmark measuring the reasoning and execution abilities of code LLMs. We differ from CRUXEval in that our benchmark is live, and our functions are more complex and human-produced (unlike Code Llama generations in CRUXEval).

Test Generation. Test generation using LLMs has been explored in (Yuan et al., 2023; Schäfer et al., 2024; Tufano et al., 2022; Watson et al., 2020). However, we decouple the test inputs and outputs which allows performing fair evaluations. Finally, some works have additionally studied other tasks and scenarios like type prediction (Mir et al., 2022; Wei et al., 2023; Malik et al., 2019), code summarization (LeClair et al., 2019; Iyer et al., 2016; Barone & Sennrich, 2017; Hasan et al., 2021; Alon et al., 2018), code security (Liguori et al., 2022; Pearce et al., 2022; Tony et al., 2023), etc.

Contamination. Data contamination and test-case leakage have received considerable attention Oren et al. (2024); Golchin & Surdeanu (2023); Weller et al. (2023); Roberts et al. (2024) recently. (Sainz et al., 2023) demonstrated contamination by simply prompting the model to highlight its contamination. Some detection methods have also been built to avoid these cases (Shi et al., 2023; Zhou et al., 2023). For code, (Riddell et al., 2024) use edit distance and AST-based semantic-similarity to detect contamination.

7 LIMITATIONS AND CONCLUSION

We describe key limitations here and provide a deeper discussion in Appendix G.

Evaluation Noise. We anticipate noise from benchmark size, prompts, and sampling. For benchmark size, while 612 problems is fairly larger than existing code benchmarks (e.g. 164 problems in HUMAN-EVAL), this set reduces on “scrolling” over the more recent problems. Next, prompting can cause large variations in model performance and we do not tune prompts across models. For sampling, we use bootstrapped PASS@1 using 10 completions, which should limit this noise considerably.

Problem Domain. LIVECODEBENCH comprises competition programming problems which might not correlate with how LLMs are used in practice. Even then, our results are well correlated with findings from human evaluations like Chatbot-Arena Chiang et al. (2024) thus providing useful signal.

Conclusion. In this work, we propose LIVECODEBENCH, a new benchmark for evaluating LLMs for code. LIVECODEBENCH provides an extensible framework that will keep on updating with new problems. Our benchmark mitigates contamination issues in existing benchmarks by introducing live evaluations and emphasizing scenarios beyond code generation to account for the broader coding abilities of LLMs. Our evaluations reveal novel findings such as contamination detection, holistic evaluations, and potential overfitting on HUMAN-EVAL. We hope LIVECODEBENCH will serve to advance understanding of current code LLMs and also guide future research through our findings.

REFERENCES

- 540
541
542 Rajas Agashe, Srinivasan Iyer, and Luke Zettlemoyer. Juice: A large scale distantly supervised
543 dataset for open domain context-based code generation. *arXiv preprint arXiv:1910.02216*, 2019.
- 544
545 Loubna Ben Allal, Raymond Li, Denis Kocetkov, Chenghao Mou, Christopher Akiki, Carlos Munoz
546 Ferrandis, Niklas Muennighoff, Mayank Mishra, Alex Gu, Manan Dey, et al. Santacoder: don't
547 reach for the stars! *arXiv preprint arXiv:2301.03988*, 2023.
- 548
549 Uri Alon, Shaked Brody, Omer Levy, and Eran Yahav. code2seq: Generating sequences from
550 structured representations of code. *arXiv preprint arXiv:1808.01400*, 2018.
- 551
552 Ben Athiwaratkun, Sanjay Krishna Gouda, Zijian Wang, Xiaopeng Li, Yuchen Tian, Ming Tan,
553 Wasi Uddin Ahmad, Shiqi Wang, Qing Sun, Mingyue Shang, et al. Multi-lingual evaluation of
554 code generation models. *arXiv preprint arXiv:2210.14868*, 2022.
- 555
556 Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan,
557 Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language
558 models. *arXiv preprint arXiv:2108.07732*, 2021.
- 559
560 Antonio Valerio Miceli Barone and Rico Sennrich. A parallel corpus of python functions and
561 documentation strings for automated code documentation and code generation. *arXiv preprint
562 arXiv:1707.02275*, 2017.
- 563
564 Xiao Bi, Deli Chen, Guanting Chen, Shanhuang Chen, Damai Dai, Chengqi Deng, Honghui Ding,
565 Kai Dong, Qiu Shi Du, Zhe Fu, et al. Deepseek llm: Scaling open-source language models with
566 longtermism. *arXiv preprint arXiv:2401.02954*, 2024.
- 567
568 Federico Cassano, John Gouwar, Daniel Nguyen, Sydney Nguyen, Luna Phipps-Costin, Donald
569 Pinckney, Ming-Ho Yee, Yangtian Zi, Carolyn Jane Anderson, Molly Q Feldman, et al. Multipl-e:
570 A scalable and extensible approach to benchmarking neural code generation. *arXiv preprint
571 arXiv:2208.08227*, 2022.
- 572
573 Bei Chen, Fengji Zhang, Anh Nguyen, Daoguang Zan, Zeqi Lin, Jian-Guang Lou, and Weizhu Chen.
574 Codet: Code generation with generated tests. *arXiv preprint arXiv:2207.10397*, 2022.
- 575
576 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared
577 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large
578 language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- 579
580 Xinyun Chen, Maxwell Lin, Nathanael Schärli, and Denny Zhou. Teaching large language models to
581 self-debug. *arXiv preprint arXiv:2304.05128*, 2023.
- 582
583 Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li, Dacheng
584 Li, Hao Zhang, Banghua Zhu, Michael Jordan, Joseph E. Gonzalez, and Ion Stoica. Chatbot arena:
585 An open platform for evaluating llms by human preference, 2024.
- 586
587 Yangruibo Ding, Zijian Wang, Wasi Uddin Ahmad, Murali Krishna Ramanathan, Ramesh Nallapati,
588 Parminder Bhatia, Dan Roth, and Bing Xiang. Cocomic: Code completion by jointly modeling
589 in-file and cross-file context. *arXiv preprint arXiv:2212.10007*, 2022.
- 590
591 Xueying Du, Mingwei Liu, Kaixin Wang, Hanlin Wang, Junwei Liu, Yixuan Chen, Jiayi Feng,
592 Chaofeng Sha, Xin Peng, and Yiling Lou. Classeval: A manually-crafted benchmark for evaluating
593 llms on class-level code generation, 2023.
- 594
595 Daniel Fried, Armen Aghajanyan, Jessy Lin, Sida Wang, Eric Wallace, Freda Shi, Ruiqi Zhong, Wen
596 tau Yih, Luke Zettlemoyer, and Mike Lewis. Incoder: A generative model for code infilling and
597 synthesis. *preprint arXiv:2204.05999*, 2022.
- 598
599 A Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu,
600 Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly
601 capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.

- 594 Shahriar Golchin and Mihai Surdeanu. Time travel in llms: Tracing data contamination in large
595 language models. *arXiv preprint arXiv:2308.08493*, 2023.
596
- 597 Alex Gu, Baptiste Rozière, Hugh Leather, Armando Solar-Lezama, Gabriel Synnaeve, and Sida I
598 Wang. Cruxeval: A benchmark for code reasoning, understanding and execution. *arXiv preprint*
599 *arXiv:2401.03065*, 2024.
- 600 Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao
601 Bi, Y Wu, YK Li, et al. Deepseek-coder: When the large language model meets programming—the
602 rise of code intelligence. *arXiv preprint arXiv:2401.14196*, 2024.
603
- 604 Masum Hasan, Tanveer Muttaqueen, Abdullah Al Ishtiaq, Kazi Sajeed Mehrab, Md Mahim Anjum
605 Haque, Tahmid Hasan, Wasi Uddin Ahmad, Anindya Iqbal, and Rifat Shahriyar. Codesc: A large
606 code-description parallel dataset. *arXiv preprint arXiv:2105.14220*, 2021.
- 607 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,
608 and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv*
609 *preprint arXiv:2103.03874*, 2021.
- 610 Yiming Huang, Zhenghao Lin, Xiao Liu, Yeyun Gong, Shuai Lu, Fangyu Lei, Yaobo Liang, Yelong
611 Shen, Chen Lin, Nan Duan, et al. Competition-level problems are effective llm evaluators. *arXiv*
612 *preprint arXiv:2312.02143*, 2023.
613
- 614 Srinivasan Iyer, Ioannis Konstas, Alvin Cheung, and Luke Zettlemoyer. Summarizing source code
615 using a neural attention model. In *54th Annual Meeting of the Association for Computational*
616 *Linguistics 2016*, pp. 2073–2083. Association for Computational Linguistics, 2016.
- 617 Naman Jain, Skanda Vaidyanath, Arun Iyer, Nagarajan Natarajan, Suresh Parthasarathy, Sriram
618 Rajamani, and Rahul Sharma. Jigsaw: Large language models meet program synthesis. In
619 *Proceedings of the 44th International Conference on Software Engineering*, pp. 1219–1231, 2022.
620
- 621 Naman Jain, Tianjun Zhang, Wei-Lin Chiang, Joseph E Gonzalez, Koushik Sen, and Ion Stoica.
622 Llm-assisted code cleaning for training accurate code generators. *arXiv preprint arXiv:2311.14904*,
623 2023.
- 624 Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik
625 Narasimhan. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint*
626 *arXiv:2310.06770*, 2023.
627
- 628 Mohammad Abdullah Matin Khan, M Saiful Bari, Xuan Long Do, Weishi Wang, Md Rizwan Parvez,
629 and Shafiq Joty. xcodeeval: A large scale multilingual multitask benchmark for code understanding,
630 generation, translation and retrieval. *arXiv preprint arXiv:2303.03004*, 2023.
- 631 Sumith Kulal, Panupong Pasupat, Kartik Chandra, Mina Lee, Oded Padon, Alex Aiken, and Percy S
632 Liang. Spoc: Search-based pseudocode to code. *Advances in Neural Information Processing*
633 *Systems*, 32, 2019.
634
- 635 Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E.
636 Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model
637 serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating*
638 *Systems Principles*, 2023.
- 639 Yuhang Lai, Chengxi Li, Yiming Wang, Tianyi Zhang, Ruiqi Zhong, Luke Zettlemoyer, Wen-tau Yih,
640 Daniel Fried, Sida Wang, and Tao Yu. Ds-1000: A natural and reliable benchmark for data science
641 code generation. In *International Conference on Machine Learning*, pp. 18319–18345. PMLR,
642 2023.
- 643 Alexander LeClair, Siyuan Jiang, and Collin McMillan. A neural model for generating natural
644 language summaries of program subroutines. In *2019 IEEE/ACM 41st International Conference*
645 *on Software Engineering (ICSE)*, pp. 795–806. IEEE, 2019.
646
- 647 Jierui Li, Szymon Tworowski, Yingying Wu, and Raymond Mooney. Explaining competitive-level
programming solutions using llms. *arXiv preprint arXiv:2307.05337*, 2023a.

- 648 Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou,
649 Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, et al. Starcoder: may the source be with
650 you! *arXiv preprint arXiv:2305.06161*, 2023b.
- 651
652 Rongao Li, Jie Fu, Bo-Wen Zhang, Tao Huang, Zhihong Sun, Chen Lyu, Guang Liu, Zhi Jin, and
653 Ge Li. Taco: Topics in algorithmic code generation dataset. *arXiv preprint arXiv:2312.14852*,
654 2023c.
- 655 Yuanzhi Li, Sébastien Bubeck, Ronen Eldan, Allie Del Giorno, Suriya Gunasekar, and Yin Tat Lee.
656 Textbooks are all you need ii: phi-1.5 technical report. *arXiv preprint arXiv:2309.05463*, 2023d.
- 657 Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom
658 Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, et al. Competition-level code generation
659 with alphacode. *Science*, 378(6624):1092–1097, 2022.
- 660
661 Pietro Liguori, Erfan Al-Hossami, Domenico Cotroneo, Roberto Natella, Bojan Cukic, and Samira
662 Shaikh. Can we generate shellcodes via natural language? an empirical study. *Automated Software
663 Engineering*, 29(1):30, 2022.
- 664 Changshu Liu, Shizhuo Dylan Zhang, and Reyhaneh Jabbarvand. Codemind: A framework to
665 challenge large language models for code reasoning. *arXiv preprint arXiv:2402.09664*, 2024.
- 666
667 Hanmeng Liu, Ruoxi Ning, Zhiyang Teng, Jian Liu, Qiji Zhou, and Yue Zhang. Evaluating the logical
668 reasoning ability of chatgpt and gpt-4. *arXiv preprint arXiv:2304.03439*, 2023a.
- 669 Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. Is your code generated by
670 chatgpt really correct? rigorous evaluation of large language models for code generation. *arXiv
671 preprint arXiv:2305.01210*, 2023b.
- 672
673 Tianyang Liu, Canwen Xu, and Julian McAuley. Repobench: Benchmarking repository-level code
674 auto-completion systems. *arXiv preprint arXiv:2306.03091*, 2023c.
- 675 Anton Lozhkov, Raymond Li, Loubna Ben Allal, Federico Cassano, Joel Lamy-Poirier, Nouamane
676 Tazi, Ao Tang, Dmytro Pykhtar, Jiawei Liu, Yuxiang Wei, Tianyang Liu, Max Tian, Denis Kocetkov,
677 Arthur Zucker, Younes Belkada, Zijian Wang, Qian Liu, Dmitry Abulkhanov, Indraneil Paul,
678 Zhuang Li, Wen-Ding Li, Megan Risdal, Jia Li, Jian Zhu, Terry Yue Zhuo, Evgenii Zheltonozhskii,
679 Nii Osa Osa Dade, Wenhao Yu, Lucas Krauß, Naman Jain, Yixuan Su, Xuanli He, Manan
680 Dey, Eduardo Abati, Yekun Chai, Niklas Muennighoff, Xiangru Tang, Muhtasham Oblokulov,
681 Christopher Akiki, Marc Marone, Chenghao Mou, Mayank Mishra, Alex Gu, Binyuan Hui, Tri
682 Dao, Armel Zebaze, Olivier Dehaene, Nicolas Patry, Canwen Xu, Julian McAuley, Torsten Scholak,
683 Sebastien Paquet, Jennifer Robinson, Carolyn Jane Anderson, Nicolas Chapados, Mostofa Patwary,
684 Nima Tajbakhsh, Yacine Jernite, Carlos Muñoz Ferrandis, Lingming Zhang, Sean Hughes, Thomas
685 Wolf, Arjun Guha, Leandro von Werra, and Harm de Vries. Starcoder 2 and the stack v2: The next
686 generation. 2024.
- 687
688 Qinyu Luo, Yining Ye, Shihao Liang, Zhong Zhang, Yujia Qin, Yaxi Lu, Yesai Wu, Xin Cong, Yankai
689 Lin, Yingli Zhang, et al. Repoagent: An llm-powered open-source framework for repository-level
690 code documentation generation. *arXiv preprint arXiv:2402.16667*, 2024.
- 691
692 Aman Madaan, Alexander Shypula, Uri Alon, Milad Hashemi, Parthasarathy Ranganathan, Yiming
693 Yang, Graham Neubig, and Amir Yazdanbakhsh. Learning performance-improving code edits.
694 *arXiv preprint arXiv:2302.07867*, 2023a.
- 695
696 Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri
697 Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. Self-refine: Iterative refinement
698 with self-feedback. *arXiv preprint arXiv:2303.17651*, 2023b.
- 699
700 Rabee Sohail Malik, Jibesh Patra, and Michael Pradel. NI2type: inferring javascript function types
701 from natural language information. In *2019 IEEE/ACM 41st International Conference on Software
Engineering (ICSE)*, pp. 304–315. IEEE, 2019.
- 702
703 Amir M Mir, Evaldas Latoškinas, Sebastian Proksch, and Georgios Gousios. Type4py: Practical
704 deep similarity learning-based type inference for python. In *Proceedings of the 44th International
705 Conference on Software Engineering*, pp. 2241–2252, 2022.

- 702 Ansong Ni, Pengcheng Yin, Yilun Zhao, Martin Riddell, Troy Feng, Rui Shen, Stephen Yin, Ye Liu,
703 Semih Yavuz, Caiming Xiong, et al. L2ceval: Evaluating language-to-code generation capabilities
704 of large language models. *arXiv preprint arXiv:2309.17446*, 2023.
- 705 Erik Nijkamp, Bo Pang, Hiroaki Hayashi, Lifu Tu, Huan Wang, Yingbo Zhou, Silvio Savarese,
706 and Caiming Xiong. Codegen: An open large language model for code with multi-turn program
707 synthesis. In *The Eleventh International Conference on Learning Representations*, 2022.
- 708 Maxwell Nye, Anders Johan Andreassen, Guy Gur-Ari, Henryk Michalewski, Jacob Austin, David
709 Bieber, David Dohan, Aitor Lewkowycz, Maarten Bosma, David Luan, et al. Show your work:
710 Scratchpads for intermediate computation with language models. *arXiv preprint arXiv:2112.00114*,
711 2021.
- 712 Theo X Olausson, Jeevana Priya Inala, Chenglong Wang, Jianfeng Gao, and Armando Solar-Lezama.
713 Demystifying gpt self-repair for code generation. *arXiv preprint arXiv:2306.09896*, 2023.
- 714 R OpenAI. Gpt-4 technical report. arxiv 2303.08774. *View in Article*, 2023.
- 715 Yonatan Oren, Nicole Meister, Niladri Chatterji, Faisal Ladhak, and Tatsunori B Hashimoto. Proving
716 test set contamination for black-box language models. In *The Twelfth International Confer-*
717 *ence on Learning Representations*, 2024. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=KS8mIvetg2)
718 [KS8mIvetg2](https://openreview.net/forum?id=KS8mIvetg2).
- 719 Shishir G Patil, Tianjun Zhang, Xin Wang, and Joseph E Gonzalez. Gorilla: Large language model
720 connected with massive apis. *arXiv preprint arXiv:2305.15334*, 2023.
- 721 Hammond Pearce, Baleegh Ahmad, Benjamin Tan, Brendan Dolan-Gavitt, and Ramesh Karri. Asleep
722 at the keyboard? assessing the security of github copilot’s code contributions. In *2022 IEEE*
723 *Symposium on Security and Privacy (SP)*, pp. 754–768. IEEE, 2022.
- 724 Baolin Peng, Michel Galley, Pengcheng He, Hao Cheng, Yujia Xie, Yu Hu, Qiuyuan Huang, Lars
725 Liden, Zhou Yu, Weizhu Chen, and Jianfeng Gao. Check your facts and try again: Improv-
726 ing large language models with external knowledge and automated feedback. *arXiv preprint*
727 *arXiv:2302.12813*, 2023.
- 728 Martin Riddell, Ansong Ni, and Arman Cohan. Quantifying contamination in evaluating code
729 generation capabilities of language models, 2024.
- 730 Tal Ridnik, Dedy Kredo, and Itamar Friedman. Code generation with alphacodium: From prompt
731 engineering to flow engineering. *arXiv preprint arXiv:2401.08500*, 2024.
- 732 Manley Roberts, Himanshu Thakur, Christine Herlihy, Colin White, and Samuel Dooley. To the
733 cutoff... and beyond? a longitudinal perspective on LLM data contamination. In *The Twelfth*
734 *International Conference on Learning Representations*, 2024. URL [https://openreview.](https://openreview.net/forum?id=m2NVG4Htxs)
735 [net/forum?id=m2NVG4Htxs](https://openreview.net/forum?id=m2NVG4Htxs).
- 736 Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi
737 Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, et al. Code llama: Open foundation models for code.
738 *arXiv preprint arXiv:2308.12950*, 2023.
- 739 Oscar Sainz, Jon Ander Campos, Iker García-Ferrero, Julen Etxaniz, and Eneko Agirre. Did
740 chatgpt cheat on your test?, Jun 2023. URL [https://hitz-zentroa.github.io/](https://hitz-zentroa.github.io/lm-contamination/blog/)
741 [lm-contamination/blog/](https://hitz-zentroa.github.io/lm-contamination/blog/).
- 742 Max Schäfer, Sarah Nadi, Aryaz Eghbali, and Frank Tip. An empirical evaluation of using large
743 language models for automated unit test generation. *IEEE Transactions on Software Engineering*,
744 50(1):85–105, 2024. doi: 10.1109/TSE.2023.3334955.
- 745 Quan Shi, Michael Tang, Karthik Narasimhan, and Shunyu Yao. Can language models solve olympiad
746 programming?, 2024.
- 747 Weijia Shi, Anirudh Ajith, Mengzhou Xia, Yangsibo Huang, Daogao Liu, Terra Blevins, Danqi Chen,
748 and Luke Zettlemoyer. Detecting pretraining data from large language models. *arXiv preprint*
749 *arXiv:2310.16789*, 2023.

- 756 Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik R Narasimhan, and Shunyu Yao. Reflexion:
757 Language agents with verbal reinforcement learning. In *Thirty-seventh Conference on Neural*
758 *Information Processing Systems*, 2023.
- 759
760 Disha Shrivastava, Hugo Larochelle, and Daniel Tarlow. Repository-level prompt generation for large
761 language models of code. In *International Conference on Machine Learning*, pp. 31693–31715.
762 PMLR, 2023.
- 763 Manav Singhal, Tushar Aggarwal, Abhijeet Awasthi, Nagarajan Natarajan, and Aditya Kanade.
764 Nofuneval: Funny how code lms falter on requirements beyond functional correctness. *arXiv*
765 *preprint arXiv:2401.15963*, 2024.
- 766
767 Benjamin Steenhoek, Michele Tufano, Neel Sundaresan, and Alexey Svyatkovskiy. Reinforce-
768 ment learning from automatic feedback for high-quality unit test generation. *arXiv preprint*
769 *arXiv:2310.02368*, 2023.
- 770
771 Ruoxi Sun, Sercan O Arik, Hootan Nakhost, Hanjun Dai, Rajarishi Sinha, Pengcheng Yin, and
772 Tomas Pfister. Sql-palm: Improved large language model adaptation for text-to-sql. *arXiv preprint*
773 *arXiv:2306.00739*, 2023.
- 774
775 Gemini Team. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of con-
776 text, 2024. URL [https://storage.googleapis.com/deepmind-media/gemini/
776 gemini_v1_5_report.pdf](https://storage.googleapis.com/deepmind-media/gemini/gemini_v1_5_report.pdf).
- 777
778 Catherine Tony, Markus Mutas, Nicolás E Díaz Ferreyra, and Riccardo Scandariato. Llmseceval: A
779 dataset of natural language prompts for security evaluations. *arXiv preprint arXiv:2303.09384*,
780 2023.
- 781
782 Michele Tufano, Shao Kun Deng, Neel Sundaresan, and Alexey Svyatkovskiy. Methods2test: A
783 dataset of focal methods mapped to test cases. In *Proceedings of the 19th International Conference*
784 *on Mining Software Repositories*, pp. 299–303, 2022.
- 785
786 Shiqi Wang, Zheng Li, Haifeng Qian, Chenghao Yang, Zijian Wang, Mingyue Shang, Varun Kumar,
787 Samson Tan, Baishakhi Ray, Parminder Bhatia, et al. Recode: Robustness evaluation of code
788 generation models. *arXiv preprint arXiv:2212.10264*, 2022a.
- 789
790 Yue Wang, Weishi Wang, Shafiq Joty, and Steven CH Hoi. Codet5: Identifier-aware unified pre-
791 trained encoder-decoder models for code understanding and generation. In *Proceedings of the*
792 *2021 Conference on Empirical Methods in Natural Language Processing*, pp. 8696–8708, 2021.
- 793
794 Yue Wang, Hung Le, Akhilesh Deepak Gotmare, Nghi DQ Bui, Junnan Li, and Steven CH Hoi.
795 Codet5+: Open code large language models for code understanding and generation. *arXiv preprint*
796 *arXiv:2305.07922*, 2023.
- 797
798 Zhiruo Wang, Shuyan Zhou, Daniel Fried, and Graham Neubig. Execution-based evaluation for
799 open-domain code generation. *arXiv preprint arXiv:2212.10481*, 2022b.
- 800
801 Cody Watson, Michele Tufano, Kevin Moran, Gabriele Bavota, and Denys Poshyvanyk. On learn-
802 ing meaningful assert statements for unit test cases. In *Proceedings of the ACM/IEEE 42nd*
803 *International Conference on Software Engineering*, pp. 1398–1409, 2020.
- 804
805 Jiayi Wei, Greg Durrett, and Isil Dillig. Typet5: Seq2seq type inference using static analysis. *arXiv*
806 *preprint arXiv:2303.09564*, 2023.
- 807
808 Orion Weller, Marc Marone, Nathaniel Weir, Dawn Lawrie, Daniel Khashabi, and Benjamin
809 Van Durme. "according to..." prompting language models improves quoting from pre-training
data. *arXiv preprint arXiv:2305.13252*, 2023.
- 810
811 Weixiang Yan, Haitian Liu, Yunkun Wang, Yunzhe Li, Qian Chen, Wen Wang, Tingyu Lin, Weishan
Zhao, Li Zhu, Shuiguang Deng, et al. Codescope: An execution-based multilingual multitask
multidimensional benchmark for evaluating llms on code understanding and generation. *arXiv*
preprint arXiv:2311.08588, 2023.

- 810 Shuo Yang, Wei-Lin Chiang, Lianmin Zheng, Joseph E. Gonzalez, and Ion Stoica. Rethinking
811 benchmark and contamination for language models with rephrased samples, 2023.
812
- 813 Michihiro Yasunaga, Xinyun Chen, Yujia Li, Panupong Pasupat, Jure Leskovec, Percy Liang,
814 Ed H Chi, and Denny Zhou. Large language models as analogical reasoners. *arXiv preprint*
815 *arXiv:2310.01714*, 2023.
- 816 Pengcheng Yin, Wen-Ding Li, Kefan Xiao, Abhishek Rao, Yeming Wen, Kensen Shi, Joshua Howland,
817 Paige Bailey, Michele Catasta, Henryk Michalewski, et al. Natural language to code generation in
818 interactive data science notebooks. *arXiv preprint arXiv:2212.09248*, 2022.
819
- 820 Zhiqiang Yuan, Yiling Lou, Mingwei Liu, Shiji Ding, Kaixin Wang, Yixuan Chen, and Xin Peng.
821 No more manual tests? evaluating and improving chatgpt for unit test generation. *arXiv preprint*
822 *arXiv:2305.04207*, 2023.
- 823 Eric Zelikman, Qian Huang, Gabriel Poesia, Noah D Goodman, and Nick Haber. Parsel: A unified
824 natural language framework for algorithmic reasoning. *arXiv preprint arXiv:2212.10561*, 2022.
825
- 826 Fengji Zhang, Bei Chen, Yue Zhang, Jin Liu, Daoguang Zan, Yi Mao, Jian-Guang Lou, and Weizhu
827 Chen. Repocoder: Repository-level code completion through iterative retrieval and generation.
828 *arXiv preprint arXiv:2303.12570*, 2023a.
- 829 Kechi Zhang, Ge Li, Jia Li, Zhuo Li, and Zhi Jin. Toolcoder: Teach code generation models to use
830 apis with search tools. *arXiv preprint arXiv:2305.04032*, 2023b.
831
- 832 Kechi Zhang, Zhuo Li, Jia Li, Ge Li, and Zhi Jin. Self-edit: Fault-aware code editor for code
833 generation. In *Proceedings of the 61st Annual Meeting of the Association for Computational*
834 *Linguistics (Volume 1: Long Papers)*, pp. 769–787, Toronto, Canada, July 2023c. Association for
835 Computational Linguistics.
- 836 Kexun Zhang, Danqing Wang, Jingtao Xia, William Yang Wang, and Lei Li. Algo: Synthesizing
837 algorithmic programs with generated oracle verifiers. *arXiv preprint arXiv:2305.14591*, 2023d.
838
- 839 Qinkai Zheng, Xiao Xia, Xu Zou, Yuxiao Dong, Shan Wang, Yufei Xue, Zihan Wang, Lei Shen,
840 Andi Wang, Yang Li, et al. Codegeex: A pre-trained model for code generation with multilingual
841 evaluations on humaneval-x. *arXiv preprint arXiv:2303.17568*, 2023.
- 842 Tianyu Zheng, Ge Zhang, Tianhao Shen, Xueling Liu, Bill Yuchen Lin, Jie Fu, Wenhui Chen, and
843 Xiang Yue. Opencodeinterpreter: Integrating code generation with execution and refinement.
844 <https://arxiv.org/abs/2402.14658>, 2024.
- 845 Kun Zhou, Yutao Zhu, Zhipeng Chen, Wentong Chen, Wayne Xin Zhao, Xu Chen, Yankai Lin,
846 Ji-Rong Wen, and Jiawei Han. Don’t make your llm an evaluation benchmark cheater. *arXiv*
847 *preprint arXiv:2311.01964*, 2023.
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

A DATASET

A.1 LEGAL COMPLIANCE AND LICENSE

Our benchmark does not comprise and personal identifiable information or offensive content.

Similar to (Hendrycks et al., 2021), we scrape only the problem statements, ground-truth solutions, and test cases from competition websites – LEETCODE, ATCODER, and CODEFORCES. Further, we primarily scrape publicly visible portions of websites, avoiding any data collection that might be pay-walled or require login or interaction with the website. Following, (Hendrycks et al., 2021) we abide by Fair Use §107: “the fair use of a copyrighted work, including such use by ... scholarship, or research, is not an infringement of copyright”, where fair use is determined by “the purpose and character of the use, including whether such use is of a commercial nature or is for nonprofit educational purposes”, “the amount and substantiality of the portion used in relation to the copyrighted work as a whole”, and “the effect of the use upon the potential market for or value of the copyrighted work.” Finally, we use the collected problems for academic purposes only and in addition, do not train on the collected problems.

A.2 GENERATOR BASED TEST GENERATION

We use GPT-4-TURBO to construct input generators. The following prompts (Figures 6 and 7) provide one-shot prompt templates used for synthesizing random and adversarial input generators. These generators define a function returns the arguments sampled in some distribution. These generators are then executed to construct inputs which validated on the collected correct programs. We use separate generators for random and adversarial setting since often times programming problems have corner cases which might not be captured by randomly sampling over the inputs. We build 2 random input generators, 4 adversarial input generators and check if the sampled inputs work for the correct programs. Finally, the number of collected inputs is thresholded to 100 for efficient grading (using random selection). We find that our generators can already function well but future work can study the design space of constructing such generators.

Note that for CODEFORCES, we construct the generators in semi-autonomous manner since only 9 problems were used.

A.3 CODE EXECUTION

The code execution split of LiveCodeBench consists of 479 samples from 85 distinct problems. To encourage diversity in our benchmark while keeping our benchmark small and usable, we place a limit of six samples for each given problem. These sample programs and corresponding test cases are chosen uniformly at random from all those passing the filter.

Filtering Criteria: The specific filtering criteria are as follows:

- Compile time: length of code is between 100 and 500 characters, no syntax errors, all necessary imports are included
- Runtime: no floating point operations, true division, exp, other integer operations must have at least one argument ≤ 3 , string and list operations must have at least one argument with length ≤ 3 , finish running in 2 seconds, “reasonable” number of steps (roughly, under 1000 Python bytecode operations).

We give two examples of two programs that are filtered out in the Listings below. Our final benchmark consists of 479 samples from 85 problems, but will increase in size due to its live nature.

Dataset Statistics: As mentioned, we filter for codes between 100 and 500 characters, as well as below 1000 steps. The statistics for programs in our dataset are shown in Fig. 8.

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

```

You are an expert python competitive programmer and your goal is to
↳ construct input-generators for testing programming contest
↳ problems. You will write relevant generators and finally
↳ construct 'construct_inputs' function that returns a list of
↳ diverse inputs sampled from the generator. Remember to strictly
↳ follow the instructions and constraints present in the problem
↳ statement

<USER>
{EXAMPLE_PROBLEM}

Construct a random input generator.
</USER>

<ASSISTANT>
import numpy as np
def random_input_generator(weight_min, weight_max, size_min, size_max):
    weights_size = np.random.randint(size_min, size_max+1)
    weights = np.random.randint(weight_min, weight_max, size=
↳ weights_size).tolist()
    k = np.random.randint(1, len(weights)+1)
    return weights, k

def construct_inputs():
    inputs_list = []
    ## small inputs
    for i in range(15):
        inputs_list.append(random_input_generator(1, 10**3, 1, 10))
    ## medium inputs
    for i in range(15):
        inputs_list.append(random_input_generator(1, 10**6, 1, 10**3))
    ## large inputs
    for i in range(15):
        inputs_list.append(random_input_generator(1, 10**9, 1, 10**5))
    return inputs_list
</ASSISTANT>

<USER>
{PROBLEM}

Construct a random input generator. Use the format used in the above
↳ example by returning a single function that builds diverse
↳ inputs named 'construct_inputs'
</USER>

```

Prompt for random input generation

Figure 6: Random Input Generator Prompt

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

```

You are an expert python competitive programmer and your goal is to
↳ construct input-generators for testing programming contest
↳ problems. You will write relevant generators and finally
↳ construct 'construct_inputs' function that returns a list of
↳ diverse inputs sampled from the generator. Remember to strictly
↳ follow the instructions and constraints present in the problem
↳ statement

<USER>
{EXAMPLE_PROBLEM}

Construct an adversarial input generator.
</USER>

<ASSISTANT>
import numpy as np
## case 1 - alternating large and small weights
def generate_adversarial_inputs_1(weight_size, max_weight, k):
    weights = [1 if i%2==0 else max_weight for i in range(weight_size)]
    return weights, k

## case 2 - equal_weights
def adversarial_input_generator_2(weight_size, max_weight, k):
    weights = [max_weight for _ in range(weight_size)]
    return weights, k

# Case 3 - Large weights at the ends
def adversarial_input_generator_3(weight_size, max_weight, k):
    pass ## truncatede
def construct_inputs():
    inputs_list = []

    weight_sizes = [10, 1000, 100000]
    max_weights = [10**3, 10**6, 10**9]

    for weight_size in weight_sizes:
        for max_weight in max_weights:
            ks = [1, 2, 5, weight_size//2, weight_size-1, weight_size]
            for k in ks:
                inputs_list.append(generate_adversarial_inputs_1(
↳ weight_size, max_weight, k))
                # truncated
    return inputs_list
</ASSISTANT>

<USER>
{PROBLEM}

Construct an adversarial input generator. Use the format used in the
↳ above example by returning a single function that builds diverse
↳ inputs named 'construct_inputs'
</USER>

```

Prompt for adversarial input generation

Figure 7: Adversarial Input Generator Prompt

```

1026
1027 def check(x, t):
1028     if x == '':
1029         return t == 0
1030     if t < 0:
1031         return False
1032     for i in range(len(x)):
1033         if check(x[:i], t - int(x[i])):
1034             return True
1035     return False
1036
1037 @cache
1038 def punishmentNumber(n: int) -> int:
1039     if n == 0:
1040         return 0
1041     ans = punishmentNumber(n-1)
1042     if check(str(n * n), n):
1043         ans += n * n
1044     return ans
1045
1046 assert punishmentNumber(n = 37) == 1478

```

Program filtered because of multiplication

```

1047 dp = [True for _ in range(int(1e6 + 5))]
1048 MAXN = int(1e6 + 5)
1049 p = []
1050 dp[0] = False
1051 dp[1] = False
1052 for i in range(2, MAXN):
1053     if not dp[i]: continue
1054     p.append(i)
1055     for j in range(2 * i, MAXN, i):
1056         dp[j] = False
1057 def findPrimePairs(n: int) -> List[List[int]]:
1058     res = []
1059     for i in range(1, n):
1060         if n % 2 == 1 and i > n//2: break
1061         if n % 2 == 0 and i > n//2: break
1062         if dp[i] and dp[n - i]:
1063             res.append([i, n - i])
1064     return res
1065
1066 assert findPrimePairs(n = 2) == []

```

Program filtered because of control flow

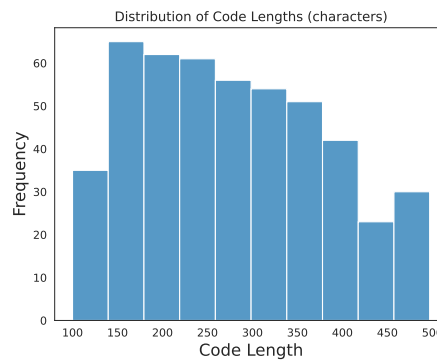
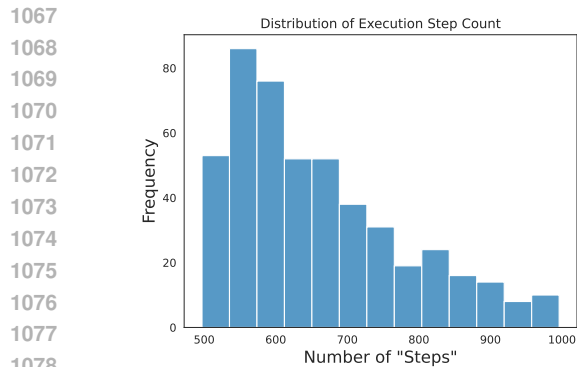
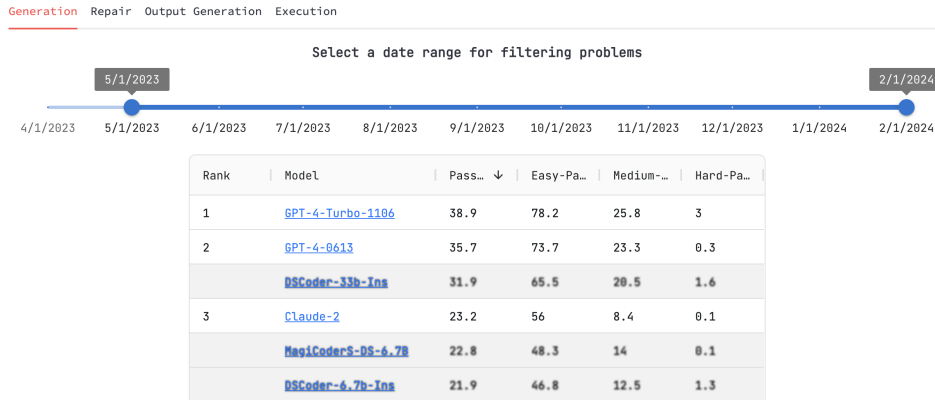


Figure 8: Distribution of code lengths and number of execution steps

B UI

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

LiveCodeBench



LiveCodeBench

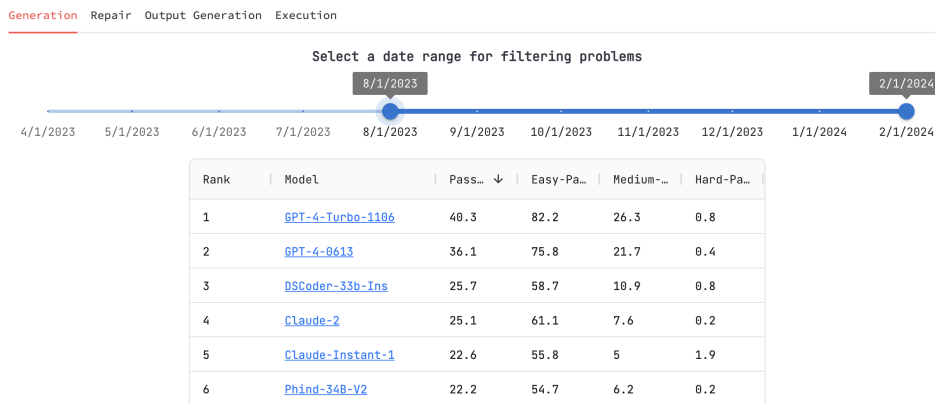


Figure 9: UI of LIVECODEBENCH showing two views – May-Jan and Sep-Jan. The contaminated models are blurred and the performance difference is visible across the two views. The scroller on the top allows selecting different periods of time highlighting the live nature of the benchmark.

C CURATION

C.1 PLATFORM SPECIFIC CURATION

We describe the curation process for each platform.

LEETCODE. We collect problems from all weekly and biweekly contests on LEETCODE that have taken place after April’23. For each problem, we collect the problems, public tests, and user solutions. The platform also provides a difficulty label for each problem which we use to tag the problems as EASY, MEDIUM, and HARD. Since LEETCODE provides a starter code for each problem, we also collect it and provide it to the LLM in the STDIN format. Since the hidden tests are not directly available, we use our generator-based test input generation approach (Section A.2) and also collect the auto grader failing tests for some of the recent problems.

ATCODER. We collect problems from the `abc` (beginner round) contests on ATCODER that have taken place after April’23. We deliberately avoid the more challenging `arc` and `agc` contests which are designed for more advanced Olympiad participants. The problems are assigned numeric difficulty ratings, and we exclude `abc` problems with a rating of more than 500. We also use these numeric ratings to tag the problems as EASY, MEDIUM, and HARD. Specifically, we use the rating brackets $[0 - 200)$, $[200 - 400)$, and $[400 - 500]$ to perform the classification. ATCODER provides public and hidden tests for each problem which we directly use in the benchmark.

CODEFORCES. We have collected problems from the Division 3 and Division 4 contests on CODEFORCES. Notably, we find that even with this filter, the problems are harder than the other two platforms. CODEFORCES also provides difficulty ratings for the problems which we use to tag the problems as EASY, MEDIUM, and HARD using the rating brackets $\{800\}$, $(800 - 1000]$, and $(1000 - 1300]$ respectively. Due to the higher difficulty, we only consider a small fraction of problems from CODEFORCES and semi-automatically construct test case generators, as they do not provide complete tests on the platform (long tests are truncated).

Table 2 provides various statistics about the problems that we have collected for LIVECODEBENCH.

C.2 SCENARIO-SPECIFIC BENCHMARK CONSTRUCTION

Code Generation and Self-Repair. We use the natural language problem statement as the problem statement for these scenarios. For LEETCODE, as noted above, an additional starter code is provided for the functional input format. For ATCODER and CODEFORCES problems, we use the standard input format (similar to (Hendrycks et al., 2021)). The collected or generated tests are then used to evaluate the correctness of the generated programs. Our final dataset consists of 511 problem instances across the three platforms.

Code Execution. We draw inspiration from the benchmark creation procedure used in (Gu et al., 2024). First, we collect a large pool of ~ 2000 *correct, human-submitted solutions* from the

| Platform | Total Count | #Easy | #Medium | #Hard | Average Tests |
|---------------|-------------|-------|---------|-------|---------------|
| LCB (May-end) | 511 | 182 | 206 | 123 | 17.0 |
| LCB (Sep-end) | 349 | 125 | 136 | 88 | 18.0 |
| ATCODER | 267 | 99 | 91 | 77 | 15.6 |
| LEETCODE | 235 | 79 | 113 | 43 | 19.0 |
| CODEFORCES | 9 | 4 | 2 | 3 | 11.1 |
| LCB-Easy | 182 | 182 | 0 | 0 | 16.1 |
| LCB-Medium | 206 | 0 | 206 | 0 | 17.4 |
| LCB-Hard | 123 | 0 | 0 | 123 | 18.0 |

Table 2: The statistics of problems collected in LIVECODEBENCH (LCB). We present the number of problems, their difficulty distributions and the average number of tests per problem. We present the results on the following subsets of LIVECODEBENCH (used throughout this manuscript) - (a) problems in the May-Feb and Sep-Feb time windows, (b) problems sourced from the three platforms, and (c) problems in the LCB-Easy, LCB-Medium, and LCB-Hard subsets.

1188 LEETCODE subset. However, many of these programs have multiple nested loops, complex numerical
1189 computations, and a large number of execution steps. Therefore, we apply compile-time and run-time
1190 filters to ensure samples are reasonable, and we double-check this with a manual inspection. More
1191 details on the filtering criteria and statistics of the dataset can be found in Appendix A.3. Our final
1192 dataset consists of 479 samples from 85 problems.

1193 **Test Case Output Prediction.** We use the natural language problem statement from the LEETCODE
1194 platform and the example test inputs to construct our test case output prediction dataset. Since
1195 the example test inputs in the problems are reasonable test cases for humans to reason about and
1196 understand the problems, they also serve as ideal test inputs for LLMs to process. Our final dataset
1197 consists of 442 problem instances from a total of 181 LEETCODE problems.
1198

1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241

1242 D EXPERIMENTAL SETUP

1243

1244 D.1 MODELS

1245

1246 We describe the details of models considered in our study in Table 3.

1247

| 1248 Model ID | 1249 Short Name | 1250 Approximate Cutoff Date | 1251 Link |
|---|------------------------|-------------------------------------|------------------------------|
| 1252 deepseek-ai/deepseek-coder-33b-instruct | DSCoder-33b-Ins | 08/30/2023 | deepseek-coder-33b-instruct |
| 1253 deepseek-ai/deepseek-coder-6.7b-instruct | DSCoder-6.7b-Ins | 08/30/2023 | deepseek-coder-6.7b-instruct |
| 1254 deepseek-ai/deepseek-coder-1.3b-instruct | DSCoder-1.3b-Ins | 08/30/2023 | deepseek-coder-1.3b-instruct |
| 1255 codellama/CodeLlama-70b-Instruct-hf | CodeLlama-70b-Ins | 01/01/2023 | CodeLlama-70b-Instruct-hf |
| 1256 openbmb/Eurus-70b-sft | Eurus-70B-SFT (n=1) | 01/01/2023 | Eurus-70b-sft |
| 1257 openbmb/Eurux-8x22b-nca | Eurux-8x22b-NCA (n=1) | 04/30/2023 | Eurux-8x22b-nca |
| 1258 codellama/CodeLlama-34b-Instruct-hf | CodeLlama-34b-Ins | 01/01/2023 | CodeLlama-34b-Instruct-hf |
| 1259 codellama/CodeLlama-13b-Instruct-hf | CodeLlama-13b-Ins | 01/01/2023 | CodeLlama-13b-Instruct-hf |
| 1260 codellama/CodeLlama-7b-Instruct-hf | CodeLlama-7b-Ins | 01/01/2023 | CodeLlama-7b-Instruct-hf |
| 1261 meta-llama/Meta-Llama-3-8B-Instruct | LLama3-8b-Ins | 01/01/2023 | Meta-Llama-3-8B-Instruct |
| 1262 meta-llama/Meta-Llama-3-70B-Instruct | LLama3-70b-Ins | 01/01/2023 | Meta-Llama-3-70B-Instruct |
| 1263 Phind/Phind-CodeLlama-34B-v2 | Phind-34B-V2 | 01/01/2023 | Phind-CodeLlama-34B-v2 |
| 1264 Smaug-2-72B | Smaug-2-72B | 01/01/2023 | Smaug-2-72B |
| 1265 Qwen-1.5-72B-Chat | Qwen-1.5-72B-Chat | 01/01/2023 | Qwen-1.5-72B-Chat |
| 1266 Qwen/CodeQwen1.5-7B | CodeQwen15-7B | 08/30/2023 | CodeQwen1.5-7B |
| 1267 Qwen/CodeQwen1.5-7B-Chat | CodeQwen15-7B-Chat | 08/30/2023 | CodeQwen1.5-7B-Chat |
| 1268 gpt-3.5-turbo-0301 | GPT-3.5-Turbo-0301 | 10/01/2021 | gpt-3.5-turbo-0301 |
| 1269 gpt-3.5-turbo-0125 | GPT-3.5-Turbo-0125 | 10/01/2021 | gpt-3.5-turbo-0125 |
| 1270 gpt-4-0613 | GPT-4-0613 | 10/01/2021 | gpt-4-0613 |
| 1271 gpt-4-1106-preview | GPT-4-Turbo-1106 | 04/30/2023 | gpt-4-1106-preview |
| 1272 gpt-4-turbo-2024-04-09 | GPT-4-Turbo-2024-04-09 | 04/30/2023 | gpt-4-turbo-2024-04-09 |
| 1273 gpt-4o-2024-05-13 | GPT-4O-2024-05-13 | 10/30/2023 | gpt-4o-2024-05-13 |
| 1274 claude-2 | Claude-2 | 12/31/2022 | claude-2 |
| 1275 claude-instant-1 | Claude-Instant-1 | 12/31/2022 | claude-instant-1 |
| 1276 claude-3-opus-20240229 | Claude-3-Opus | 04/30/2023 | claude-3-opus-20240229 |
| 1277 claude-3-sonnet-20240229 | Claude-3-Sonnet | 04/30/2023 | claude-3-sonnet-20240229 |
| 1278 claude-3-haiku-20240307 | Claude-3-Haiku | 04/30/2023 | claude-3-haiku-20240307 |
| 1279 codestral-latest | Codestral-Latest | 01/31/2024 | codestral-latest |
| 1280 gemini-pro | Gemini-Pro | 04/30/2023 | gemini-pro |
| 1281 gemini-1.5-pro-latest | Gemini-Pro-1.5-May | 04/30/2023 | gemini-1.5-pro-latest |
| 1282 gemini-1.5-flash-latest | Gemini-Flash-1.5-May | 04/30/2023 | gemini-1.5-flash-latest |
| 1283 ise-uiuc/Magicoder-S-DS-6.7B | MagiCoderS-DS-6.7B | 08/30/2023 | Magicoder-S-DS-6.7B |
| 1284 ise-uiuc/Magicoder-S-CL-7B | MagiCoderS-CL-7B | 01/01/2023 | Magicoder-S-CL-7B |
| 1285 bigcode/starcoder2-3b | StarCoder2-3b | 01/01/2023 | starcoder2-3b |
| 1286 bigcode/starcoder2-7b | StarCoder2-7b | 01/01/2023 | starcoder2-7b |
| 1287 bigcode/starcoder2-15b | StarCoder2-15b | 01/01/2023 | starcoder2-15b |

| | | | | |
|------|--------------------------------------|--------------------|------------|-----------------------------|
| 1296 | codellama/CodeLlama-70b-hf | CodeLlama-70b-Base | 01/01/2023 | CodeLlama-70b-hf |
| 1297 | codellama/CodeLlama-34b-hf | CodeLlama-34b-Base | 01/01/2023 | CodeLlama-34b-hf |
| 1298 | codellama/CodeLlama-13b-hf | CodeLlama-13b-Base | 01/01/2023 | CodeLlama-13b-hf |
| 1299 | codellama/CodeLlama-7b-hf | CodeLlama-7b-Base | 01/01/2023 | CodeLlama-7b-hf |
| 1300 | deepseek-ai/deepseek-coder-33b-base | DSCoder-33b-Base | 08/30/2023 | deepseek-coder-33b-base |
| 1301 | deepseek-ai/deepseek-coder-6.7b-base | DSCoder-6.7b-Base | 08/30/2023 | deepseek-coder-6.7b-base |
| 1302 | deepseek-ai/deepseek-coder-1.3b-base | DSCoder-1.3b-Base | 08/30/2023 | deepseek-coder-1.3b-base |
| 1303 | google/codegemma-7b | CodeGemma-7b-Base | 01/01/2023 | codegemma-7b |
| 1304 | google/codegemma-2b | CodeGemma-2b-Base | 01/01/2023 | codegemma-2b |
| 1305 | google/gemma-7b | Gemma-7b-Base | 01/01/2023 | gemma-7b |
| 1306 | google/gemma-2b | Gemma-2b-Base | 01/01/2023 | gemma-2b |
| 1307 | meta-llama/Meta-Llama-3-70B | LLama3-70b-Base | 01/01/2023 | Meta-Llama-3-70B |
| 1308 | meta-llama/Meta-Llama-3-8B | LLama3-8b-Base | 01/01/2023 | Meta-Llama-3-8B |
| 1309 | mistral-large-latest | Mistral-Large | 01/01/2023 | mistral-large-latest |
| 1310 | open-mixtral-8x22b | Mixtral-8x22B-Ins | 01/01/2023 | open-mixtral-8x22b |
| 1311 | open-mixtral-8x7b | Mixtral-8x7B-Ins | 01/01/2023 | open-mixtral-8x7b |
| 1312 | m-a-p/OpenCodeInterpreter-DS-33B | OC-DS-33B | 08/30/2023 | OpenCodeInterpreter-DS-33B |
| 1313 | m-a-p/OpenCodeInterpreter-DS-6.7B | OC-DS-6.7B | 08/30/2023 | OpenCodeInterpreter-DS-6.7B |
| 1314 | m-a-p/OpenCodeInterpreter-DS-1.3B | OC-DS-1.3B | 08/30/2023 | OpenCodeInterpreter-DS-1.3B |
| 1315 | command-r | Command-R | 01/01/2023 | command-r |
| 1316 | command-r+ | Command-R+ | 01/01/2023 | command-r+ |

Table 3: Language Models Overview

We use variety of GPUs based on availability for running local models (A6000, L4, A100).

D.2 CODE GENERATION

Below we provide the prompt format (with appropriate variants adding special tokens accommodating each instruct-tuned model) used for this scenario.

D.3 SELF REPAIR

Below we provide the prompt format (with appropriate variants adding special tokens accommodating each instruct-tuned model) used for this scenario.

D.4 CODE EXECUTION

Below we provide the prompts for code execution with and without CoT. The prompts are modified versions of those from (Gu et al., 2024) to fit the format of the samples in our benchmark.

D.5 TEST OUTPUT PREDICTION

Below we provide the prompt format (with appropriate variants adding special tokens accommodating each instruct-tuned model) used for this scenario.

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

You are an expert Python programmer. You will be given a question (↪ problem specification) and will generate a correct Python ↪ program that matches the specification and passes all tests. You ↪ will NOT return anything except for the program

```

### Question:\n{question.question_content}

{ if question.starter_code }
  ### Format: {PromptConstants.FORMATTING_MESSAGE}

  ```python
 {question.starter_code}
  ```
  { else }
  ### Format: {PromptConstants.FORMATTING_WITHOUT_STARTER_MESSAGE}

  ```python
 # YOUR CODE HERE
  ```
  { endif }

### Answer: (use the provided format with backticks)

```

Code Generation Prompt

```

{if check_result.result_status is "Wrong Answer"}
The above code is incorrect and does not pass the testcase.
Input: {wrong_testcase_input}
Output: {wrong_testcase_output}
Expected: {wrong_testcase_expected}

{elif check_result.result_status is "Time Limit Exceeded"}
The above code is incorrect and exceeds the time limit.
Input: {wrong_testcase_input}

{elif check_result.result_status is "Runtime Error"}
The above code is incorrect and has a runtime error.
Input: {wrong_testcase_input}
Error Message: {wrong_testcase_error_message}

{endif}

```

Self Repair Error Feedback Pseudocode

1404
 1405
 1406
 1407
 1408
 1409
 1410
 1411
 1412
 1413
 1414
 1415
 1416
 1417
 1418
 1419
 1420
 1421
 1422
 1423
 1424
 1425
 1426
 1427
 1428
 1429
 1430
 1431
 1432
 1433
 1434
 1435
 1436
 1437
 1438
 1439
 1440
 1441
 1442
 1443
 1444
 1445
 1446
 1447
 1448
 1449
 1450
 1451
 1452
 1453
 1454
 1455
 1456
 1457

You are a helpful programming assistant and an expert Python programmer
 ↪ . You are helping a user write a program to solve a problem. The
 ↪ user has written some code, but it has some errors and is not
 ↪ passing the tests. You will help the user by first giving a
 ↪ concise (at most 2-3 sentences) textual explanation of what is
 ↪ wrong with the code. After you have pointed out what is wrong
 ↪ with the code, you will then generate a fixed version of the
 ↪ program. You must put the entire fixed program within code
 ↪ delimiters only for once.

```

### Question:\n{question.question_content}

### Answer: ```python
{code.code_to_be_corrected}
```

Format: {PromptConstants.FORMATTING_CHECK_ERROR_MESSAGE}

Answer: (use the provided format with backticks)

```

Self-Repair Prompt

You are given a Python function and an assertion containing an input to  
 ↪ the function. Complete the assertion with a literal (no  
 ↪ unsimplified expressions, no function calls) containing the  
 ↪ output when executing the provided code on the given input, even  
 ↪ if the function is incorrect or incomplete. Do NOT output any  
 ↪ extra information. Provide the full assertion with the correct  
 ↪ output in [ANSWER] and [/ANSWER] tags, following the examples.

```

[PYTHON]
def repeatNumber(number : int) -> int:
 return number
assert repeatNumber(number = 17) == ??
[/PYTHON]
[ANSWER]
assert repeatNumber(number = 17) == 17
[/ANSWER]

[PYTHON]
def addCharacterA(string : str) -> str:
 return string + "a"
assert addCharacterA(string = "x9j") == ??
[/PYTHON]
[ANSWER]
assert addCharacterA(string = "x9j") == "x9ja"
[/ANSWER]

[PYTHON]
{code}
assert {input} == ??
[/PYTHON]
[ANSWER]

```

Code Execution Prompt

1458  
 1459  
 1460  
 1461  
 1462  
 1463  
 1464  
 1465  
 1466  
 1467  
 1468  
 1469  
 1470  
 1471  
 1472  
 1473  
 1474  
 1475  
 1476  
 1477  
 1478  
 1479  
 1480  
 1481  
 1482  
 1483  
 1484  
 1485  
 1486  
 1487  
 1488  
 1489  
 1490  
 1491  
 1492  
 1493  
 1494  
 1495  
 1496  
 1497  
 1498  
 1499  
 1500  
 1501  
 1502  
 1503  
 1504  
 1505  
 1506  
 1507  
 1508  
 1509  
 1510  
 1511

You are given a Python function and an assertion containing an `input` to  
 ↪ the function. Complete the assertion with a literal (no  
 ↪ unsimplified expressions, no function calls) containing the  
 ↪ output when executing the provided code on the given `input`, even  
 ↪ if the function is incorrect or incomplete. Do NOT output any  
 ↪ extra information. Execute the program step by step before  
 ↪ arriving at an answer, and provide the full assertion with the  
 ↪ correct output in `[ANSWER]` and `[/ANSWER]` tags, following the  
 ↪ examples.

```
[PYTHON]
def performOperation(s):
 s = s + s
 return "b" + s + "a"
assert performOperation(s = "hi") == ??
[/PYTHON]
```

[THOUGHT]  
 Let's execute the code step by step:

1. The function `performOperation` is defined, which takes a single  
 ↪ argument `s`.
2. The function is called with the argument `"hi"`, so within the  
 ↪ function, `s` is initially `"hi"`.
3. Inside the function, `s` is concatenated with itself, so `s` becomes "  
 ↪ `hihi`".
4. The function then returns a new string that starts with `"b"`,  
 ↪ followed by the value of `s` (which is now `"hihi"`), and ends with  
 ↪ `"a"`.
5. The return value of the function is therefore `"bhihia"`.

```
[/THOUGHT]
[ANSWER]
assert performOperation(s = "hi") == "bhihia"
[/ANSWER]
```

```
[PYTHON]
{code}
assert {input} == ??
[/PYTHON]
[THOUGHT]
```

Code Execution Prompt with CoT

1512  
1513  
1514  
1515  
1516  
1517  
1518  
1519  
1520  
1521  
1522  
1523  
1524  
1525  
1526  
1527  
1528  
1529  
1530  
1531  
1532  
1533  
1534  
1535  
1536  
1537  
1538  
1539  
1540  
1541  
1542  
1543  
1544  
1545  
1546  
1547  
1548  
1549  
1550  
1551  
1552  
1553  
1554  
1555  
1556  
1557  
1558  
1559  
1560  
1561  
1562  
1563  
1564  
1565

```
Instruction: You are a helpful programming assistant and an expert
↪ Python programmer. You are helping a user to write a test case
↪ to help to check the correctness of the function. The user has
↪ written a input for the testcase. You will calculate the output
↪ of the testcase and write the whole assertion statement in the
↪ markdown code block with the correct output.
```

Problem:

```
{problem_statement}
```

Function:

```
```\n{function_signature}\n```
```

Please complete the following test case:

```
```\nassert {function_name}({testcase_input}) == # TODO\n```
```

```
Response:
```

Test Output Prediction Prompt

## E RESULTS

### E.1 CONTAMINATION

Figure 10 demonstrates contamination in DEEPSEEK in self repair and test output prediction scenarios.

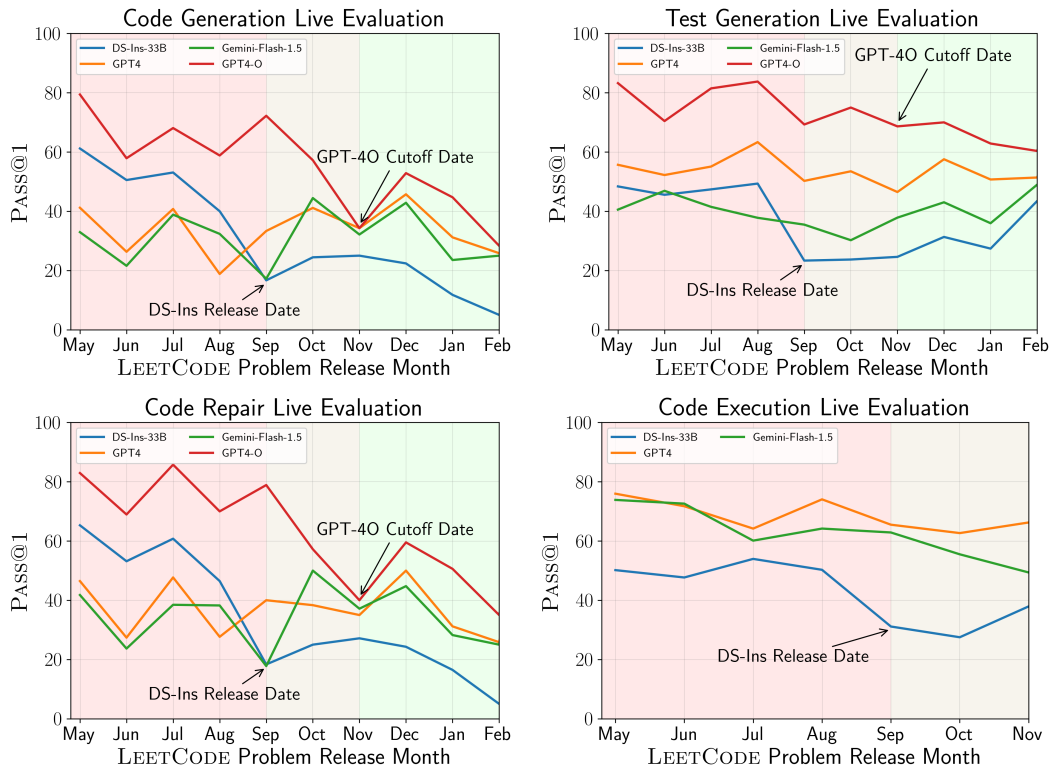


Figure 10: Contamination in DS-B models across self-repair and code execution (without COT) scenarios over time. Note that code execution currently runs between May and November

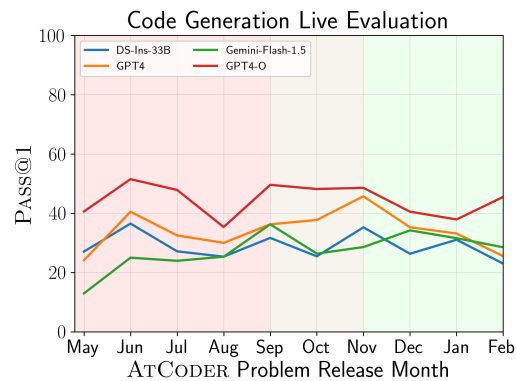


Figure 11: Performance on problems released over different months for ATCODER

### E.2 ALL RESULTS

Below we provide the tables comprising of results across different LIVECODEBENCH scenarios.

**Comparing Closed Models.** We evaluate a range of closed (API access) models ranging from different model families like GPTs, CLAUDES, GEMINI, and MISTRAL. We find the GPT-4-TURBO and

1620  
1621  
1622  
1623  
1624  
1625  
1626  
1627  
1628  
1629  
1630  
1631  
1632  
1633  
1634  
1635  
1636  
1637  
1638  
1639  
1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1649  
1650  
1651  
1652  
1653  
1654  
1655  
1656  
1657  
1658  
1659  
1660  
1661  
1662  
1663  
1664  
1665  
1666  
1667  
1668  
1669  
1670  
1671  
1672  
1673

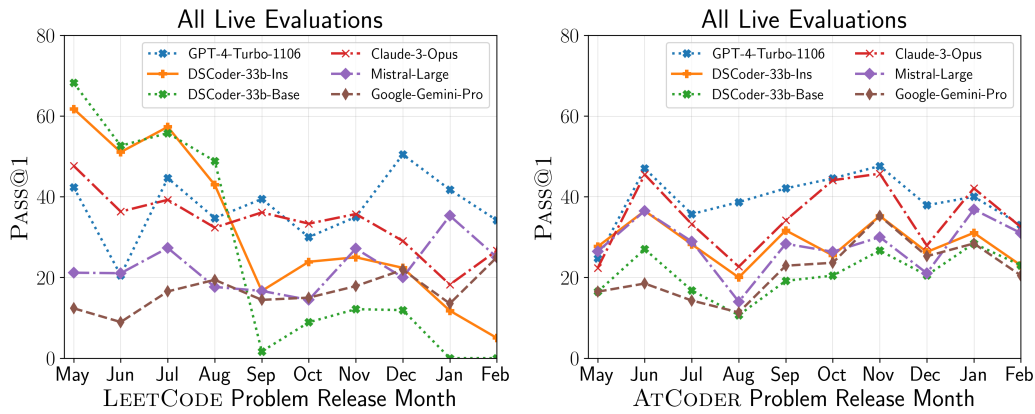


Figure 12: Live evaluation over time for various models on code generation scenario in LIVECODEBENCH. We consider many recently released models and do not find significant performance variations across months except for DS-B models.

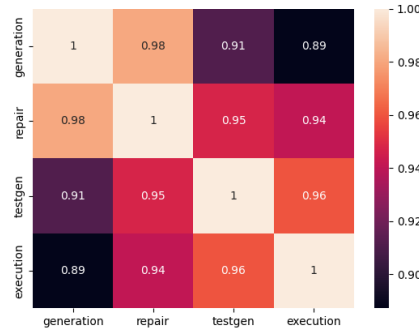


Figure 13: Correlations across different scenarios studied in LIVECODEBENCH

CLAUDE-3-OPUS rank at the top across all scenarios followed by MISTRAL-L and CLAUDE-3-SONNET models. Finally, GEMINI-PRO and GPT-3.5-TURBO lie on the lower end of the models. The relative differences between the models vary across the scenarios. For example, GPT-4-TURBO demonstrates remarkable improvement from self-repair (24.5% to 36.9% on the LCB-Medium problems) while GEMINI-PRO only improves from 8.5% to 9.4%. Similarly, as identified above, CLAUDE-3-OPUS and MISTRAL-L perform considerably better on test output prediction and code execution scenarios.

1674  
1675  
1676  
1677  
1678  
1679  
1680  
1681  
1682  
1683  
1684  
1685  
1686  
1687  
1688  
1689  
1690  
1691  
1692  
1693  
1694  
1695  
1696  
1697  
1698  
1699  
1700  
1701  
1702  
1703  
1704  
1705  
1706  
1707  
1708  
1709  
1710  
1711  
1712  
1713  
1714  
1715  
1716  
1717  
1718  
1719  
1720  
1721  
1722  
1723  
1724  
1725  
1726  
1727

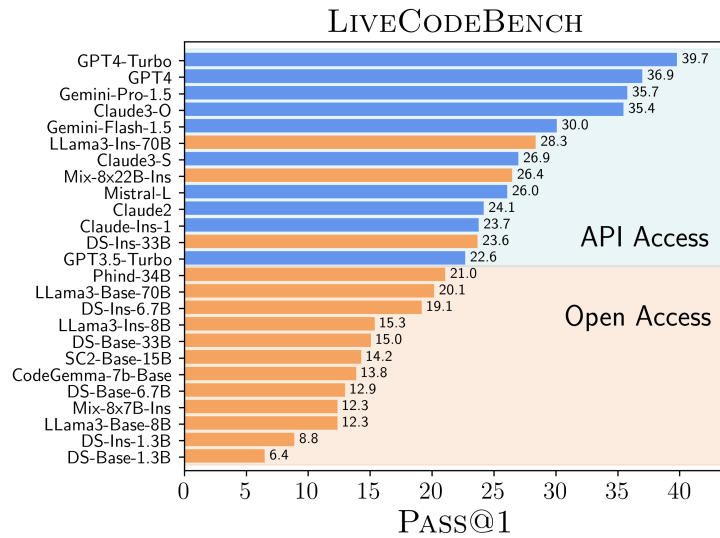


Figure 14: Comparison of open access and (closed) API access models on LIVECODEBENCH-Easy code generation scenario. We find that closed-access models consistently outperform the open models with only strong instruction-tuned variants of > 30B models (specifically L3-INS-70B, MIXTRAL and DS-INS-33B models) crossing the performance gap.



| Model Name                 | Easy  | Medium | Hard | Total |
|----------------------------|-------|--------|------|-------|
| Claude-2                   | 61.80 | 4.90   | 0.20 | 22.30 |
| Claude-3-Haiku             | 63.00 | 4.30   | 1.10 | 22.80 |
| Claude-3-Opus              | 78.80 | 16.30  | 3.20 | 32.80 |
| Claude-3-Sonnet            | 67.60 | 6.20   | 1.10 | 25.00 |
| Claude-Instant-1           | 60.70 | 4.30   | 1.10 | 22.10 |
| CodeGemma-2b-Base          | 18.30 | 0.40   | 0.00 | 6.30  |
| CodeGemma-7b-Base          | 35.70 | 2.60   | 0.10 | 12.80 |
| CodeLlama-13b-Base         | 24.60 | 0.90   | 0.00 | 8.50  |
| CodeLlama-13b-Ins          | 36.60 | 2.40   | 0.00 | 13.00 |
| CodeLlama-34b-Base         | 32.20 | 1.80   | 0.10 | 11.40 |
| CodeLlama-34b-Ins          | 33.70 | 2.40   | 1.10 | 12.40 |
| CodeLlama-70b-Base         | 15.80 | 1.20   | 0.00 | 5.70  |
| CodeLlama-70b-Ins          | 7.80  | 0.60   | 0.00 | 2.80  |
| CodeLlama-7b-Base          | 19.00 | 0.40   | 0.00 | 6.50  |
| CodeLlama-7b-Ins           | 28.60 | 2.50   | 0.00 | 10.40 |
| CodeQwen1.5-7B             | 40.40 | 4.80   | 0.00 | 15.10 |
| CodeQwen1.5-7B-Chat        | 39.20 | 13.10  | 0.50 | 17.60 |
| Codestral-Latest           | 69.00 | 18.70  | 0.90 | 29.50 |
| Command-R                  | 39.00 | 3.60   | 0.00 | 14.20 |
| Command-R+                 | 56.60 | 6.80   | 0.00 | 21.10 |
| DSCoder-1.3b-Base          | 17.30 | 0.70   | 0.00 | 6.00  |
| DSCoder-1.3b-Ins           | 22.90 | 1.50   | 0.00 | 8.10  |
| DSCoder-33b-Base           | 39.40 | 2.30   | 0.00 | 13.90 |
| DSCoder-33b-Ins            | 55.60 | 9.00   | 0.70 | 21.80 |
| DSCoder-6.7b-Base          | 34.30 | 1.40   | 0.10 | 11.90 |
| DSCoder-6.7b-Ins           | 46.40 | 5.80   | 0.70 | 17.60 |
| GPT-3.5-Turbo-0125         | 56.80 | 10.80  | 0.10 | 22.60 |
| GPT-3.5-Turbo-0301         | 53.40 | 8.80   | 0.20 | 20.80 |
| GPT-4-0613                 | 78.40 | 21.20  | 2.30 | 33.90 |
| GPT-4-Turbo-1106           | 84.40 | 24.00  | 0.50 | 36.30 |
| GPT-4-Turbo-2024-04-09     | 85.30 | 33.00  | 5.10 | 41.10 |
| GPT-4o-2024-05-13          | 88.30 | 33.20  | 4.20 | 41.90 |
| Gemini-Flash-1.5-May       | 68.10 | 12.60  | 2.70 | 27.80 |
| Gemini-Pro-1.5-April (n=1) | 56.50 | 14.30  | 3.60 | 24.80 |
| Gemini-Pro-1.5-May         | 76.00 | 19.40  | 3.50 | 33.00 |
| Gemma-2b-Base              | 6.10  | 0.00   | 0.00 | 2.00  |
| Gemma-7b-Base              | 27.00 | 0.90   | 0.00 | 9.30  |
| LLama3-70b-Base            | 52.20 | 3.20   | 0.60 | 18.60 |
| LLama3-70b-Ins             | 60.70 | 15.80  | 1.40 | 26.00 |
| LLama3-8b-Base             | 32.90 | 1.50   | 0.00 | 11.50 |
| LLama3-8b-Ins              | 38.60 | 3.50   | 0.50 | 14.20 |
| MagiCoderS-CL-7B           | 32.80 | 2.40   | 0.00 | 11.70 |
| MagiCoderS-DS-6.7B         | 49.20 | 7.50   | 0.00 | 18.90 |
| Mistral-Large              | 60.20 | 10.90  | 0.90 | 24.00 |
| Mixtral-8x22B-Ins          | 59.80 | 12.70  | 0.00 | 24.20 |
| Mixtral-8x7B-Ins           | 31.60 | 2.60   | 0.00 | 11.40 |
| OC-DS-1.3B                 | 11.30 | 0.10   | 0.00 | 3.80  |
| OC-DS-33B                  | 53.90 | 5.10   | 0.00 | 19.70 |
| OC-DS-6.7B                 | 46.30 | 4.50   | 0.00 | 16.90 |
| Phind-34B-V2               | 53.40 | 4.70   | 0.10 | 19.40 |
| StarCoder2-15b             | 37.30 | 2.20   | 0.00 | 13.20 |
| StarCoder2-3b              | 28.20 | 0.70   | 0.00 | 9.60  |
| StarCoder2-7b              | 29.90 | 1.20   | 0.00 | 10.40 |

Table 4: Code Generation Performances

1782  
 1783  
 1784  
 1785  
 1786  
 1787  
 1788  
 1789  
 1790  
 1791  
 1792  
 1793  
 1794  
 1795  
 1796  
 1797  
 1798  
 1799  
 1800  
 1801  
 1802  
 1803  
 1804  
 1805  
 1806  
 1807  
 1808  
 1809  
 1810  
 1811  
 1812  
 1813  
 1814  
 1815  
 1816  
 1817  
 1818  
 1819  
 1820  
 1821  
 1822  
 1823  
 1824  
 1825  
 1826  
 1827  
 1828  
 1829  
 1830  
 1831  
 1832  
 1833  
 1834  
 1835

| Model Name                 | Easy  | Medium | Hard | Total |
|----------------------------|-------|--------|------|-------|
| Claude-2                   | 66.20 | 10.30  | 0.40 | 25.60 |
| Claude-3-Haiku             | 66.50 | 8.70   | 2.50 | 25.90 |
| Claude-3-Opus              | 83.10 | 23.70  | 6.70 | 37.80 |
| Claude-3-Sonnet            | 72.60 | 11.80  | 2.20 | 28.90 |
| Claude-Instant-1           | 64.40 | 7.10   | 2.20 | 24.60 |
| CodeLlama-13b-Ins          | 43.10 | 3.00   | 0.00 | 15.30 |
| CodeLlama-34b-Ins          | 31.50 | 3.50   | 1.80 | 12.30 |
| CodeLlama-7b-Ins           | 31.90 | 3.10   | 1.50 | 12.10 |
| Codestral-Latest           | 72.50 | 25.90  | 3.30 | 33.90 |
| DSCoder-1.3b-Ins           | 29.50 | 2.10   | 0.00 | 10.60 |
| DSCoder-33b-Ins            | 60.70 | 8.10   | 1.50 | 23.40 |
| DSCoder-6.7b-Ins           | 49.90 | 5.70   | 1.10 | 18.90 |
| GPT-3.5-Turbo-0125         | 59.30 | 11.90  | 0.50 | 23.90 |
| GPT-3.5-Turbo-0301         | 58.40 | 11.60  | 0.70 | 23.60 |
| GPT-4-0613                 | 79.30 | 25.00  | 2.40 | 35.60 |
| GPT-4-Turbo-1106           | 86.90 | 36.90  | 4.00 | 42.60 |
| GPT-4-Turbo-2024-04-09     | 88.70 | 39.70  | 8.40 | 45.60 |
| GPT-4O-2024-05-13          | 92.60 | 46.40  | 8.20 | 49.10 |
| Gemini-Flash-1.5-May       | 73.40 | 16.40  | 4.40 | 31.40 |
| Gemini-Pro                 | 53.80 | 9.40   | 0.20 | 21.10 |
| Gemini-Pro-1.5-April (n=1) | 71.80 | 19.40  | 5.50 | 32.20 |
| Gemini-Pro-1.5-May         | 84.80 | 30.10  | 7.30 | 40.70 |
| LLama3-70b-Ins             | 69.60 | 19.00  | 1.80 | 30.10 |
| LLama3-8b-Ins              | 47.10 | 6.10   | 0.00 | 17.70 |
| MagiCoderS-CL-7B           | 36.50 | 3.10   | 0.00 | 13.20 |
| MagiCoderS-DS-6.7B         | 50.60 | 8.60   | 0.00 | 19.70 |
| Mistral-Large              | 71.20 | 15.60  | 3.60 | 30.10 |
| OC-DS-1.3B                 | 20.00 | 0.40   | 0.00 | 6.80  |
| OC-DS-33B                  | 58.90 | 7.20   | 1.30 | 22.50 |
| OC-DS-6.7B                 | 50.90 | 6.30   | 0.20 | 19.10 |
| Phind-34B-V2               | 62.00 | 6.50   | 0.90 | 23.10 |

Table 5: Self Repair Performances

1836  
1837  
1838  
1839  
1840  
1841  
1842  
1843  
1844  
1845  
1846  
1847  
1848  
1849  
1850  
1851  
1852  
1853  
1854  
1855  
1856  
1857  
1858  
1859  
1860  
1861  
1862  
1863  
1864  
1865  
1866  
1867  
1868  
1869  
1870  
1871  
1872  
1873  
1874  
1875  
1876  
1877  
1878  
1879  
1880  
1881  
1882  
1883  
1884  
1885  
1886  
1887  
1888  
1889

| Model Name                 | Pass@1 |
|----------------------------|--------|
| Claude-2                   | 32.70  |
| Claude-3-Haiku             | 32.90  |
| Claude-3-Opus              | 58.70  |
| Claude-3-Sonnet            | 34.10  |
| Claude-Instant-1           | 25.40  |
| CodeLlama-13b-Ins          | 24.40  |
| CodeLlama-34b-Ins          | 23.00  |
| CodeLlama-70b-Ins          | 16.10  |
| CodeLlama-7b-Ins           | 15.30  |
| Codestral-Latest           | 41.80  |
| DSCoder-1.3b-Ins           | 12.50  |
| DSCoder-33b-Ins            | 28.30  |
| DSCoder-6.7b-Ins           | 26.50  |
| GPT-3.5-Turbo-0125         | 35.40  |
| GPT-3.5-Turbo-0301         | 32.50  |
| GPT-4-0613                 | 52.90  |
| GPT-4-Turbo-1106           | 55.70  |
| GPT-4-Turbo-2024-04-09     | 66.10  |
| GPT-4O-2024-05-13          | 68.90  |
| Gemini-Flash-1.5-May       | 38.10  |
| Gemini-Pro                 | 29.50  |
| Gemini-Pro-1.5-April (n=1) | 49.60  |
| Gemini-Pro-1.5-May         | 44.80  |
| LLama3-70b-Ins             | 41.40  |
| LLama3-8b-Ins              | 24.40  |
| MagiCoderS-CL-7B           | 21.30  |
| MagiCoderS-DS-6.7B         | 27.10  |
| Mistral-Large              | 46.50  |
| Mixtral-8x22B-Ins          | 44.70  |
| Mixtral-8x7B-Ins           | 31.80  |
| OC-DS-1.3B                 | 7.80   |
| OC-DS-33B                  | 11.30  |
| OC-DS-6.7B                 | 18.30  |
| Phind-34B-V2               | 27.20  |

Table 6: Test Output Prediction Performances

1890  
1891  
1892  
1893  
1894  
1895  
1896  
1897  
1898  
1899  
1900  
1901  
1902  
1903  
1904  
1905  
1906  
1907  
1908  
1909  
1910  
1911  
1912  
1913  
1914  
1915  
1916  
1917  
1918  
1919  
1920  
1921  
1922  
1923  
1924  
1925  
1926  
1927  
1928  
1929  
1930  
1931  
1932  
1933  
1934  
1935  
1936  
1937  
1938  
1939  
1940  
1941  
1942  
1943

| Model Name                   | Pass@1 | Pass@1 (COT) |
|------------------------------|--------|--------------|
| Claude-2                     | 31.50  | 43.80        |
| Claude-3-Haiku               | 0.70   | 28.30        |
| Claude-3-Opus                | 36.50  | 80.10        |
| Claude-3-Sonnet              | 29.30  | 39.40        |
| Claude-Instant-1             | 20.00  | 34.80        |
| ClLama-13b-Ins               | 23.50  | 14.10        |
| ClLama-34b-Ins               | 28.90  | 24.50        |
| ClLama-7b-Ins                | 20.60  | 14.20        |
| CodeLlama-70b-Ins            | 31.20  | -1.00        |
| Codestral-Latest             | 37.90  | 41.80        |
| DSCoder-1.3b-Base            | 19.00  | 13.40        |
| DSCoder-1.3b-Ins             | 18.10  | 17.00        |
| DSCoder-33b-Base             | 29.90  | 29.10        |
| DSCoder-33b-Ins              | 26.60  | 31.70        |
| DSCoder-6.7b-Base            | 23.50  | 25.10        |
| DSCoder-6.7b-Ins             | 23.10  | 23.80        |
| GPT-3.5-Turbo-0301           | 33.90  | 34.80        |
| GPT-4-0613                   | 44.30  | 64.80        |
| GPT-4-Turbo-1106             | 40.50  | 83.60        |
| GPT-4-Turbo-2024-04-09       | 45.90  | 83.80        |
| GPT-4O-2024-05-13            | 39.10  | 91.00        |
| Gemini-Flash-1.5-May         | 21.40  | 57.10        |
| Gemini-Pro                   | 27.70  | 37.40        |
| Gemini-Pro-1.5 (April) (n=1) | 30.30  | 64.40        |
| Gemini-Pro-1.5-May           | 42.10  | 72.10        |
| LLama3-70b-Ins               | 29.60  | 55.50        |
| LLama3-8b-Ins                | 18.40  | 29.40        |
| MagiCoderS-CL-7B             | 21.20  | -1.00        |
| MagiCoderS-DS-6.7B           | 27.20  | -1.00        |
| Mistral-Large                | 36.60  | 54.40        |
| Phind-34B-V2                 | 26.90  | -1.00        |
| StarCoder                    | 20.30  | -1.00        |
| WCoder-34B-V1                | 28.40  | -1.00        |

Table 7: Code Execution Performances

## 1944 F QUALITATIVE EXAMPLES

1945

## 1946 F.1 CODE EXECUTION

1947

1948 We show 5 examples from the code execution task that GPT-4 (gpt-4-1106-preview) still  
 1949 struggles to execute, even with CoT.

1950

```
1951 def countWays(nums: List[int]) -> int:
1952 nums.sort()
1953 n = len(nums)
1954 ans = 0
1955 for i in range(n + 1):
1956 if i and nums[i-1] >= i: continue
1957 if i < n and nums[i] <= i: continue
1958 ans += 1
1959 return ans
1960 assert countWays(nums = [6, 0, 3, 3, 6, 7, 2, 7]) == 3
1961 # GPT-4 + CoT Outputs: 1, 2, 4, 5
```

1962

Mistake 1

1963

1964

```
1965 def minimumCoins(prices: List[int]) -> int:
1966
1967 @cache
1968 def dfs(i, free_until):
1969 if i >= len(prices):
1970 return 0
1971
1972 res = prices[i] + dfs(i + 1, min(len(prices) - 1, i + i + 1))
1973
1974 if free_until >= i:
1975 res = min(res, dfs(i + 1, free_until))
1976
1977 return res
1978
1979 dfs.cache_clear()
1980 return dfs(0, -1)
1981 assert minimumCoins(prices = [3, 1, 2]) == 4
1982 # GPT-4 + CoT Outputs: 1, 3, 5, 6
```

1983

Mistake 2

1984

1985

```
1986 def sortVowels(s: str) -> str:
1987 q = deque(sorted((ch for ch in s if vowel(ch))))
1988 res = []
1989 for ch in s:
1990 if vowel(ch):
1991 res.append(q.popleft())
1992 else:
1993 res.append(ch)
1994 return ''.join(res)
1995 assert sortVowels(s = 'lEetcOde') == 'lEOtcede'
1996 # GPT-4 + CoT Outputs: "leetecode", "lEetecOde", "leetcede", "leetcEde"
1997 ↔ ", "leetcOde"
```

1998

Mistake 3

1999

1998  
1999  
2000  
2001  
2002  
2003  
2004  
2005  
2006  
2007  
2008  
2009  
2010  
2011  
2012  
2013  
2014  
2015  
2016  
2017  
2018  
2019  
2020  
2021  
2022  
2023  
2024  
2025  
2026  
2027  
2028  
2029  
2030  
2031  
2032  
2033  
2034  
2035  
2036  
2037  
2038  
2039  
2040  
2041  
2042  
2043  
2044  
2045  
2046  
2047  
2048  
2049  
2050  
2051

```
def relocateMarbles(nums: List[int], moveFrom: List[int], moveTo: List[
 ↪ int]) -> List[int]:
 nums = sorted(list(set(nums)))
 dd = {}
 for item in nums:
 dd[item] = 1
 for a,b in zip(moveFrom, moveTo):
 del dd[a]
 dd[b] = 1
 ll = dd.keys()
 return sorted(ll)
assert relocateMarbles(nums = [1, 6, 7, 8], moveFrom = [1, 7, 2],
 ↪ moveTo = [2, 9, 5]) == [5, 6, 8, 9]
GPT-4 + CoT Outputs: [2, 6, 8, 9], [2, 5, 6, 8, 9], KeyError
```

Mistake 4

```
def minimumSum(nums: List[int]) -> int:
 left, right, ans = [inf], [inf], inf
 for num in nums:
 left.append(min(left[-1], num))
 for num in nums[::-1]:
 right.append(min(right[-1], num))
 right.reverse()
 for i, num in enumerate(nums):
 if left[i] < num and right[i + 1] < num:
 ans = min(ans, num + left[i] + right[i + 1])
 return ans if ans < inf else -1
assert minimumSum(nums = [6, 5, 4, 3, 4, 5]) == -1
GPT-4 + CoT Outputs: 10, 11, 12
```

Mistake 5

## G LIMITATIONS

**Benchmark Size.** LIVECODEBENCH code generation scenario currently hosts over 500 instances from problems released since May 2023. To account for contamination in DEEPSEEK, we only perform evaluations on problems released after the model cutoff date. This leads to only 349 problems used in our final evaluations which might add noise due to problem set samples. We currently estimate 1 – 2% performance variations in LCB code generation due to this issue (measured by bootstrapping 349 sized problem sets from the 349 sized dataset). Other scenarios, i.e. self-repair, code execution, and test output prediction comprise 238, 188, and 254 problems would have similar performance variations. We thus recommend exercising proper judgement when comparing models with small performance differences. Note that HUMANEVAL has 164 problems and would also struggle with similar issues.

This issue is also exacerbated for newer models, with more recent cutoff dates, as they might only have access to a smaller evaluation set. We propose two solutions addressing this issue as we evolve LIVECODEBENCH. First, we will use other competition platforms for problem collection, allowing larger number of recent problems to be added to the benchmark. In addition, we also hope supplement this with an unreleased private test set constructed specifically for model evaluation. These problems will use a similar flavor to current problems and will be used when models are submitted for evaluation to the LIVECODEBENCH platform. This would reduce the reliance on public accessible problems and provide a more robust evaluation of the models while providing community public access to similar problems, similar to strategies employed by popular platforms like KAGGLE.

**Focus on PYTHON.** LIVECODEBENCH currently only focuses on PYTHON which might not provide enough signal about model capabilities in other languages. However, since we collected problem statements and serialized tests, adding new programming languages would be straightforward once appropriate evaluation engines are used.

**Robustness to Prompts.** Recent works have identified huge performance variances that can be caused due to insufficient prompt. Here, we either do not tune prompts across models or make minor adjustments based on the system prompts and delimiter tokens. This can lead to performance variance in our results. Our findings and model comparison orders generalize across LIVECODEBENCH scenarios and mostly match the performance trends observed on HUMANEVAL making this a less prominent issue.

This issue can be particularly observed open models on the code execution scenario with COT prompting. Interestingly, often the open models perform even worse in comparison to the direct code execution baseline. Note that we used same prompts for the closed models all of which show noticeable improvement from COT. While the used prompts might be sub-optimal, this highlights how open-models perform worse against the closed models at performing chain-of-thought.

**Problem Domain.** Programming is a vast domain and occurs in various forms such as programming puzzles, competition programming, and real-world software development. Different domains might have individual requirements, constraints, challenges, and difficulty levels. LIVECODEBENCH currently focuses on competition problems sourced from three platforms. This might not be representative of the “most general” notion of LLM programming capabilities. Particularly, real-world usage of LLMs is drawn upon open-ended and unconstrained problems rasied by users. We therefore recommend using LIVECODEBENCH as a starting point for evaluating LLMs and further using domain-specific evaluations to measure and compare LLMs in specific settings as required.