# INFOSYNTH: INFORMATION-GUIDED BENCHMARK SYNTHESIS FOR LLMS

**Anonymous authors** 

Paper under double-blind review

#### **ABSTRACT**

Large language models (LLMs) have demonstrated significant advancements in reasoning and code generation. However, efficiently creating new benchmarks to evaluate these capabilities remains a challenge. Traditional benchmark creation relies on manual human effort, a process that is both expensive and timeconsuming. Furthermore, existing benchmarks often contaminate LLM training data, necessitating novel and diverse benchmarks to accurately assess their genuine capabilities. This work introduces InfoSynth, a novel framework for automatically generating and evaluating reasoning benchmarks guided by informationtheoretic principles. We propose metrics based on KL-divergence and entropy to quantify benchmark novelty and diversity without relying on costly model evaluations. Building on this framework, we develop an end-to-end pipeline that synthesizes robust Python coding problems from seed datasets using genetic algorithms and iterative code feedback. Our method generates accurate test cases and solutions to new problems 97% of the time, and the synthesized benchmarks consistently exhibit higher novelty and diversity compared to their seed datasets. Moreover, our algorithm provides a method for controlling the novelty/diversity and difficulty of generated problems. InfoSynth offers a scalable, self-verifying pipeline for constructing high-quality, novel and diverse benchmarks for LLMs.

# 1 Introduction

Large language models (LLMs) have demonstrated impressive capabilities in code generation and reasoning. However, rigorously evaluating these reasoning abilities remains a significant challenge. While substantial effort has been invested in creating robust math and coding benchmarks (Austin et al., 2021; Chen et al., 2021; Cobbe et al., 2021; Jain et al., 2024; Liu et al., 2024a; Zhuo et al., 2024), their development often demands considerable human labor or extensive computational resources for problem and solution validation. Some existing approaches utilize a judge LLM to generate and verify new problems (Ding et al., 2024; Li et al., 2025b; Majumdar et al., 2024). However, this method can yield erroneous benchmarks, as the judge LLM's capacity to reliably solve the generated problems is not guaranteed. This paper focuses on Python coding problems, whose solutions can be verified through execution in a code environment. Our novel pipeline leverages this executability to ensure the robustness of the generated problems.

Beyond the challenge of ensuring robustness, state-of-the-art (SOTA) reasoning models often overfit to their training data, leading to poor performance on out-of-distribution problems (Huang et al., 2025). Furthermore, recent studies reveal that LLM training data is frequently contaminated by existing evaluation benchmarks, which can artificially inflate reported performance (Deng et al., 2023a;b; Golchin & Surdeanu, 2023). For instance, Zhang et al. (2024) show that LLMs experience accuracy drops of up to 8% on their novel GSM1k dataset, despite its similarity in difficulty to the widely used GSM8k. This underscores the critical need for new, contamination-free reasoning benchmarks to genuinely assess the capabilities of LLMs.

To address these pressing issues, our work emphasizes two crucial benchmark properties: *novelty* and *diversity*. While this work does not directly address the task of creating contamination-free benchmarks, we provide an improved method of generating benchmarks that cover more diverse and novel coding tasks. A novel benchmark should comprise problems distinct from existing datasets, thereby preventing models from achieving high scores through mere memorization of previously

seen examples. Conversely, a diverse benchmark should encompass a broad spectrum of dissimilar problems, enhancing its resilience against model overfitting and providing a more comprehensive evaluation. Clearly, robust, novel, and diverse benchmarks are essential for the reliable evaluation of LLM reasoning abilities. Our work seeks to answer two fundamental questions: (1) How can we effectively measure the novelty and diversity of benchmarks? (2) How can we efficiently generate benchmarks that possess these desirable properties while ensuring their correctness and robustness?

Our main contributions can be summarized as:

- We introduce an information-theoretic framework to quantify and compare the *novelty* and *diver*sity of benchmarks, offering a principled approach to benchmark assessment without reliance on model evaluations.
- We propose and validate an end-to-end pipeline, InfoSynth, for efficiently synthesizing novel, diverse, and verifiably correct Python coding problems from seed datasets with genetic algorithms and iterative code feedback.
- Through extensive experiments, we demonstrate that InfoSynth exhibits superior robustness compared to existing data generation methods. Our pipeline provides a method for increasing the novelty/diversity of generated problems and controlling their difficulty.

#### 2 Related Work

**Synthetic Problem Generation.** Previous work has explored generating novel synthetic datasets from high-quality seed benchmarks. Majumdar et al. (2024); Wang et al. (2023); Xu et al. (2024); Zhao et al. (2025) show that LLMs can generate new instructions from existing ones. Chen et al. (2022); Liu et al. (2023; 2024b); Xu et al. (2025); Zeng et al. (2025) successfully used LLMs to generate unit tests for solution verification. Our end-to-end pipeline extends existing methods with code-execution environments to ensure robustness, novelty, and diversity.

**Benchmark Quality Assessment.** Efficient, concrete analysis of benchmark quality remains an open problem. Prior work defines metrics for novelty, separability, and difficulty via test-taker performance (Li et al., 2025a;b), proposes adaptive selection of novel problems to reduce evaluation cost (Truong et al., 2025), and develops similarity and difficulty scores for coding tasks (Tambon et al., 2024). A key limitation is reliance on SOTA model performance, making these methods accurate but computationally expensive. Our approach analyzes novelty and diversity more efficiently, without requiring costly model evaluations.

# 3 DESIRABLE BENCHMARK PROPERTIES

We propose a framework for characterizing the novelty and diversity of benchmarks. Our new novelty metric uses the KL-Divergence to capture how different the benchmark is from existing datasets, with the broader goal of creating benchmarks that are contamination-free. Previous work has used the KL-Divergence in a similar way; Schulman et al. (2017) uses it to measure differences in reinforcement learning policies and Kingma & Welling (2014) uses it to regularize output distributions for Variational Auto Encoders. Similarly, our proposed diversity metric uses Shannon entropy to capture how much variety exists among the problems; more diverse datasets provide a broader characterization of an LLMs reasoning abilities.

#### 3.1 AN INFORMATION-THEORY BASED FRAMEWORK FOR BENCHMARK ANALYSIS

Formally, a baseline dataset can be modeled as samples  $X = \{x_i\} \subseteq \mathbb{R}^d$  drawn from some true distribution p(x), and the new dataset that we want to compare against the baseline can be modeled as samples  $Y = \{y_i\} \subseteq \mathbb{R}^d$  drawn from a distribution q(x). Here,  $x_i, y_i$  represent the embedding vectors of the problem statements in an embedding space  $\mathbb{R}^d$ . We define the **novelty** of the new dataset Y to be the KL-divergence between the distributions

Novelty
$$(Y|X) = D_{KL}(q \parallel p) = \int_{\mathbb{R}^d} q(\mathbf{x}) \log \frac{q(\mathbf{x})}{p(\mathbf{x})} d\mathbf{x}$$
 (1)

Note that we take the KL-divergence of p with q as the null hypothesis because we want to reward datasets where q(x) is large and p(x) is small, indicating that the dataset contains problems not in

the distribution of the seed dataset. Given a dataset  $X = \{x_i\}, x_i \sim p(x)$ , we define the **diversity** of the dataset to be the differential entropy of its distribution.

Diversity(X) = 
$$-\int_{\mathbb{R}^d} p(\mathbf{x}) \log p(\mathbf{x}) d\mathbf{x}$$
 (2)

Intuitively, the KL-Divergence captures the fact that novel datasets should have different embeddings from existing datasets. Similarly, diverse datasets should have embeddings that are fairly spread out; clusters indicate problems that are likely to be similar and not diverse. A "perfectly diverse" dataset should resemble a uniform distribution over the embedding space so that it covers a large class of problems. Since the uniform distribution is entropy-maximizing, our metric captures this intuitive characterization of diversity.

In practice, obtaining the full distribution of the embedding space is intractable. Instead, we use statistical estimators for the KL-Divergence and differential entropy. Given samples  $x_1, ..., x_n, y_1, ..., y_m \in \mathbb{R}^d$  where x and y are drawn from p(x), q(x) respectively, we use the k-NN based estimator by Wang et al. (2009)

$$D_{KL}(q||p) = \frac{d}{m} \sum_{i=1}^{m} \log \frac{\nu_k(i)}{\rho_k(i)} + \log \frac{n}{m-1}$$
 (3)

where  $\nu_k(i)$  is the distance from  $y_i$  to its k-th nearest neighbor in  $\{x_j\}$  and  $\rho_k(i)$  is the distance from  $y_i$  to its k-th nearest neighbor in  $\{y_j \mid j \neq i\}$ ; k is a hyperparameter.

Similarly, we can estimate the differential entropy of a dataset. Given samples  $x_1, ..., x_N \in \mathbb{R}^d$ , the Kozachenko-Leonenko estimator is

$$h(X) = \psi(N) - \psi(k) + \log V_d + \frac{d}{N} \sum_{i=1}^{N} \log \rho_k(i)$$
 (4)

where  $\psi$  is the digamma function,  $V_d$  is the volume of the unit ball in  $\mathbb{R}^d$ , and  $\rho_k(i)$  is the distance between  $x_i$  and its k-th nearest neighbor in  $\{x_j \mid j \neq i\}$ ; k is a hyperparameter.

Computing the embeddings, KL-divergence, and entropy for a text dataset is significantly faster and cheaper than computing test-taker statistics. Hence, we provide a way to cheaply estimate the quality of new benchmarks. Moreover, using these metrics, algorithm development can be formulated as an optimization problem that tries to maximize the novelty and diversity of the new dataset.

#### 3.2 EMPIRICAL VALIDATION

In this section, we empirically verify the correctness of our metrics on existing datasets. We use all-mpnet-base-v2 (Song et al., 2020) to embed questions in  $\mathbb{R}^{768}$ . Because estimating entropy in high dimensions is difficult due to the curse of dimensionality, we project down to a lower dimension using UMAP (McInnes et al., 2018). Note that any two datasets that we want to compare must be projected down in the same UMAP call, preserve their relative geometry. We renormalize embeddings after projection so that the distance between embeddings corresponds to cosine similarity.

#### 3.2.1 KL-DIVERGENCE METRIC VALIDATION

We use a dataset of 3511 Leetcode problems with concept labels for every problem (kaysss, 2025). We extract three smaller datasets from this: problems tagged "Hash Map" (686 problems), problems tagged "Graph" (160 problems), and problems tagged "String" (786 problems). We also use the MBPP test dataset (500 problems). We run UMAP with 80 nearest neighbors and a minimum distance of 0.1; we compute this over 10 independent UMAP runs using k=4 for the novelty k-NN parameter. The plots in Figure 1 show 95% confidence intervals.

As shown in Figure 1, our results align with intuition, confirming KL divergence as a measure of benchmark novelty. In the second graph, the estimator becomes negative despite KL being theoretically nonnegative. This is because we are comparing a subset against a superset: subset—superset distances are small, but intra-subset distances are large causing the estimator to be negative. This is not an issue, as we only care about relative differences. For well-generated datasets, negativity should not occur since the new set should not be a strict subset of the seed; we include this case to illustrate that our metric still reflects intuition.

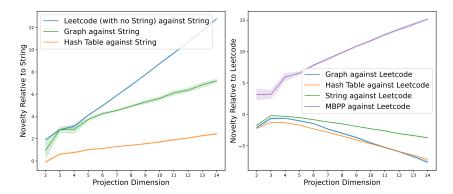


Figure 1: The left plot shows that the full Leetcode dataset has higher novelty than its Hash Table and String subsets as expected. The right plot shows that the MBPP dataset has high novelty against the Leetcode dataset, whereas the Leetcode subdatasets have very little relative novelty as expected.

#### 3.2.2 DIFFERENTIAL ENTROPY METRIC VALIDATION

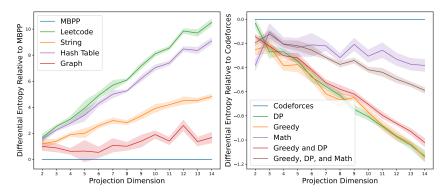


Figure 2: Left: Leetcode vs. MBPP entropy. MBPP shows lower entropy due to simpler, more repetitive problems. Right: Codeforces vs. subsets. Full datasets have higher entropy than topic-specific subsets, except Math, which overlaps with others (e.g., DP, Greedy) and thus appears highly dispersed when isolated. All plots show 95% confidence intervals for the variance.

In addition to prior datasets, we analyze 4,000 random Codeforces problems (open-r1, 2024). UMAP is run with 80 neighbors and min. distance 0.1 over 10 trials. Because the Kozachenko–Leonenko estimator depends on k-NN distances, larger datasets yield artificially lower entropy due to smaller inter-point gaps. To compare fairly, we sample a fixed N points per dataset (without replacement) and average entropy across iterations to reduce variance.

For Leetcode vs. MBPP we use  $N{=}150$  points per dataset over 250 trials ( $k{=}4$ ), and for Codeforces  $N{=}800$  over 250 trials ( $k{=}21$ ) (Figure 2). "Diversity relative to X" is plotted as a difference for visualization only; diversity itself is a unary function. Results show entropy aligns with intuition: datasets expected to be more diverse exhibit higher entropy.

#### 3.3 Choosing k and d

Kraskov et al. (2004) show that the bias-variance tradeoff for dataset size N depends on k/N: larger k increases bias but reduces variance and captures global structure. For KL divergence, we use  $k \approx 4$  to emphasize local differences; for entropy, larger k better captures global diversity, especially with scattered clusters; we find  $k/N \in [0.02, 0.04]$  effective. Overall, diversity and entropy rankings are consistent across dimensions, and we recommend projecting to  $d \in [8, 12]$  for dataset comparison.

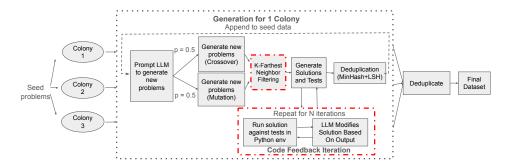


Figure 3: Generation Pipeline. Each colony receives a subset of the seed problems and applies mutation or crossover to them at each iteration. For each problem, it generates solutions and tests which go through multiple iterations of testing to ensure correctness. Deduplication removes similar problems within each colony; the remaining ones are used as seed data for the next iteration. The colony outputs are merged and deduplicated to produce the final dataset.

# 4 A NOVEL BENCHMARK SYNTHESIS PIPELINE

We introduce a novel end-to-end genetic algorithm for generating novel and diverse datasets from a seed dataset. The detailed algorithm is given in Appendix A, but we describe the main ideas here. Figure 3 provides a visual outline of the pipeline.

#### 4.1 Data Generation Pipeline

**Mutation and Crossover.** Starting from the seed data, at each iteration, we randomly apply either crossover or mutation to generate new coding instructions (Majumdar et al., 2024). One key change from previous work is that our mutation prompts ask the model to modify an existing problem in three different difficulty variations: easier, equally difficult, and more difficult to encourage diversity. We demonstrate the benefit of having multiple mutation difficulties in 5.2. Crossover prompts combine existing questions into new ones. Prompts are given in Appendix C.1, C.2. This example shows how mutation varies problem difficulty. Appendix E contains more examples.

**Seed Question:** Write a Python function to find the sum of an array.

**Hard Mutation Variant:** Write a Python function to find the sum of an array, where the array may contain nested lists of integers at any depth.

This example shows how crossover creates an interesting, novel question by combining two unrelated ones. Appendix E contains more examples.

**Seed Questions:** 1. Write a function to rotate a given list by a specified number of items to the right direction. 2. Write a function to find the maximum sum that can be formed which has no three consecutive elements present.

**Crossover Variant:** Write a function to rotate a list by a specified number of steps to the right, ensuring that the sum of any three consecutive elements in the newly rotated list does not exceed a given threshold.

**k-Farthest Neighbor Selection.** A key improvement of InfoSynth is that in order to increase novelty and diversity, we filter problems by cosine similarity to those already generated. In mutation, we produce easy, medium, and hard variants, retaining the two of three with lowest similarity to the seed and generated set. In crossover, we likewise generate three problems and keep the two least similar to the dataset.

**Iterative Code Feedback.** For each new problem, the model generates a Python solution and test cases (prompts in Appendix C.3, C.4). Candidate solutions are executed in an isolated environment, and the results are fed back to the model, which iteratively refines its solution and tests until all

Table 1: Dataset statistics and quality measures. Gen. Size: Initial generation size; Filtered: Problems removed via filtering; Avg. Tests: Avg. # test cases per problem; Human Correct: Human-verified correctness (%); Coverage: Test coverage (%); Hours: Person-hours spent generating.

Dataset	Gen. Size	# Filtered	Avg. Tests	% Human Correct	% Coverage	Hours
MBPP-New	1002	539	8.30	97%	99%	13
MBPP-Guided	992	572	8.86	98%	99%	14
MBPP-Hard	1007	223	10.35	96%	100%	14
MBPP-Hard-Guided	994	471	8.86	96%	100%	15
Leetcode-New	997	170	8.22	98%	99%	25
Leetcode-Guided	991	179	8.66	97%	100%	27

tests pass or a maximum number of iterations is reached. A key improvement of InfoSynth over prior methods is feeding the entire feedback history at each step, giving the model richer context; Section 5.3 explains how this induces chain-of-thought reasoning (Wei et al., 2022). Importantly, problems failing self-verification are excluded from the final dataset but still serve as seeds for the next generation round to encourage diversity.

**Deduplication.** We use the MinHash + LSH algorithm with 250 permutations and a 0.75 similarity threshold to remove textually similar problems, similar to that done by Majumdar et al. (2024).

**Postprocessing.** Generated problem descriptions are not always well-aligned with their test cases. For example, a problem may not describe how to handle edge-cases such as null inputs or empty arrays. In some problems, it is unreasonable to expect a test-taker to infer the desired behavior (ex., should we return None, or -1 on an empty array input?). An example of such a problem is given in Appendix D. For each problem-test pair, the model is prompted to rephrase the question to incorporate details on handling obscure edge-cases. The prompt is given in Appendix C.5.

# 4.2 EXPERIMENTAL SETUP

We generate six datasets using GPT-40 (Hurst et al., 2024) as the generator. The first dataset, called MBPP-Guided, is seeded with MBPP. The second dataset, MBPP-Hard-Guided, is also seeded with MBPP, but during mutation, the model is specifically prompted to make the questions more difficult. The third dataset, Leetcode-New, is seeded with a Leetcode-Guided dataset developed by Xia et al. (2025). For each of these three datasets, we perform the generation process again, this time without using K-farthest neighbor selection, resulting in a total of six datasets. These additional datasets are referred to as MBPP-New, MBPP-Hard, and Leetcode-New, respectively.

# 5 RESULTS AND ANALYSIS

We categorize generated problems as: (1) **Passing** (solution passes all tests), (2) **Failing** (fails  $\geq$ 1 test), (3) **Erroring** (syntax/runtime error), and (4) **Unparsable** (malformed, e.g., missing [solution]/[test] tags, more common for smaller models). Table 1 reports benchmark statistics: test cases equal the number of assert statements, and test coverage is the fraction of code lines executed when all tests are run. We also evaluate SOTA models on all datasets (Table 2); Qwen2.5 models use 4-bit quantization. Since MBPP contains vague/misleading problems (Austin et al., 2021), we post-processed it for fairer comparison. Thus, Table 2 focuses on post-processed results as these provide the most fair comparison, omitting filtered versions (see Section 5.4). For each benchmark, we randomly sample 100 problems and manually verify that solutions and test cases are correct. A solution is correct only if its syntax and logic solve the described problem, and a test is correct only if its assertions match the task.

#### 5.1 NOVELTY AND DIVERSITY ANALYSIS

Figure 4a presents the novelty and diversity of MBPP-New and MBPP-Hard relative to MBPP-Original, while Figure 4b shows the same comparison for Leetcode-New relative to Leetcode-Original. Overall, our pipeline produces datasets that are more novel and diverse than the original seeds. However, filtering and post-processing reduce novelty compared to the initial generation. We

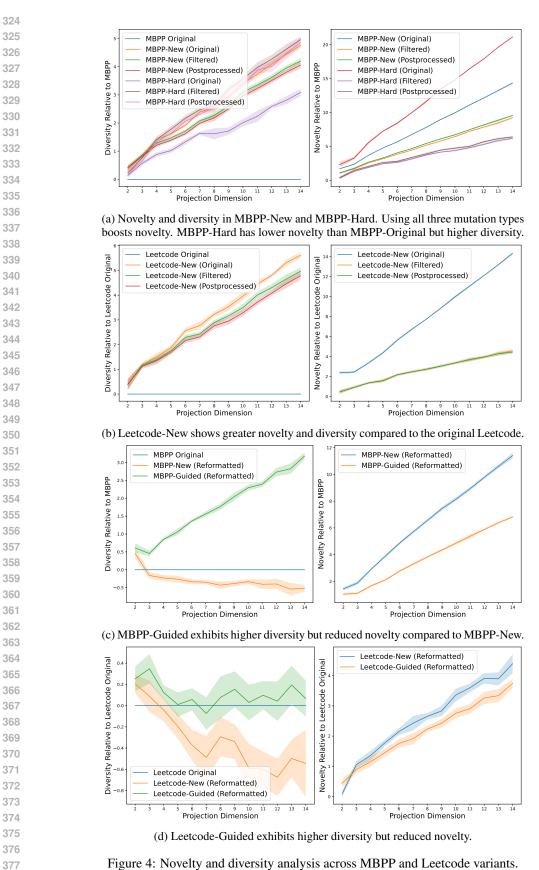


Figure 4: Novelty and diversity analysis across MBPP and Leetcode variants.

attribute this to LLM memorization (Huang et al., 2025; Kiyomaru et al., 2024), since more novel problems are often out-of-distribution and harder for the model to solve. All UMAP simulations use 80 neighbors and a minimum distance of 0.1, except for Figure 4d, which uses 30 neighbors due to smaller dataset size. In general, our results are not sensitive to UMAP hyperparameters.

Figures 4c and 4d show that k-farthest-neighbor filtering improves dataset novelty and diversity. This comes at the cost of generating easier problems (Table 2), highlighting another advantage of InfoSynth in controlling the novelty–diversity–difficulty tradeoff. Empirically, this arises because the generator struggles to produce difficult problems unless they are conceptually aligned with the seeds. We also observe a tradeoff between novelty and diversity: highly diverse datasets tend to concentrate around low-density regions in the seed-embedding distribution, which increases novelty but reduces diversity. Our results also show that filtering and post-processing reliably improve diversity.

#### 5.2 EFFECT OF VARYING MUTATION DIFFICULTIES

Table 2 shows MBPP-Hard scores are 8%-15% lower than MBPP-Original across most models, suggesting hard mutations effectively raise difficulty. This comes with a tradeoff of the dataset having reduced diversity and novelty as the problems tend to be concentrated around fewer, but more challenging topics. Hence, the set of mutation difficulties can be carefully chosen in InfoSynth in order to control the difficulty of the produced benchmark.

#### 5.3 EFFECT OF ITERATIVE CODE FEEDBACK

We find that passing solution-test pairs increase by 20% over 5 feedback iterations, showing the effectiveness of code iteration in producing robust problems. Error rates drop as the LLM fixes syntax/runtime issues, though the unparsable rate rises slightly due to occasional formatting failures. Appendix B shows feedback curves. Three iterations are typically ideal; further iterations yield marginal gains not worth the extra inference cost. We also find that iterative feedback acts as chain-of-thought (CoT) reasoning (Wei et al., 2022), as the model leverages the full feedback history to refine solutions/tests, lowering both error and failure rates. An example of this is in Appendix F.

#### 5.4 EFFECT OF THE POSTPROCESSING STEP

Appendix D shows two post-processed examples, where the model resolves ambiguous edge cases and sometimes rephrases statements more concisely without losing information. Table 2 shows 5–15% accuracy gains across most test-taker models, confirming that post-processing reduces ambiguity. Manual verification of 100 problems per dataset further shows 100% of post-processed problems are correctly reformatted without altering the core question.

#### 5.5 RELATING DIVERSITY AND TOPIC COVERAGE

For each problem in the original MBPP, MBPP-New, and MBPP-Guided, we prompted GPT-40-mini (Hurst et al., 2024) to provide up to 3 topic labels describing the problem similar to the approach used by Zhao et al. (2025). The list of available topic labels was the same as the list of topics available in the Leetcode dataset (kaysss, 2025). The prompt is given in Appendix G.

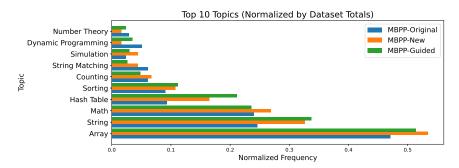


Figure 5: Fraction of problems relating to each topic for the 10 most common topics

432

Table 2: Test-taker performance on datasets

Filtered Postprocessed Model Dataset %Pass %Fail %Err %Pass %Fail %Err 45.99 MBPP-Original 52.67 1.34 60.43 37.701.87 51.58 MBPP-New 34.32 4.08 61.60 45.83 2.60 Qwen2.5-7b-Instruct MBPP-Guided 43.88 54.55 1.57 MBPP-Hard 19.82 72.52 7.66 31.08 6.76 62.16 MBPP-Hard-Guided 39.49 56.48 4.03 Leetcode-Original 11.84 83.77 4.39 25.29 0.59 Leetcode-New 74.12 25.88 72.35 1.76 29.05 70.39 Leetcode-Guided 0.56 MBPP-Original 46.79 48.93 4.28 52.41 42.51 5.08 MBPP-New 31.35 55.47 38.78 49.54 11.69 13.17 MBPP-Guided 40.56 49.13 Qwen2.5-3b-Coder 10.31 MBPP-Hard 54.95 21.62 23.42 24.77 53.15 22.07 MBPP-Hard-Guided 29.30 57.75 12.95 Leetcode-Original 2.63 85.53 11.84 Leetcode-New 18.24 62.94 18.82 25.53 64.12 12.35 Leetcode-Guided 15.08 75.42 9.50 58.02 66.04 30.48 0.00 MBPP-Original 36.10 5.88 MBPP-New 56.96 40.63 2.41 67.35 29.68 2.97 MBPP-Guided 5.59 GPT-4.1-Mini 72.0322.38 MBPP-Hard 44.14 50.45 5.41 55.86 38.29 5.86 MBPP-Hard-Guided 68.58 27.60 3.82 Leetcode-Original 32.89 54.82 12.28 45.28 38.24 40.59 16.47 50.00 9.41 Leetcode-New Leetcode-Guided 48.60 46.37 5.03 64.97 35.03 MBPP-Original 0.00 68.72 31.02 0.26 MBPP-New 53.99 45.64 0.3763.64 36.18 0.19Gemini-2.0-Flash MBPP-Guided 26.92 72.031.05 MBPP-Hard 44.59 52.25 3.15 50.00 45.05 4.95 MBPP-Hard-Guided 62.63 34.82 2.55 Leetcode-Original 32.46 64.47 3.07 Leetcode-New 44.71 54.12 1.18 51.18 46.47 2.35 50.28 3.35 Leetcode-Guided 46.37 MBPP-Original 63.37 0.00 70.86 29.14 0.00 36.63 MBPP-New 44.71 55.29 0.00 64.7535.25 0.00MBPP-Guided Claude 3.7 Sonnet 74.83 24.83 0.35 MBPP-Hard 45.05 52.25 2.70 57.21 40.09 2.70 MBPP-Hard-Guided 66.24 33.12 0.64 Leetcode-Original 31.14 67.11 1.75 Leetcode-New 44.71 54.71 0.59 44.71 55.29 0.00 Leetcode-Guided 48.60 51.40 0.00 MBPP-Original 66.58 33.42 0.00 70.05 29.68 0.27 29.68 MBPP-New 58.26 41.19 0.56 70.130.19o4-mini MBPP-Guided 77.27 22.55 0.17 MBPP-Hard 47.75 49.55 2.70 32.88 63.06 4.05 MBPP-Hard-Guided 72.6126.11 1.27 Leetcode-Original 38.60 58.77 2.63 Leetcode-New 40.59 58.82 0.59 50.59 47.06 2.35 Leetcode-Guided 46.93 51.96 1.18

471 472 473

474

475

476

477

468

469

470

Figure 5 shows that InfoSynth increases the number of problems that use each concept for most topics, creating more diverse and widely covering problems. We find that for some topics with lesser coverage in the original MBPP dataset, our pipeline produces many more problems covering those topics. Evidently, MBPP-Guided has more topic coverage than MBPP-New across most categories, demonstrating the effectiveness of k-farthest-neighbor-filtering in encouraging diversity.

478 479 480

#### 6 Conclusion

481 482

483

484

485

In this paper, we introduced InfoSynth, a novel framework to efficiently calculate the diversity and novelty of new benchmark in a cheaper, more-efficient manner. Moreover, we demonstrate the effectiveness of InfoSynth in generating high-quality, novel, and diverse synthetic coding dataset from seed data. We hope that future work will leverage our designs to create robust, novel, and diverse benchmarks.

# 7 ETHICS STATEMENT

This work does not involve human subjects, personally identifiable data, or sensitive attributes. All datasets used (MBPP, LeetCode, and Codeforces) are publicly available, and our generated benchmarks were produced through synthetic problem generation and automated verification. We have carefully ensured that no private or proprietary code was included. Potential risks include the misuse of generated benchmarks for unfair evaluation or dataset contamination in future model training; to mitigate this, we document our pipeline in detail and encourage responsible use. We adhered to the ICLR Code of Ethics throughout the research and submission process.

# 8 REPRODUCIBILITY STATEMENT

We provide detailed descriptions of our dataset generation pipeline (Section 4), evaluation setup (Section 5), and algorithmic parameters (Appendix A). All prompts used for mutation, crossover, and verification are listed in Appendix C, and examples of generated problems are in Appendix E. Post-processing steps and deduplication methods are described in Section 4 and Appendix D. To facilitate reproducibility, we plan to release the full codebase and generated datasets upon publication.

# REFERENCES

- Jacob Austin, Augustus Odena, Maxwell I. Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie J. Cai, Michael Terry, Quoc V. Le, and Charles Sutton. Program Synthesis with Large Language Models. *CoRR*, abs/2108.07732, 2021. URL https://arxiv.org/abs/2108.07732.
- Bei Chen, Fengji Zhang, Anh Nguyen, Daoguang Zan, Zeqi Lin, Jian-Guang Lou, and Weizhu Chen. Code T: Code Generation with Generated Tests. *CoRR*, abs/2207.10397, 2022. doi: 10. 48550/ARXIV.2207.10397. URL https://doi.org/10.48550/arXiv.2207.10397.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating Large Language Models Trained on Code. *CoRR*, abs/2107.03374, 2021. URL https://arxiv.org/abs/2107.03374.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training Verifiers to Solve Math Word Problems. *CoRR*, abs/2110.14168, 2021. URL https://arxiv.org/abs/2110.14168.
- Chunyuan Deng, Yilun Zhao, Xiangru Tang, Mark Gerstein, and Arman Cohan. Investigating Data Contamination in Modern Benchmarks for Large Language Models. *CoRR*, abs/2311.09783, 2023a. doi: 10.48550/ARXIV.2311.09783. URL https://doi.org/10.48550/arXiv.2311.09783.
- Chunyuan Deng, Yilun Zhao, Xiangru Tang, Mark Gerstein, and Arman Cohan. Investigating Data Contamination in Modern Benchmarks for Large Language Models. *CoRR*, abs/2311.09783, 2023b. doi: 10.48550/ARXIV.2311.09783. URL https://doi.org/10.48550/arXiv.2311.09783.
- Yuyang Ding, Xinyu Shi, Xiaobo Liang, Juntao Li, Qiaoming Zhu, and Min Zhang. Unleashing Reasoning Capability of LLMs via Scalable Question Synthesis from Scratch. *CoRR*, abs/2410.18693, 2024. doi: 10.48550/ARXIV.2410.18693. URL https://doi.org/10.48550/arXiv.2410.18693.
- Shahriar Golchin and Mihai Surdeanu. Data Contamination Quiz: A Tool to Detect and Estimate Contamination in Large Language Models. *CoRR*, abs/2311.06233, 2023. doi: 10.48550/ARXIV. 2311.06233. URL https://doi.org/10.48550/arXiv.2311.06233.
- Kaixuan Huang, Jiacheng Guo, Zihao Li, Xiang Ji, Jiawei Ge, Wenzhe Li, Yingqing Guo, Tianle Cai, Hui Yuan, Runzhe Wang, Yue Wu, Ming Yin, Shange Tang, Yangsibo Huang, Chi Jin, Xinyun

- Chen, Chiyuan Zhang, and Mengdi Wang. MATH-Perturb: Benchmarking LLMs' Math Reasoning Abilities against Hard Perturbations. *CoRR*, abs/2502.06453, 2025. doi: 10.48550/ARXIV. 2502.06453. URL https://doi.org/10.48550/arXiv.2502.06453.
  - Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card, 2024. URL https://arxiv.org/abs/2410.21276.
  - Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. LiveCodeBench: Holistic and Contamination Free Evaluation of Large Language Models for Code. *CoRR*, abs/2403.07974, 2024. doi: 10.48550/ARXIV.2403.07974. URL https://doi.org/10.48550/arXiv.2403.07974.
  - kaysss. leetcode-problem-set. https://huggingface.co/datasets/kaysss/ leetcode-problem-set, 2025. Accessed: 2025-05-11.
  - Diederik P. Kingma and Max Welling. Auto-encoding Variational Bayes. In Yoshua Bengio and Yann LeCun (eds.), 2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings, 2014. URL http://arxiv.org/abs/1312.6114.
  - Hirokazu Kiyomaru, Issa Sugiura, Daisuke Kawahara, and Sadao Kurohashi. A Comprehensive Analysis of Memorization in Large Language Models. In Saad Mahamood, Minh Le Nguyen, and Daphne Ippolito (eds.), *Proceedings of the 17th International Natural Language Generation Conference, INLG 2024, Tokyo, Japan, September 23 27, 2024*, pp. 584–596. Association for Computational Linguistics, 2024. doi: 10.18653/V1/2024.INLG-MAIN.45. URL https://doi.org/10.18653/v1/2024.inlg-main.45.
  - Alexander Kraskov, Harald Stögbauer, and Peter Grassberger. Estimating mutual information. *Physical Review E*, 69(6), 2004. ISSN 1550-2376. doi: 10.1103/physreve.69.066138. URL http://dx.doi.org/10.1103/PhysRevE.69.066138.
  - Xiang Lisa Li, Farzaan Kaiyom, Evan Zheran Liu, Yifan Mai, Percy Liang, and Tatsunori Hashimoto. AutoBencher: Towards Declarative Benchmark Construction. In *The Thirteenth International Conference on Learning Representations, ICLR* 2025, Singapore, April 24-28, 2025. OpenReview.net, 2025a. URL https://openreview.net/forum?id=ymt4crbbXh.
  - Yisen Li, Lingfeng Yang, Wenxuan Shen, Pan Zhou, Yao Wan, Weiwei Lin, and Dongping Chen. CrowdSelect: Synthetic Instruction Data Selection with Multi-LLM Wisdom. *CoRR*, abs/2503.01836, 2025b. doi: 10.48550/ARXIV.2503.01836. URL https://doi.org/10.48550/arXiv.2503.01836.
  - Jiate Liu, Yiqin Zhu, Kaiwen Xiao, Qiang Fu, Xiao Han, Wei Yang, and Deheng Ye. RLTF: Reinforcement Learning from Unit Test Feedback. *CoRR*, abs/2307.04349, 2023. doi: 10.48550/ARXIV.2307.04349. URL https://doi.org/10.48550/arXiv.2307.04349.
  - Jiawei Liu, Thanh Nguyen, Mingyue Shang, Hantian Ding, Xiaopeng Li, Yu Yu, Varun Kumar, and Zijian Wang. Learning Code Preference via Synthetic Evolution. *CoRR*, abs/2410.03837, 2024a. doi: 10.48550/ARXIV.2410.03837. URL https://doi.org/10.48550/arXiv.2410.03837.
  - Zhihan Liu, Shenao Zhang, and Zhaoran Wang. DSTC: Direct Preference Learning with Only Self-generated Tests and Code to Improve Code LMs. *CoRR*, abs/2411.13611, 2024b. doi: 10.48550/ARXIV.2411.13611. URL https://doi.org/10.48550/arXiv.2411.13611.
- Somshubra Majumdar, Vahid Noroozi, Sean Narenthiran, Aleksander Ficek, Jagadeesh Balam, and Boris Ginsburg. Genetic Instruct: Scaling up Synthetic Generation of Coding Instructions for Large Language Models. *CoRR*, abs/2407.21077, 2024. doi: 10.48550/ARXIV.2407.21077. URL https://doi.org/10.48550/arXiv.2407.21077.
  - Leland McInnes, John Healy, Nathaniel Saul, and Lukas Großberger. UMAP: Uniform Manifold Approximation and Projection. *J. Open Source Softw.*, 3(29):861, 2018. doi: 10.21105/JOSS. 00861. URL https://doi.org/10.21105/joss.00861.

- open-r1. Codeforces Problems Dataset. https://huggingface.co/datasets/ open-r1/codeforces, 2024. Accessed: 2025-05-11.
  - John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal Policy Optimization Algorithms. *CoRR*, abs/1707.06347, 2017. URL http://arxiv.org/abs/1707.06347.
  - Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. Mpnet: Masked and permuted pretraining for language understanding. In *Advances in Neural Information Processing Systems*, volume 33, 2020.
  - Florian Tambon, Amin Nikanjam, Cyrine Zid, Foutse Khomh, and Giuliano Antoniol. TaskEval: Assessing Difficulty of Code Generation Tasks for Large Language Models. *ArXiv preprint*, abs/2407.21227, 2024. URL https://arxiv.org/abs/2407.21227.
  - Sang T. Truong, Yuheng Tu, Percy Liang, Bo Li, and Sanmi Koyejo. Reliable and Efficient Amortized Model-based Evaluation. *CoRR*, abs/2503.13335, 2025. doi: 10.48550/ARXIV.2503.13335. URL https://doi.org/10.48550/arXiv.2503.13335.
  - Qing Wang, Sanjeev R. Kulkarni, and Sergio Verdú. Divergence estimation for multidimensional densities via k-nearest-neighbor distances. *IEEE Trans. Inf. Theory*, 55(5):2392–2405, 2009. doi: 10.1109/TIT.2009.2016060. URL https://doi.org/10.1109/TIT.2009.2016060.
  - Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. Self-Instruct: Aligning Language Models with Self-generated Instructions. In Anna Rogers, Jordan L. Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, *ACL 2023, Toronto, Canada, July 9-14, 2023*, pp. 13484–13508. Association for Computational Linguistics, 2023. doi: 10.18653/V1/2023.ACL-LONG.754. URL https://doi.org/10.18653/v1/2023.acl-long.754.
  - Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed H. Chi, Quoc Le, and Denny Zhou. Chain of Thought Prompting Elicits Reasoning in Large Language Models. *CoRR*, abs/2201.11903, 2022. URL https://arxiv.org/abs/2201.11903.
  - Yunhui Xia, Wei Shen, Yan Wang, Jason Klein Liu, Huifeng Sun, Siyue Wu, Jian Hu, and Xiaolong Xu. LeetCodeDataset: A Temporal Dataset for Robust Evaluation and Efficient Training of Code LLMs. *CoRR*, abs/2504.14655, 2025. doi: 10.48550/ARXIV.2504.14655. URL https://doi.org/10.48550/arXiv.2504.14655.
  - Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, Qingwei Lin, and Daxin Jiang. WizardLM: Empowering Large Pre-trained Language Models to Follow Complex Instructions. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024.* OpenReview.net, 2024. URL https://openreview.net/forum?id=CfXh93NDgH.
  - Zhangchen Xu, Yang Liu, Yueqin Yin, Mingyuan Zhou, and Radha Poovendran. KodCode: A Diverse, Challenging, and Verifiable Synthetic Dataset for Coding. *CoRR*, abs/2503.02951, 2025. doi: 10.48550/ARXIV.2503.02951. URL https://doi.org/10.48550/arXiv.2503.02951.
  - Huaye Zeng, Dongfu Jiang, Haozhe Wang, Ping Nie, Xiaotong Chen, and Wenhu Chen. ACE-CODER: Acing Coder RL via Automated Test-case Synthesis. *CoRR*, abs/2502.01718, 2025. doi: 10.48550/ARXIV.2502.01718. URL https://doi.org/10.48550/arXiv.2502.01718.
  - Hugh Zhang, Jeff Da, Dean Lee, Vaughn Robinson, Catherine Wu, Will Song, Tiffany Zhao, Pranav Raja, Dylan Slack, Qin Lyu, Sean Hendryx, Russell Kaplan, Michele Lunati, and Summer Yue. A Careful Examination of Large Language Model Performance on Grade School Arithmetic. *CoRR*, abs/2405.00332, 2024. doi: 10.48550/ARXIV.2405.00332. URL https://doi.org/10.48550/arXiv.2405.00332.

Xueliang Zhao, Wei Wu, Jian Guan, and Lingpeng Kong. PromptCoT: Synthesizing Olympiad-level Problems for Mathematical Reasoning in Large Language Models. *CoRR*, abs/2503.02324, 2025. doi: 10.48550/ARXIV.2503.02324. URL https://doi.org/10.48550/arXiv.2503.02324.

Terry Yue Zhuo, Minh Chien Vu, Jenny Chim, Han Hu, Wenhao Yu, Ratnadira Widyasari, Imam Nur Bani Yusuf, Haolan Zhan, Junda He, Indraneil Paul, Simon Brunner, Chen Gong, Thong Hoang, Armel Randy Zebaze, Xiaoheng Hong, Wen-Ding Li, Jean Kaddour, Ming Xu, Zhihan Zhang, Prateek Yadav, Naman Jain, Alex Gu, Zhoujun Cheng, Jiawei Liu, Qian Liu, Zijian Wang, David Lo, Binyuan Hui, Niklas Muennighoff, Daniel Fried, Xiaoning Du, Harm de Vries, and Leandro von Werra. BigCodeBench: Benchmarking Code Generation with Diverse Function Calls and Complex Instructions. *CoRR*, abs/2406.15877, 2024. doi: 10.48550/ARXIV.2406. 15877. URL https://doi.org/10.48550/arxiv.2406.15877.

703 704

705

706

708

709

710

# A ALGORITHM FOR PROBLEM GENERATION

For MBPP-New we used N=1000,  $N_c=10$ ,  $B_s=30$ , C=2,  $B_c=5$ ,  $N_{it}=5$ . For MBPP-Hard we used N=500,  $N_c=10$ ,  $B_s=15$ , C=1,  $B_c=4$ ,  $N_{it}=5$ . For Leetcode-New we used N=1000,  $N_c=10$ ,  $B_s=30$ , C=2,  $B_c=4$ ,  $N_{it}=5$ .

For MBPP-Guided we used  $N=1000,\ N_c=10,\ B_s=30,\ C=3,\ B_c=5,\ N_{it}=3.$  For MBPP-Hard-Guided we used  $N=500,\ N_c=10,\ B_s=15,\ C=3,\ B_c=4,\ N_{it}=3.$  For Leetcode-Guided we used  $N=1000,\ N_c=10,\ B_s=30,\ C=3,\ B_c=4,\ N_{it}=3.$ 

```
711
         Algorithm 1: Problem Generation with Evolutionary Strategies
712
713
         Input: N: total number of problems to generate
714
         N_c: number of colonies
         B_s: number of problems to sample as the seed data for each colony
715
716
         C: number of problems to generate per crossover operation
         B_c: crossover seed batch size; this is the number of problems fed into the model for it to
717
         combine
718
         N_{\rm it}: number of code feedback iterations
719
720
         Generate (seedData, numSamples):
            Initialize problems \leftarrow \emptyset
721
            for colony \leftarrow 1 to N_c do
722
                N_s \leftarrow N/N_c
723
                colonySeedData \leftarrow Random sample of size B_s from seedData
724
                problems \leftarrow problems \cup EvolveColony (colonySeedData, N_s)
725
                deduplicate(problems)
726
            end
727
            return problems
728
         EvolveColony (seedData, N_s):
729
            newProblems \leftarrow \emptyset
730
            while |newProblems| < N_s do
731
                operation \leftarrow "mutation" with p = 0.5, "crossover" with p = 0.5
732
                if operation == "mutation" then
733
                    problem ← random sample of size 1 from seedData
734
                    problems \leftarrow problems \cup mutate(problem)
735
                else if operation == "crossover" then
                    batch \leftarrow random sample of size B_c from seedData
737
                    problems \leftarrow problems \cup crossover(C, batch)
738
                end
                if k-farthest-neighbor-filtering-enabled then
739
                    U \leftarrow \text{newProblems} \cup \text{seedData}
740
                    problems \leftarrow select K problems with least cosine similarity to U
741
742
                newProblems \leftarrow newProblems \cup problems
743
                deduplicate(newProblems)
744
                seedData \leftarrow seedData \cup problems
745
                deduplicate(seedData)
746
747
            return newProblems
748
         GenerateSolutionsTests(problem):
749
            GenerateTests(problem)
750
            GenerateSolutions(problem)
751
            for i = 1 to N_{it} do
752
                Run tests against solutions
753
                Feed output back into LLM to modify tests and solution
754
            end
```

# B CODE-FEEDBACK RESULTS

 Figures 6, 7 show how the proportion of problems that pass, fail, error, and are unparsable changes as a function of the number of code feedback iterations. The results shown are consistent across all datasets; in general, more than 3 feedback iterations provides minimal gains in pass rate.

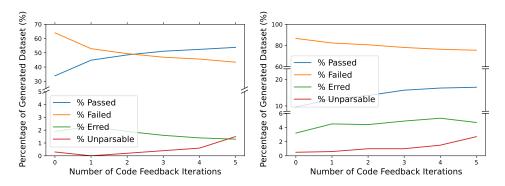


Figure 6: Self-Verification for MBPP-New Figure 7: Self-Verification for Leetcode-New

# C PROBLEM GENERATION PIPELINE PROMPTS

We note that our prompts share some similarity with those used by Majumdar et al. (2024).

#### C.1 MUTATION PROMPTS

```
782
783
      % easv
784
      Please decrease the difficulty of the given programming test
785
      question a bit.
      The new problem should be conceptually similar to the given
786
      question,
787
      but should not simply paraphrase it. Do not provide any hints,
788
      solutions
789
      or outputs. Only one new instruction is allowed.
790
      Original Question: {instruction}
791
      New Question:
792
793
      % medium
794
      Please create a new programming problem of the same difficulty as
795
796
      given programming test question. The new problem should be
      conceptually
797
      similar to the given question, but should not simply paraphrase
798
      it.
799
      Do not provide any hints, solutions or outputs. Only one new
800
      instruction is allowed.
801
      Original Question: {instruction}
802
      New Question:
804
      % hard
805
      Please increase the difficulty of the given programming test
806
      question a bit.
807
      Do not provide any hints, solutions or outputs. Only one new
      instruction is allowed.
808
      Original Question: {instruction}
809
      New Question:
```

```
810
      C.2 Crossover Prompt
811
      I will provide you with a set of coding questions. Please give me
813
      a new coding
      question that combines core concepts from two or more of the given
814
      questions.
815
      Please ensure that the new question is novel and does not simply
816
      paraphrase
817
      any of the problems I am giving you. Do not include any extra
818
      information
819
      that would help a test-taker other than the problem statement
820
      itself.
821
822
      Question 1:
823
      {instruction 1}
824
825
      Question 2:
      {instruction 2}
826
827
      . . .
828
      New Question:
829
830
      C.3 SOLUTION & TEST GENERATION PROMPT
831
832
      We note that our prompts are similar to those used by Xu et al. (2025).
833
834
      You are an expert in Python coding.
835
      ## Task:
836
      Please answer the question and generate unit tests to verify your
837
      answer.
838
839
      ## Output Format:
      Your solution and unit tests should be presented in the format
840
      within the
841
      specified sections below. Ensure your code is within code blocks.
842
      For the
843
      tests, use pytest style by defining individual test functions
844
      (without
845
      classes) and using assert statements. Your tests should be
      implementation
847
      independent. Ensure that you include the <|Solution Begin|>,
848
      <|Solution End|>, <|Test Begin|>, and <|Test End|> tags as
849
      depicted. The
      solution function must be named solution.
850
851
      <|Solution Begin|>
852
      {Solution Code in Python}
853
      <|Solution End|>
854
855
      <|Test Begin|>
856
      {Unit Test Code in Python}
857
      <|Test End|>
858
859
      ## Example
860
      Below is an example output format implementing a simple a + b
861
      function.
      <|Solution Begin|>
862
      def add(a, b):
863
          "\"Returns the sum of a and b.\"\"\"
```

```
864
          return a + b
865
      <|Solution End|>
867
      <|Test Begin|>
868
      from solution import add
      def test_add_positive_numbers():
869
          assert add(2, 3) == 5
870
      def test_add_with_zero():
871
          assert add(0, 5) == 5
872
          assert add(5, 0) == 5
873
      def test_add_negative_numbers():
874
          assert add(-1, -1) == -2
875
      def test_add_mixed_sign_numbers():
876
          assert add(-1, 3) == 2
877
      <|Test End|>
878
879
      ## Question:
880
      {problem}
881
882
      C.4 SOLUTION & TEST GENERATION WITH ITERATIVE FEEDBACK PROMPT
883
      You are an expert in Python coding.
884
      ## Task:
885
      Please answer the question and generate unit tests to verify your
886
      answer. The
887
      entire chat history of your previous attempts to generate
888
      questions and unit
889
      tests is presented below in the "Chat History" section, along with
890
      the output
891
      of running your solution against your tests in a code execution
892
      environment.
893
      Please modify only your tests and/or solution to be more correct.
894
      ## Output Format:
895
      <Same as "Output Format" section above>
896
897
      ## Chat History:
898
      Attempt 1 Solution:
899
      {attempt 1 solution}
900
901
      Attempt 1 Code Execution Output:
902
      {attempt 1 code output}
903
904
      Attempt 2 Solution:
      {Attempt 2 solution}
905
906
      Attempt 2 Code Execution Output:
907
      {attempt 2 code output}
908
909
910
      ## Question:
911
      {problem}
912
913
      C.5 POSTPROCESSING PROMPTS
914
915
      You are an expert in Python coding. Here is a coding problem
      with associated test cases. Please rephrase the question so it
916
      describes what the user should output for edge-cases without
917
      changing the essence of the problem. Add as little information
```

as possible, only describing what the user should output for edge-cases that cannot be inferred from the problem description. Do not include anything except for the rewritten problem in your response and do not include the test cases.

Question:
{question}
Tests:
{tests}

# D EXAMPLE OF POST-PROCESSED PROBLEMS

#### **MBPP-New Dataset Original Problem:**

Write a function that filters a list of usernames stored in a dictionary, returning only those associated with students who fall within a specified age range.

# **Post-Processed Version:**

Write a function that filters a list of usernames stored in a dictionary, returning only those associated with students who fall within a specified age range. Ensure that the function returns an empty list when there are no students or when the input dictionary is empty.

# **Leetcode-New Dataset Original Problem:**

Alice and Bob are engaged in a strategic game on an infinite 2D plane with n points provided by their coordinates in two integer arrays, xCoord and yCoord. They take turns, starting with Alice, attacking a point on the plane to capture it, with the condition that once a point is attacked, it is removed permanently from the game, and they have to remove exactly 1 point per turn.

The winner is the player who either removes the point that leaves no rectangle capable of being formed using the remaining points on their turn or can force the scenario by optimal play such that the opponent has no such move left on their subsequent turns.

Given the arrays xCoord and yCoord, along with knowledge of optimal strategies for both players, determine if Alice, who starts the game, can always guarantee a win. Return true if Alice has a winning strategy, or false if Bob can always force a win even with Alice starting first.

# **Post-Processed Version:**

Alice and Bob are playing a strategic game on an infinite 2D plane with n points defined by their coordinates in two integer arrays, xCoord and yCoord. They alternately take turns, with Alice starting first, to attack and permanently remove exactly one point at a time. The objective is for a player to leave no possibility of forming a rectangle using any four of the remaining points. A player wins if they achieve this or if they can force a scenario where the opponent has no such moves left. Given the arrays xCoord and yCoord, determine if Alice has a guaranteed winning strategy. Return true if Alice can always win, and false if Bob can always force a win or if there are no points to start with.

# E EXAMPLES OF GENERATED PROBLEMS

We note that for conciseness, the examples shown in this section are before the postprocessing step. In this section, we illuminate the different mechanisms through which mutation and crossover generate interesting problems.

#### E.1 MUTATION OPERATION: ADDING OR REMOVING CONSTRAINTS

This example shows how mutation creates three variants of the question by adding requirements or removing constraints from the problem. We can see that all three problems require similar conceptual understanding, but the harder ones simply require more code or bookkeeping.

Original Question:

Write a function to filter the height and width of students which are stored in a dictionary.

#### Easy Mutation:

Write a function that filters a list of usernames stored in a dictionary, returning only those associated with students who fall within a specified age range. Ensure that the function returns an empty list when there are no students or when the input dictionary is empty.

#### Medium Mutation:

Create a function that filters user profiles based on a dictionary. The function should return a list of user IDs for profiles where the age is within a specified range (inclusive) and the profile only contains lowercase alphabetic characters. If there are no valid profiles that match the criteria, the function should return an empty list.

# Hard Mutation:

Create a function to filter user passwords from a dictionary, returning only those that are valid for students whose dimensions (height, width, and weight) are within a given range. Each valid password must include at least one uppercase letter, one lowercase letter, one digit, and one special character. Ensure the solution appropriately handles and returns an empty dictionary for cases where no users are provided or all entries are invalid due to dimension or password criteria.

# E.2 MUTATION OPERATION EXAMPLE: CREATIVE MODIFICATION

This example shows how mutation creates three variants of the question by creatively modifying the central idea of the problem itself. This differs from the previous example; in this case, the harder questions require a fundamental understanding of new topics.

#### Original Question:

Write a function to find the perimeter of a rectangle.

#### Easy Mutation:

1022 Write a function to find the area of a rectangle.

# 1024 Medium Mutation:

1025 Write a function to calculate the area of a trapezoid given its base lengths and height.

1026 Hard Mutation:

Write a function to find the area of a rhombus given its diagonals and verify if the rhombus is also a square by using its side lengths.

# E.3 CROSSOVER OPERATION: COMBINING CONCEPTS

This example shows how crossover creates an interesting, novel question by combining two unrelated concepts. A key difference between crossover and mutation is that crossover does not introduce any new concepts or content into the generated problem as it draws from existing ones. This demonstrates the necessity of combining mutation and crossover into one pipeline; mutation introduces new concepts into the dataset while crossover takes existing concepts and uses them to create richer problems.

#### **Seed Problems:**

Ouestion 1:

There is a 50 x 50 chessboard with one knight and some pawns on it. You are given two integers kx and ky where (kx, ky) denotes the position of the knight, and a 2D array positions where positions[i] = [xi, yi] denotes the position of the pawns on the chessboard.

Alice and Bob play a turn-based game, where Alice goes first. In each player's turn:

The player selects a pawn that still exists on the board and captures it with the knight in the fewest possible moves. Note that the player can select any pawn, it might not be one that can be captured in the least number of moves.

In the process of capturing the selected pawn, the knight may pass other pawns without capturing them. Only the selected pawn can be captured in this turn.

Alice is trying to maximize the sum of the number of moves made by both players until there are no more pawns on the board, whereas Bob tries to minimize them.

Return the maximum total number of moves made during the game that Alice can achieve, assuming both players play optimally.

Note that in one move, a chess knight has eight possible positions it can move to, as illustrated below. Each move is two cells in a cardinal direction, then one cell in an orthogonal direction.

Question 2:

You are given an array points where points[i] = [xi, yi] represents the coordinates of a point on an infinite plane. Your task is to find the maximum area of a rectangle that:

1070 Ca

Can be formed using four of these points as its corners. Does not contain any other point inside or on its border. Has its edges parallel to the axes.

Return the maximum area that you can obtain or -1 if no such rectangle is possible.

## **Crossover Generated Question:**

Alice and Bob are playing a strategic game on an infinite 2D plane with n points defined by their coordinates in two integer arrays, xCoord and yCoord. They alternately take turns, with Alice

starting first, to attack and permanently remove exactly one point at a time. The objective is for a player to leave no possibility of forming a rectangle using any four of the remaining points. A player wins if they achieve this or if they can force a scenario where the opponent has no such moves left. Given the arrays xCoord and yCoord, determine if Alice has a guaranteed winning strategy. Return true if Alice can always win, and false if Bob can always force a win or if there are no points to start with.

# E.4 Crossover Operation: Multistep Reasoning

This example shows how crossover sometimes creates problems that require multistep reasoning. The new problem effectively requires a test-taker to solve both the seed problems.

#### **Seed Problems:**

1096 Question 1:

Write a python function to find the sum of fourth power of first  $\boldsymbol{n}$  odd natural numbers.

Question 2:

Write a python function to find the sum of xor of all pairs of numbers in the given array.

# **Crossover Generated Question:**

Write a python function to compute the combined metric of two properties for a given array of integers:

- 1. The sum of the fourth power of the first 'k' odd numbers from the array.
- 2. The sum of the bitwise XOR of all pairs of numbers formed by selecting one element from the first 'k' even numbers and one element from the last 'k' elements in the array.

The function should take two arguments — the array of integers and an integer 'k'. Ensure that 'k' is less than or equal to half the length of the array and that the array contains a sufficient number of odd and even numbers.

# F EXAMPLE OF CHAIN-OF-THOUGHT THINKING

This example shows how the model is able to reason about why its solution and tests are erroring due to receiving feedback from the code execution environment.

**Question:** Write a function that takes a list of strings and returns the second shortest string that contains all vowels, sorted alphabetically. If no such string exists, return an empty string.

```
Initial Solution:
```

```
def solution(strings):
    """

1128
    Returns the second shortest string containing all vowels,
1130    sorted alphabetically.
1131    """

1132    all_vowels = set('aeiou')

1133
    # Filter strings that contain all vowels
```

```
1134
          filtered_strings = [s for s in strings if all_vowels <=</pre>
1135
      set(s)]
1137
          # Sort them by length first and then alphabetically
1138
          filtered_strings.sort(key=lambda x: (len(x), x))
1139
          # Return the second shortest string or empty string if it
1140
      doesn't exist
1141
          return filtered_strings[1] if len(filtered_strings) > 1 else
1142
1143
      {Tests Omitted for Conciseness}
1144
      Solution after First Feedback iteration:
1145
1146
      Let's revise the test cases to ensure they correctly test the
1147
      function's behavior and check the solution implementation
1148
      and its handling of edge cases. First, we'll make sure the
1149
      solution properly handles strings with varying cases (uppercase
1150
      and lowercase) and that it verifies the full vowel set in a
1151
      case-insensitive manner. Additionally, let's adjust the tests to
1152
      cover more scenarios and better reflect the conditions described
1153
      in the question. Here's the revised version of the solution and
     unit tests:
1154
1155
      {Solutions Omitted for Conciseness}
1156
      {Tests Omitted for Conciseness}
1157
1158
      Solution after Second Feedback iteration:
1159
      We need to ensure that our solution and tests align correctly with
1160
      the task. Let's refine both the solution and the tests so that
1161
      they pass correctly according to the given problem statement. The
1162
      solution needs to ensure that we filter strings containing all
1163
     vowels in any case, sort them, and return the second shortest. The
1164
      test cases should reflect this logic.
1165
1166
      Here's the refined solution and tests:
1167
      {Solutions Omitted for Conciseness}
1168
      {Tests Omitted for Conciseness}
1169
1170
1171
      G TOPIC-LABELING PROMPT
1172
1173
      You are a precise tagger for coding interview problems.
      Given a problem statement and solution, choose up to THREE topics
1175
      from the bank below that best describe the core techniques/data
1176
      structures needed to solve it. Return JSON ONLY, no extra text,
1177
1178
      {"topics": ["Topic1", "Topic2"]}
1179
1180
1181
      Rules:
1182
     - Use EXACT spellings from the bank (case/spacing must match).
1183
      - Prefer the most specific tag available (e.g., "Binary Tree" over
1184
      "Tree", "Shortest Path" over "Graph" when appropriate).
1185
      - If the solution critically relies on a data structure (e.g.,
      "Heap (Priority Queue)"), include it.
1186
      - If multiple techniques are essential (e.g., DP + Bitmask),
1187
```

include both.

```
1188
      - Do NOT exceed 3 topics; order them by importance.
1189
      - If nothing fits, choose the closest general tag (e.g., "Graph",
      "Array", "Math") | never invent tags.
1191
1192
      Topic Bank (allowed values only):
1193
      Array; String; Hash Table; Dynamic Programming; Math; Sorting;
      Greedy; Depth-First Search; Binary Search; Database; Matrix; Tree;
1195
      Breadth-First Search; Bit Manipulation; Two Pointers; Prefix Sum;
1196
      Heap (Priority Queue); Simulation; Binary Tree; Graph; Stack;
1197
      Counting; Sliding Window; Design; Enumeration; Backtracking; Union
1198
      Find; Linked List; Number Theory; Ordered Set; Monotonic Stack;
1199
      Segment Tree; Trie; Combinatorics; Bitmask; Divide and Conquer;
1200
      Queue; Recursion; Geometry; Binary Indexed Tree; Memoization; Hash
1201
      Function; Binary Search Tree; Shortest Path; String Matching;
1202
      Topological Sort; Rolling Hash; Game Theory; Interactive;
      Data Stream; Monotonic Queue; Brainteaser; Doubly-Linked List;
1203
      Randomized; Merge Sort; Counting Sort; Iterator; Concurrency;
1204
      Probability and Statistics; Quickselect; Suffix Array; Line Sweep;
1205
      Minimum Spanning Tree; Bucket Sort; Shell; Reservoir Sampling;
1206
      Strongly Connected Component; Eulerian Circuit; Radix Sort;
1207
      Rejection Sampling; Biconnected Component
1208
1209
1210
      Input:
1211
      [Problem]
1212
      {problem}
1213
1214
      [Solution]
1215
      {solution}
1216
1217
1218
      Output:
1219
      JSON with key "topics" and UP TO 3 strings from the bank. No
1220
      prose, no explanations.
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241
```