

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 time-consuming. 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 information-theoretic 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 may not reliably solve the generated problems. This paper focuses on Python coding problems, whose solutions can be verified by executing them 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 have revealed that LLM training data is frequently contaminated by existing evaluation benchmarks, which can artificially inflate reported performance (Deng et al., 2024a;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

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

060 Our main contributions can be summarized as:

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

071 2 RELATED WORK

072
 073 **Synthetic Problem Generation.** Previous work has explored generating novel synthetic datasets
 074 from high-quality seed benchmarks. Majumdar et al. (2024); Wang et al. (2023); Xu et al. (2024);
 075 Zhao et al. (2025) show that LLMs can generate new instructions from existing ones. Chen et al.
 076 (2023); Liu et al. (2023; 2024b); Xu et al. (2025); Zeng et al. (2025) successfully used LLMs to
 077 generate unit tests for solution verification. Our end-to-end pipeline extends existing methods with
 078 code-execution environments to ensure robustness, novelty, and diversity.

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

087 3 DESIRABLE BENCHMARK PROPERTIES

088
 089 We propose a framework for characterizing the novelty and diversity of benchmarks. Our new
 090 novelty metric uses the KL-Divergence to capture how different the benchmark is from existing
 091 datasets, with the broader goal of creating benchmarks that are contamination-free. Previous work
 092 has used the KL-Divergence in a similar way; Schulman et al. (2017) use it to measure differences
 093 in reinforcement learning policies and Kingma & Welling (2014) use it to regularize output distribu-
 094 tions for Variational Autoencoders. Similarly, our proposed diversity metric uses Shannon entropy
 095 to capture how much variety exists among the problems; more diverse datasets provide a broader
 096 characterization of an LLM’s reasoning abilities.

097 3.1 AN INFORMATION-THEORY BASED FRAMEWORK FOR BENCHMARK ANALYSIS

098
 099 Formally, a baseline dataset can be modeled as samples $X = \{x_i\} \subseteq \mathbb{R}^d$ drawn from some true
 100 distribution $p(x)$, and the new dataset that we want to compare against the baseline can be modeled
 101 as samples $Y = \{y_i\} \subseteq \mathbb{R}^d$ drawn from a distribution $q(x)$. Here, x_i, y_i represent the embedding
 102 vectors of the problem statements in an embedding space \mathbb{R}^d . We define the **novelty** of the new
 103 dataset Y to be the KL-divergence between the distributions

$$104 \text{Novelty}(Y|X) = D_{KL}(q||p) = \int_{\mathbb{R}^d} q(\mathbf{x}) \log \frac{q(\mathbf{x})}{p(\mathbf{x})} d\mathbf{x} \quad (1)$$

105 Note that we take the KL-divergence of p with q as the null hypothesis because we want to reward
 106 datasets where $q(x)$ is large and $p(x)$ is small, indicating that the dataset contains problems not in
 107

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.

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

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, as 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 maximizes entropy, 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 $x_1, \dots, x_n, y_1, \dots, y_m \in \mathbb{R}^d$ where x and y are drawn from $p(x), q(x)$ respectively, we use the k-NN based estimator by Wang et al. (2009)

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

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 $x_1, \dots, x_N \in \mathbb{R}^d$, the Kozachenko-Leonenko estimator is

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

where ψ 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 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 allmpnet-base-v2 (Song et al., 2020) to embed the 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, preserving their relative geometry. We renormalize embeddings after projection so that distances 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 (374 problems) Austin et al. (2021). We also use [HumanEval \(164 problems\)](#) and [500 randomly chosen problems from the APPS dataset](#) Chen et al. (2021); Hendrycks et al. (2021). 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.

As shown in Figure 1, our results align with intuition, confirming the 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. We find that in practice, for practically useful comparisons, negativity never occurs; we include this case primarily to illustrate that our metric still reflects human intuition. Note that our comparisons can still work in such degenerate cases as we only care about *relative* differences between datasets.

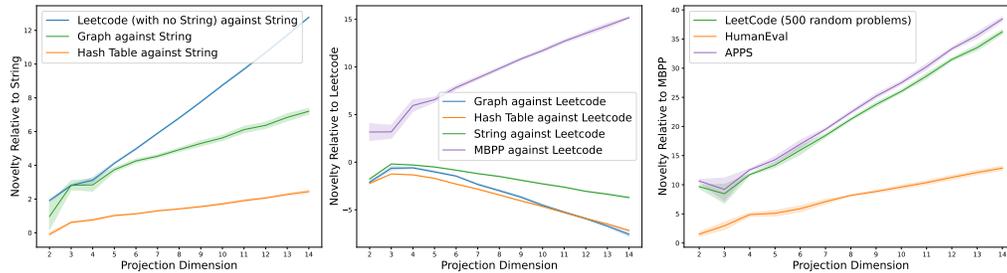


Figure 1: Left: The full Leetcode dataset has higher novelty than its Hash Table and String subsets as expected. Middle: The MBPP dataset has high novelty against the Leetcode dataset, whereas the Leetcode subdatasets have very little relative novelty as expected. Right: Leetcode and APPS have high novelty as their problems are harder and very different from MBPP, whereas HumanEval has low novelty as it is known to be more similar to MBPP. All plots show 95% confidence intervals.

3.2.2 DIFFERENTIAL ENTROPY METRIC VALIDATION

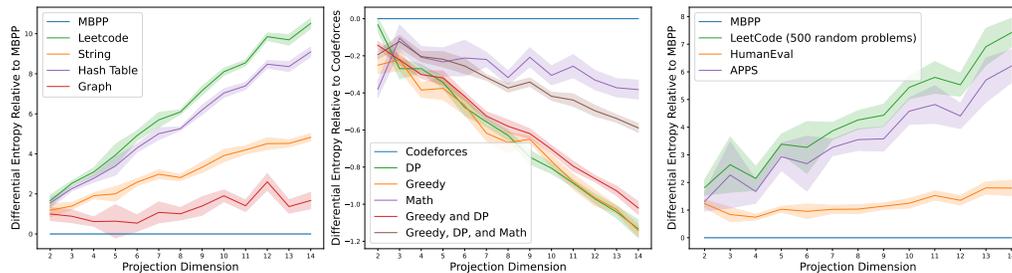


Figure 2: Left: Leetcode vs. MBPP entropy. MBPP shows lower entropy due to simpler, more repetitive problems. Middle: Codeforces vs. subsets. Full datasets have higher entropy than topic-specific subsets, except Math, which overlaps with others (e.g., DP, Greedy) and thus appears highly dispersed when isolated. Right: Leetcode and APPS have high diversity as their problems span many different computer science topics, whereas HumanEval has low diversity as the problems are easier and known to be more similar to MBPP. All plots show 95% confidence intervals.

In addition to previous 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 draw fair comparisons, we sample N points per dataset (without replacement) many times 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. The results show that the 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. We present ablations for all hyper-parameters used to compute these metrics in Appendix A.

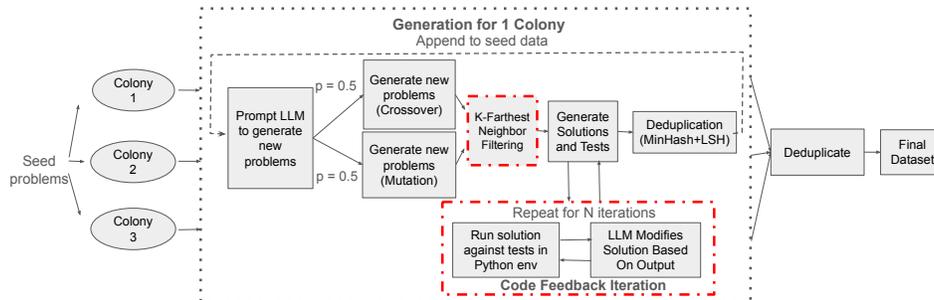


Figure 3: Generation Pipeline. Each colony receives a subset of the seed problems and applies mutation or crossover to them at each iteration. For each problem, it generates solutions and tests which go through multiple iterations of testing to ensure correctness. Deduplication removes similar problems 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 an end-to-end genetic algorithm for generating novel and diverse datasets from a seed dataset. The detailed algorithm is given in Appendix B, 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, encouraging diversity. We demonstrate the benefit of having multiple mutation difficulties in Section 5.2. Crossover prompts combine existing questions into new ones. Prompts are given in Appendix D.1, D.2. This example shows how mutation varies problem difficulty. Appendix F 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 F 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 D.3, D.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 (e.g., should we return None or -1 on an empty array input?). An example of such a problem is given in Appendix E. 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 D.5.

4.2 EXPERIMENTAL SETUP

We generate six datasets using GPT-4o (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 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.

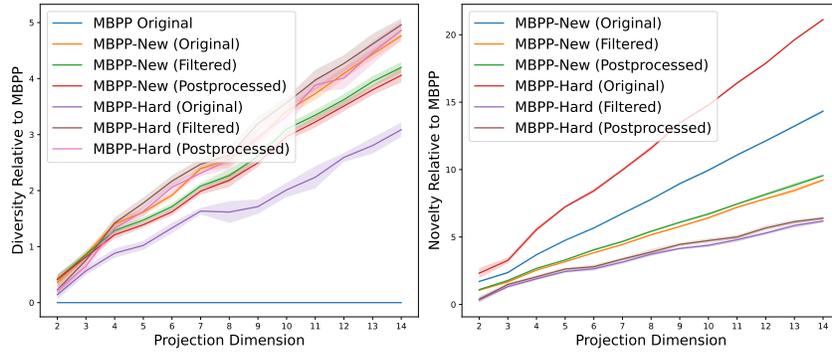
5 RESULTS AND ANALYSIS

We categorize generated problems as: (1) **Passing** (solution passes all tests), (2) **Failing** (fails ≥ 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 fairest comparison, omitting filtered versions (see Section 5.5). For each benchmark, we randomly sample 100 problems and manually verify that its solutions and test cases are fully correct and consistent with its problem statement.

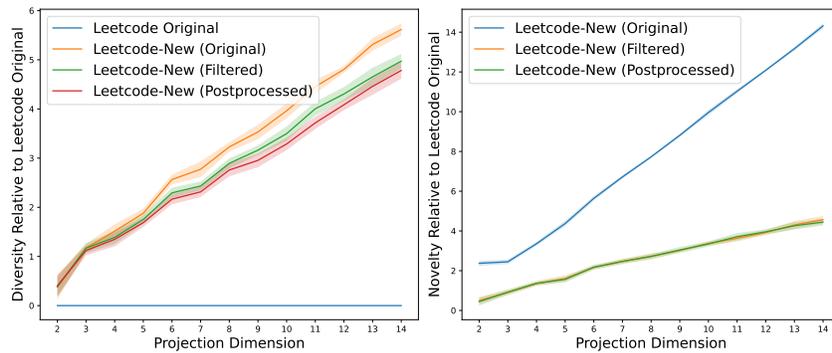
5.1 NOVELTY AND DIVERSITY ANALYSIS

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 attribute this to LLM memorization (Huang et al., 2025; Kiyomaru et al., 2024), since more novel

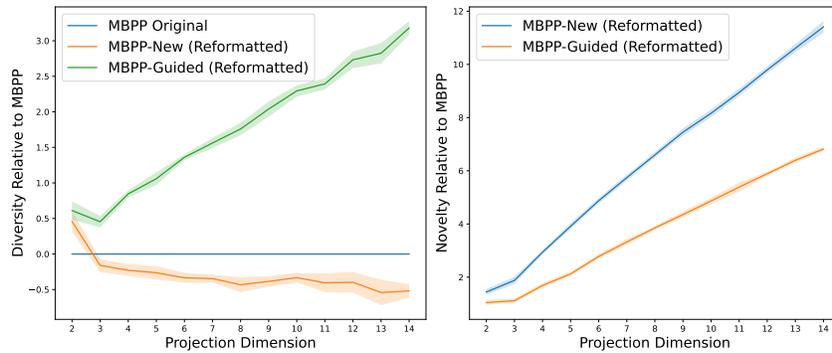
324
 325
 326
 327
 328
 329
 330
 331
 332
 333
 334
 335
 336
 337
 338
 339
 340
 341
 342
 343
 344
 345
 346
 347
 348
 349
 350
 351
 352
 353
 354
 355
 356
 357
 358
 359
 360
 361
 362
 363
 364
 365
 366
 367
 368
 369
 370
 371
 372
 373
 374
 375
 376
 377



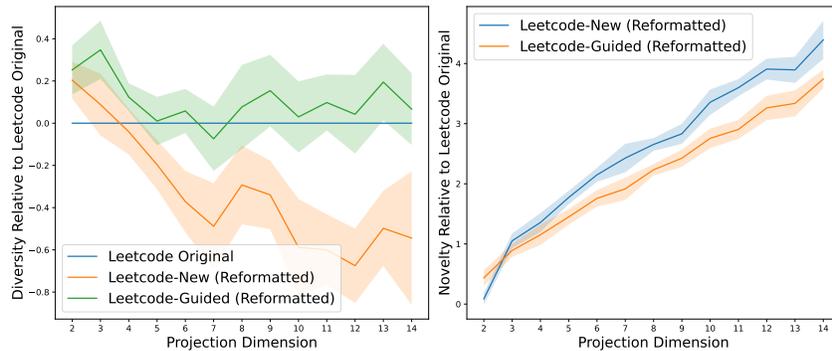
(a) Novelty and diversity in MBPP-New and MBPP-Hard. Using all three mutation types boosts novelty. MBPP-Hard has lower novelty than MBPP-Original but higher diversity.



(b) Leetcode-New shows greater novelty and diversity compared to the original Leetcode.



(c) MBPP-Guided exhibits higher diversity but reduced novelty compared to MBPP-New.



(d) Leetcode-Guided exhibits higher diversity but reduced novelty.

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

378 problems are often out-of-distribution and harder for the model to solve. All UMAP simulations use
 379 80 neighbors and a minimum distance of 0.1, except for Figure 4d, which uses 30 neighbors due to
 380 smaller dataset size. In general, our results are not sensitive to UMAP hyperparameters.

381 Figures 4c and 4d show that k-farthest-neighbor filtering improves dataset novelty and diversity.
 382 This comes at the cost of generating easier problems (Table 2), highlighting another advantage of
 383 InfoSynth in controlling the novelty–diversity–difficulty tradeoff. Empirically, this arises because
 384 the generator struggles to produce difficult problems unless they are conceptually aligned with the
 385 seeds. We also observe a tradeoff between novelty and diversity: highly diverse datasets tend to con-
 386 centrate around low-density regions in the seed-embedding distribution, which increases novelty but
 387 reduces diversity. Our results also show that filtering and post-processing reliably improve diversity.
 388

389 5.2 EFFECT OF VARYING MUTATION DIFFICULTIES

390 Table 2 shows MBPP-Hard scores are 8%-15% lower than MBPP-Original across most models,
 391 suggesting hard mutations effectively raise difficulty. This comes with a tradeoff of the dataset
 392 having reduced diversity and novelty as the problems tend to be concentrated around fewer, but
 393 more challenging topics. Hence, the set of mutation difficulties can be carefully chosen in InfoSynth
 394 in order to control the difficulty of the produced benchmark.
 395

396 5.3 EFFECT OF ITERATIVE CODE FEEDBACK

397 We find that passing solution-test pairs increase by 20% over 5 feedback iterations, showing the
 398 effectiveness of code iteration in producing robust problems. Error rates drop as the LLM fixes
 399 syntax/runtime issues, though the unparseable rate rises slightly due to occasional formatting failures.
 400 Appendix C shows feedback curves. Three iterations are typically ideal; further iterations yield
 401 marginal gains not worth the extra inference cost. We also find that iterative feedback acts as chain-
 402 of-thought (CoT) reasoning (Wei et al., 2022), as the model leverages the full feedback history to
 403 refine solutions/tests, lowering both error and failure rates. An example of this is in Appendix G.
 404
 405

406 5.4 CATEGORIZING FILTERED-OUT PROBLEMS

407 In general, we don’t observe any major biases in the topics that make it past the filtering step,
 408 however, there are 2 notable types of problems that get filtered more often than others:
 409

- 410 • Problems that are a crossover of problems that were already generated by crossover. This tends to
 411 create very difficult problems with many constraints, making it difficult for the generator LLM to
 412 generate accurate solutions and tests that adhere to every step. We interpret this as a limitation of
 413 current models’ ability to perform long-horizon tasks, rather than a limitation of our pipeline.
- 414 • Problems involving significant numerical calculations (e.g. computing the nth Delannoy number);
 415 although the generator LLM often generates correct solutions for such problems, it struggles to
 416 make tests with accurate numerical assertions.

417 5.5 EFFECT OF THE POSTPROCESSING STEP

418 Appendix E shows two post-processed examples, where the model resolves ambiguous edge cases
 419 and sometimes rephrases statements more concisely without losing information. Table 2 shows
 420 5–15% accuracy gains across most test-taker models, confirming that post-processing reduces am-
 421 biguity. Manual verification of 100 problems per dataset further shows 100% of post-processed
 422 problems are correctly reformatted without altering the core question.
 423
 424

425 5.6 RELATING DIVERSITY AND TOPIC COVERAGE

426 For each problem in MBPP Original, MBPP-New, and MBPP-Guided, we prompted GPT-4o-mini
 427 (Hurst et al., 2024) to provide up to 3 topic labels describing the problem, similar to the approach
 428 used by Zhao et al. (2025). The list of available topic labels was taken from the Leetcode dataset
 429 (kaysss, 2025). Figure 5 shows the results. The prompt is given in Appendix H.
 430

431 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

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

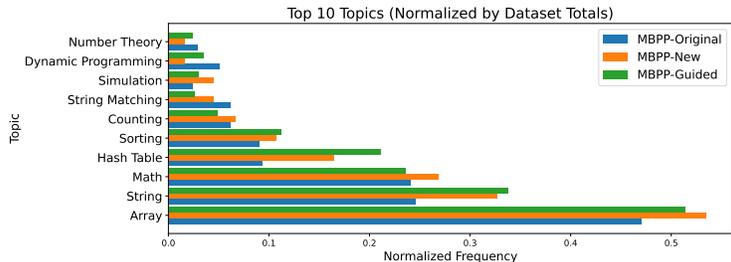


Figure 5: Fraction of problems relating to each topic for the 10 most common topics

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.

5.7 COMPARISON TO PREVIOUS GENERATION METHODS

We compare InfoSynth with GeneticInstruct and KodCode in generating novel and diverse problems. All three pipelines use Leetcode-based seeds, so novelty and diversity are measured against Leetcode Original. Since GeneticInstruct and KodCode use multiple seed datasets, we pool all reformatted problems generated by InfoSynth to form a combined dataset. Figure 6 shows the results.

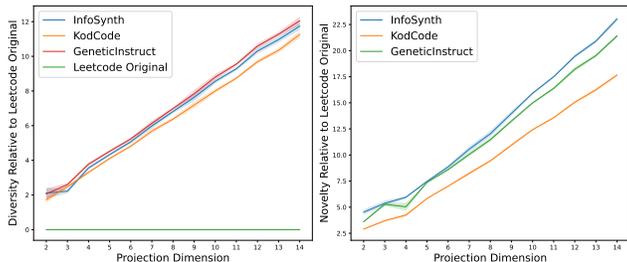


Figure 6: Novelty and diversity of various problem generation pipelines

We find that InfoSynth and GeneticInstruct exhibit very similar diversity, both slightly higher than KodCode. Additionally, InfoSynth stands out as the most novel dataset. Importantly, this comparison measures novelty against the Leetcode seed dataset *used by InfoSynth*, rather than the seed datasets for GeneticInstruct or KodCode, highlighting InfoSynth’s superior ability to generate problems that differ substantially from its seed data.

5.8 CHOOSING AN EMBEDDING MODEL

We test various embeddings model on our datasets. Figures 7, 8 show that the relative novelty and diversity of datasets remains similar across embedding models despite some fluctuations in the magnitudes of those differences.

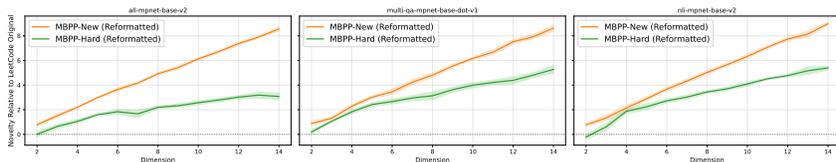


Figure 7: Novelty of datasets for various embedding models

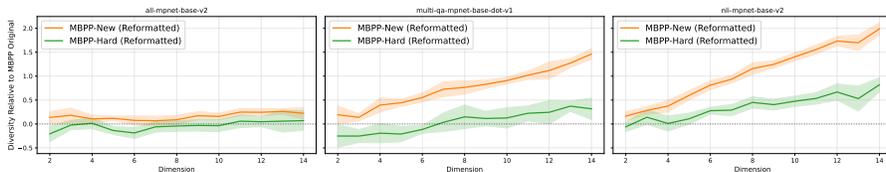


Figure 8: Diversity of datasets for various embedding models

Table 2: Test-taker performance on datasets

Model	Dataset	Filtered			Postprocessed		
		%Pass	%Fail	%Err	%Pass	%Fail	%Err
Qwen2.5-7b-Instruct	MBPP-Original	52.67	45.99	1.34	60.43	37.70	1.87
	MBPP-New	34.32	61.60	4.08	45.83	51.58	2.60
	MBPP-Guided	-	-	-	54.55	43.88	1.57
	MBPP-Hard	19.82	72.52	7.66	31.08	62.16	6.76
	MBPP-Hard-Guided	-	-	-	39.49	56.48	4.03
	Leetcode-Original	-	-	-	11.84	83.77	4.39
	Leetcode-New	25.29	74.12	0.59	25.88	72.35	1.76
	Leetcode-Guided	-	-	-	29.05	70.39	0.56
Qwen2.5-3b-Coder	MBPP-Original	46.79	48.93	4.28	52.41	42.51	5.08
	MBPP-New	31.35	55.47	13.17	38.78	49.54	11.69
	MBPP-Guided	-	-	-	40.56	49.13	10.31
	MBPP-Hard	21.62	54.95	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
GPT-4.1-Mini	MBPP-Original	58.02	36.10	5.88	66.04	30.48	0.00
	MBPP-New	56.96	40.63	2.41	67.35	29.68	2.97
	MBPP-Guided	-	-	-	72.03	22.38	5.59
	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
	Leetcode-New	45.28	38.24	16.47	50.00	40.59	9.41
	Leetcode-Guided	-	-	-	48.60	46.37	5.03
Gemini-2.0-Flash	MBPP-Original	64.97	35.03	0.00	68.72	31.02	0.26
	MBPP-New	53.99	45.64	0.37	63.64	36.18	0.19
	MBPP-Guided	-	-	-	72.03	26.92	1.05
	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
	Leetcode-Guided	-	-	-	46.37	50.28	3.35
Claude 3.7 Sonnet	MBPP-Original	63.37	36.63	0.00	70.86	29.14	0.00
	MBPP-New	55.29	44.71	0.00	64.75	35.25	0.00
	MBPP-Guided	-	-	-	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
o4-mini	MBPP-Original	66.58	33.42	0.00	70.05	29.68	0.27
	MBPP-New	58.26	41.19	0.56	70.13	29.68	0.19
	MBPP-Guided	-	-	-	77.27	22.55	0.17
	MBPP-Hard	47.75	49.55	2.70	63.06	32.88	4.05
	MBPP-Hard-Guided	-	-	-	72.61	26.11	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

6 CONCLUSION

In this paper, we introduced InfoSynth, a novel framework to efficiently calculate the diversity and novelty of new benchmarks in a more-efficient and cost-effective manner. Moreover, we demonstrated the effectiveness of InfoSynth in generating high-quality, novel, and diverse synthetic coding datasets from seed data. We hope that future work will leverage our ideas 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 B). All prompts used for mutation, crossover, and verification are listed in Appendix D, and examples of generated problems are in Appendix F. Post-processing steps and deduplication methods are described in Section 4 and Appendix E. 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. *ArXiv preprint*, 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. Codet: Code generation with generated tests. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. URL <https://openreview.net/pdf?id=ktrw68Cmu9c>.
- 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. *ArXiv preprint*, 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. *ArXiv preprint*, 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. In Kevin Duh, Helena Gomez, and Steven Bethard (eds.), *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 8706–8719, Mexico City, Mexico, 2024a. Association for Computational Linguistics. URL <https://aclanthology.org/2024.naacl-long.482>.
- Chunyuan Deng, Yilun Zhao, Xiangru Tang, Mark Gerstein, and Arman Cohan. Investigating data contamination in modern benchmarks for large language models. In Kevin Duh, Helena Gomez, and Steven Bethard (eds.), *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 8706–8719, Mexico City, Mexico, 2024b. Association for Computational Linguistics. URL <https://aclanthology.org/2024.naacl-long.482>.
- Yuyang Ding, Xinyu Shi, Xiaobo Liang, Juntao Li, Qiaoming Zhu, and Min Zhang. Unleashing Reasoning Capability of LLMs via Scalable Question Synthesis from Scratch. *ArXiv preprint*, abs/2410.18693, 2024. URL <https://arxiv.org/abs/2410.18693>.
- Shahriar Golchin and Mihai Surdeanu. Data Contamination Quiz: A Tool to Detect and Estimate Contamination in Large Language Models. *ArXiv preprint*, abs/2311.06233, 2023. URL <https://arxiv.org/abs/2311.06233>.

- 594 Dan Hendrycks, Steven Basart, Saurav Kadavath, Mantas Mazeika, Akul Arora, Ethan Guo, Collin
595 Burns, Samir Puranik, Horace He, Dawn Song, and Jacob Steinhardt. Measuring coding challenge
596 competence with apps. *ArXiv preprint*, abs/2105.09938, 2021. URL [https://arxiv.org/
597 abs/2105.09938](https://arxiv.org/abs/2105.09938).
- 598 Kaixuan Huang, Jiacheng Guo, Zihao Li, Xiang Ji, Jiawei Ge, Wenzhe Li, Yingqing Guo, Tianle
599 Cai, Hui Yuan, Runzhe Wang, Yue Wu, Ming Yin, Shange Tang, Yangsibo Huang, Chi Jin,
600 Xinyun Chen, Chiyuan Zhang, and Mengdi Wang. MATH-Perturb: Benchmarking LLMs’ Math
601 Reasoning Abilities against Hard Perturbations. *ArXiv preprint*, abs/2502.06453, 2025. URL
602 <https://arxiv.org/abs/2502.06453>.
- 603 Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Os-
604 trow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card, 2024. URL
605 <https://arxiv.org/abs/2410.21276>.
- 606 Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando
607 Solar-Lezama, Koushik Sen, and Ion Stoica. LiveCodeBench: Holistic and Contamination Free
608 Evaluation of Large Language Models for Code. *ArXiv preprint*, abs/2403.07974, 2024. URL
609 <https://arxiv.org/abs/2403.07974>.
- 610 kaysss. leetcode-problem-set. [https://huggingface.co/datasets/kaysss/
611 leetcode-problem-set](https://huggingface.co/datasets/kaysss/leetcode-problem-set), 2025. Accessed: 2025-05-11.
- 612 Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. In Yoshua Bengio and Yann
613 LeCun (eds.), *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB,
614 Canada, April 14-16, 2014, Conference Track Proceedings*, 2014. URL [http://arxiv.org/
615 abs/1312.6114](http://arxiv.org/abs/1312.6114).
- 616 Hirokazu Kiyomaru, Issa Sugiura, Daisuke Kawahara, and Sadao Kurohashi. A Comprehensive
617 Analysis of Memorization in Large Language Models. In Saad Mahamood, Minh Le Nguyen,
618 and Daphne Ippolito (eds.), *Proceedings of the 17th International Natural Language Generation
619 Conference, INLG 2024, Tokyo, Japan, September 23 - 27, 2024*, pp. 584–596. Association for
620 Computational Linguistics, 2024. doi: 10.18653/V1/2024.INLG-MAIN.45. URL [https://
621 doi.org/10.18653/v1/2024.inlg-main.45](https://doi.org/10.18653/v1/2024.inlg-main.45).
- 622 Alexander Kraskov, Harald Stögbauer, and Peter Grassberger. Estimating mutual information.
623 *Physical Review E*, 69(6), 2004. ISSN 1550-2376. doi: 10.1103/physreve.69.066138. URL
624 <http://dx.doi.org/10.1103/PhysRevE.69.066138>.
- 625 Xiang Lisa Li, Farzaan Kaiyom, Evan Zheran Liu, Yifan Mai, Percy Liang, and Tatsunori
626 Hashimoto. AutoBench: Towards Declarative Benchmark Construction. In *The Thirteenth In-
627 ternational Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025*.
628 OpenReview.net, 2025a. URL <https://openreview.net/forum?id=ymt4crbbXh>.
- 629 Yisen Li, Lingfeng Yang, Wenxuan Shen, Pan Zhou, Yao Wan, Weiwei Lin, and Dongping Chen.
630 CrowdSelect: Synthetic Instruction Data Selection with Multi-LLM Wisdom. *ArXiv preprint*,
631 abs/2503.01836, 2025b. URL <https://arxiv.org/abs/2503.01836>.
- 632 Jiatae Liu, Yiqin Zhu, Kaiwen Xiao, Qiang Fu, Xiao Han, Wei Yang, and Deheng Ye. RLTF: Re-
633 inforcement Learning from Unit Test Feedback. *ArXiv preprint*, abs/2307.04349, 2023. URL
634 <https://arxiv.org/abs/2307.04349>.
- 635 Jiawei Liu, Thanh Nguyen, Mingyue Shang, Hantian Ding, Xiaopeng Li, Yu Yu, Varun Kumar, and
636 Zijian Wang. Learning Code Preference via Synthetic Evolution. *ArXiv preprint*, abs/2410.03837,
637 2024a. URL <https://arxiv.org/abs/2410.03837>.
- 638 Zhihan Liu, Shenao Zhang, and Zhaoran Wang. DSTC: Direct Preference Learning with Only Self-
639 generated Tests and Code to Improve Code LMs. *ArXiv preprint*, abs/2411.13611, 2024b. URL
640 <https://arxiv.org/abs/2411.13611>.
- 641 Somshubra Majumdar, Vahid Noroozi, Sean Narenthiran, Aleksander Ficek, Jagadeesh Balam, and
642 Boris Ginsburg. Genetic Instruct: Scaling up Synthetic Generation of Coding Instructions for
643 Large Language Models. *ArXiv preprint*, abs/2407.21077, 2024. URL [https://arxiv.org/
644 abs/2407.21077](https://arxiv.org/abs/2407.21077).

- 648 Leland McInnes, John Healy, Nathaniel Saul, and Lukas Großberger. UMAP: Uniform Manifold
649 Approximation and Projection. *J. Open Source Softw.*, 3(29):861, 2018. doi: 10.21105/JOSS.
650 00861. URL <https://doi.org/10.21105/joss.00861>.
- 651 open-rl. Codeforces Problems Dataset. [https://huggingface.co/datasets/
652 open-rl/codeforces](https://huggingface.co/datasets/open-rl/codeforces), 2024. Accessed: 2025-05-11.
- 654 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal Policy
655 Optimization Algorithms. *ArXiv preprint*, abs/1707.06347, 2017. URL [https://arxiv.
656 org/abs/1707.06347](https://arxiv.org/abs/1707.06347).
- 657 Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. MpNet: Masked and permuted
658 pre-training for language understanding. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Had-
659 sell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), *Advances in Neural Information Process-
660 ing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS
661 2020, December 6-12, 2020, virtual*, 2020. URL [https://proceedings.neurips.cc/
662 paper/2020/hash/c3a690be93aa602ee2dc0ccab5b7b67e-Abstract.html](https://proceedings.neurips.cc/paper/2020/hash/c3a690be93aa602ee2dc0ccab5b7b67e-Abstract.html).
- 664 Florian Tambon, Amin Nikanjam, Cyrine Zid, Foutse Khomh, and Giuliano Antoniol. TaskEval:
665 Assessing Difficulty of Code Generation Tasks for Large Language Models. *ArXiv preprint*,
666 abs/2407.21227, 2024. URL <https://arxiv.org/abs/2407.21227>.
- 667 Sang T. Truong, Yuheng Tu, Percy Liang, Bo Li, and Sanmi Koyejo. Reliable and Efficient
668 Amortized Model-based Evaluation. *ArXiv preprint*, abs/2503.13335, 2025. URL [https:
669 //arxiv.org/abs/2503.13335](https://arxiv.org/abs/2503.13335).
- 670 Qing Wang, Sanjeev R. Kulkarni, and Sergio Verdú. Divergence estimation for multidimensional
671 densities via k-nearest-neighbor distances. *IEEE Trans. Inf. Theory*, 55(5):2392–2405, 2009. doi:
672 10.1109/TIT.2009.2016060. URL <https://doi.org/10.1109/TIT.2009.2016060>.
- 674 Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and
675 Hannaneh Hajishirzi. Self-instruct: Aligning language models with self-generated instructions.
676 In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual
677 Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 13484–
678 13508, Toronto, Canada, 2023. Association for Computational Linguistics. doi: 10.18653/v1/
679 2023.acl-long.754. URL <https://aclanthology.org/2023.acl-long.754>.
- 680 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi,
681 Quoc V. Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language
682 models. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh
683 (eds.), *Advances in Neural Information Processing Systems 35: Annual Conference on Neural
684 Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - De-
685 cember 9, 2022*, 2022. URL [http://papers.nips.cc/paper_files/paper/2022/
686 hash/9d5609613524ecf4f15af0f7b31abca4-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/9d5609613524ecf4f15af0f7b31abca4-Abstract-Conference.html).
- 687 Yunhui Xia, Wei Shen, Yan Wang, Jason Klein Liu, Huifeng Sun, Siyue Wu, Jian Hu, and Xiaolong
688 Xu. LeetCodeDataset: A Temporal Dataset for Robust Evaluation and Efficient Training of Code
689 LLMs. *ArXiv preprint*, abs/2504.14655, 2025. URL [https://arxiv.org/abs/2504.
690 14655](https://arxiv.org/abs/2504.14655).
- 691 Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, Qing-
692 wei Lin, and Daxin Jiang. Wizardlm: Empowering large pre-trained language models to fol-
693 low complex instructions. In *The Twelfth International Conference on Learning Representa-
694 tions, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024. URL [https:
695 //openreview.net/forum?id=CfXh93NDgH](https://openreview.net/forum?id=CfXh93NDgH).
- 696 Zhangchen Xu, Yang Liu, Yueqin Yin, Mingyuan Zhou, and Radha Poovendran. KodCode: A Di-
697 verse, Challenging, and Verifiable Synthetic Dataset for Coding. *ArXiv preprint*, abs/2503.02951,
698 2025. URL <https://arxiv.org/abs/2503.02951>.
- 700 Huaye Zeng, Dongfu Jiang, Haozhe Wang, Ping Nie, Xiaotong Chen, and Wenhui Chen. ACE-
701 CODER: Acing Coder RL via Automated Test-case Synthesis. *ArXiv preprint*, abs/2502.01718,
2025. URL <https://arxiv.org/abs/2502.01718>.

702 Hugh Zhang, Jeff Da, Dean Lee, Vaughn Robinson, Catherine Wu, William Song, Tiffany
703 Zhao, Pranav Raja, Charlotte Zhuang, Dylan Slack, Qin Lyu, Sean Hendryx, Russell Ka-
704 plan, Michele Lunati, and Summer Yue. A careful examination of large language model
705 performance on grade school arithmetic. In Amir Globersons, Lester Mackey, Danielle Bel-
706 grave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang (eds.), *Advances*
707 *in Neural Information Processing Systems 38: Annual Conference on Neural Informa-*
708 *tion Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 -*
709 *15, 2024*, 2024. URL [http://papers.nips.cc/paper_files/paper/2024/](http://papers.nips.cc/paper_files/paper/2024/hash/53384f2090c6a5cac952c598fd67992f-Abstract-Datasets_and_Benchmarks_Track.html)
710 [hash/53384f2090c6a5cac952c598fd67992f-Abstract-Datasets_and_](http://papers.nips.cc/paper_files/paper/2024/hash/53384f2090c6a5cac952c598fd67992f-Abstract-Datasets_and_Benchmarks_Track.html)
711 [Benchmarks_Track.html](http://papers.nips.cc/paper_files/paper/2024/hash/53384f2090c6a5cac952c598fd67992f-Abstract-Datasets_and_Benchmarks_Track.html).

712 Xueliang Zhao, Wei Wu, Jian Guan, and Lingpeng Kong. PromptCoT: Synthesizing Olympiad-
713 level Problems for Mathematical Reasoning in Large Language Models. *ArXiv preprint*,
714 [abs/2503.02324](https://arxiv.org/abs/2503.02324), 2025. URL <https://arxiv.org/abs/2503.02324>.

715 Terry Yue Zhuo, Minh Chien Vu, Jenny Chim, Han Hu, Wenhao Yu, Ratnadira Widayarsi, Imam
716 Nur Bani Yusuf, Haolan Zhan, Junda He, Indraneil Paul, Simon Brunner, Chen Gong, Thong
717 Hoang, Armel Randy Zebaze, Xiaoheng Hong, Wen-Ding Li, Jean Kaddour, Ming Xu, Zhi-
718 han Zhang, Prateek Yadav, Naman Jain, Alex Gu, Zhoujun Cheng, Jiawei Liu, Qian Liu, Zi-
719 jian Wang, David Lo, Binyuan Hui, Niklas Muennighoff, Daniel Fried, Xiaoning Du, Harm
720 de Vries, and Leandro von Werra. BigCodeBench: Benchmarking Code Generation with Di-
721 verse Function Calls and Complex Instructions. *ArXiv preprint*, [abs/2406.15877](https://arxiv.org/abs/2406.15877), 2024. URL
722 <https://arxiv.org/abs/2406.15877>.

723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

A SENSITIVITY ANALYSIS FOR k AND UMAP PARAMETERS

We present ablations for hyper-parameters used in our experiments. In general, our results generalize robustly across a wide range of hyper-parameters. Note that in the cases where changing a hyperparameter affects the novelty or diversity, it affects all datasets by roughly the same amount. Since our analysis looks at *relative* differences between benchmarks, our conclusions are still robust.

Specifically, we ablate three hyperparameters in Figure 9:

- k : the number of k -NN neighbors considered in the estimators
- `n_neighbors`: the number of neighbors used in the UMAP algorithm
- `min_dist`: the minimum distance used in the UMAP algorithm

B 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 $N = 1000$, $N_c = 10$, $B_s = 30$, $C = 3$, $B_c = 4$, $N_{it} = 3$.

The algorithm is given in Algorithm 1

C CODE-FEEDBACK RESULTS

Figures 10, 11 show how the proportion of problems that pass, fail, error, and are unparseable 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.

D PROBLEM GENERATION PIPELINE PROMPTS

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

D.1 MUTATION PROMPTS

`% easy`

Please decrease the difficulty of the given programming test question a bit.

The new problem should be conceptually similar to the given question,

but should not simply paraphrase it. Do not provide any hints, solutions

or outputs. Only one new instruction is allowed.

Original Question: {instruction}

New Question:

`% medium`

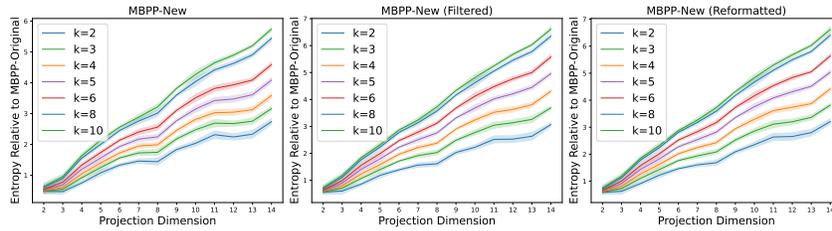
Please create a new programming problem of the same difficulty as the

given programming test question. The new problem should be conceptually

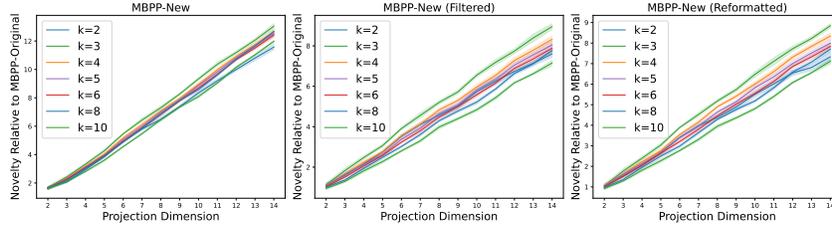
similar to the given question, but should not simply paraphrase it.

Do not provide any hints, solutions or outputs. Only one new instruction is allowed.

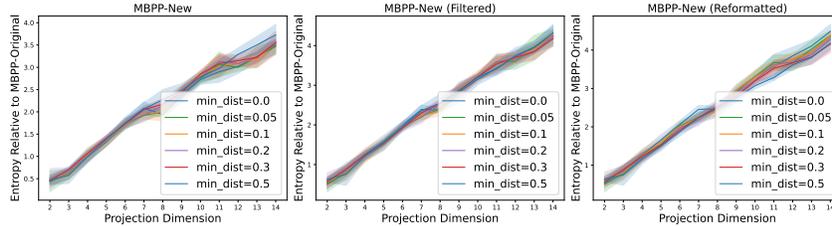
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863



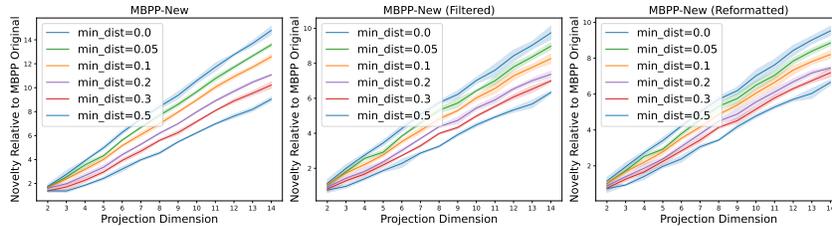
(a) Diversity changes by a similar amount for all benchmarks as k changes



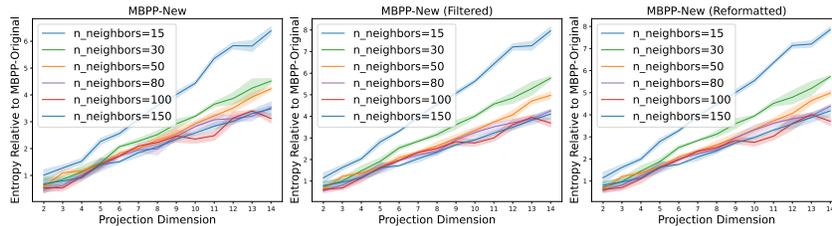
(b) Novelty is almost invariant to fluctuations in k



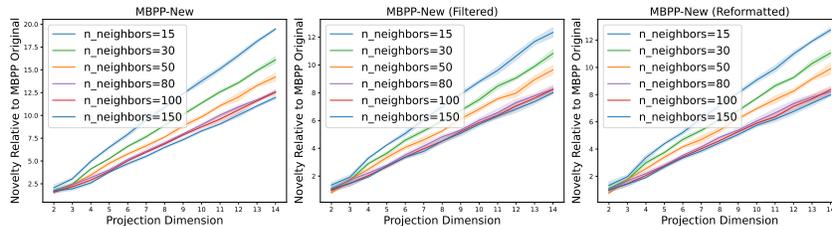
(c) Diversity is almost invariant to fluctuations in $n_neighbors$



(d) The novelty changes by a similar amount for all benchmarks as $n_neighbors$ changes



(e) Diversity changes by a similar amount for all benchmarks as min_dist changes



(f) Novelty changes by a similar amount for all benchmarks as min_dist changes

Figure 9: Sensitivity analysis across all hyperparameters.

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

Algorithm 1: Problem Generation with Evolutionary Strategies

```

Input:  $N$ : total number of problems to generate
 $N_c$ : number of colonies
 $B_s$ : number of problems to sample as the seed data for each colony
 $C$ : 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
combine
 $N_{it}$ : number of code feedback iterations

Generate (seedData, numSamples):
  Initialize problems  $\leftarrow \emptyset$ 
  for colony  $\leftarrow 1$  to  $N_c$  do
     $N_s \leftarrow N/N_c$ 
    colonySeedData  $\leftarrow$  Random sample of size  $B_s$  from seedData
    problems  $\leftarrow$  problems  $\cup$  EvolveColony (colonySeedData,  $N_s$ )
    deduplicate(problems)
  end
  return problems

EvolveColony (seedData,  $N_s$ ):
  newProblems  $\leftarrow \emptyset$ 
  while  $|newProblems| < N_s$  do
    operation  $\leftarrow$  "mutation" with  $p = 0.5$ , "crossover" with  $p = 0.5$ 
    if operation == "mutation" then
      problem  $\leftarrow$  random sample of size 1 from seedData
      problems  $\leftarrow$  problems  $\cup$  mutate(problem)
    else if operation == "crossover" then
      batch  $\leftarrow$  random sample of size  $B_c$  from seedData
      problems  $\leftarrow$  problems  $\cup$  crossover( $C$ , batch)
    end
    if k-farthest-neighbor-filtering-enabled then
       $U \leftarrow newProblems \cup seedData$ 
      problems  $\leftarrow$  select  $K$  problems with least cosine similarity to  $U$ 
    end
    newProblems  $\leftarrow newProblems \cup problems$ 
    deduplicate(newProblems)
    seedData  $\leftarrow seedData \cup problems$ 
    deduplicate(seedData)
  end
  return newProblems

GenerateSolutionsTests (problem):
  GenerateTests(problem)
  GenerateSolutions(problem)
  for  $i = 1$  to  $N_{it}$  do
    Run tests against solutions
    Feed output back into LLM to modify tests and solution
  end

```

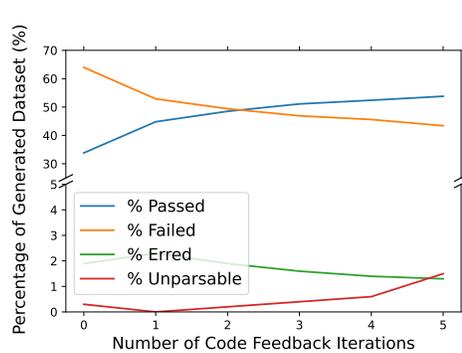


Figure 10: Self-Verification for MBPP-New

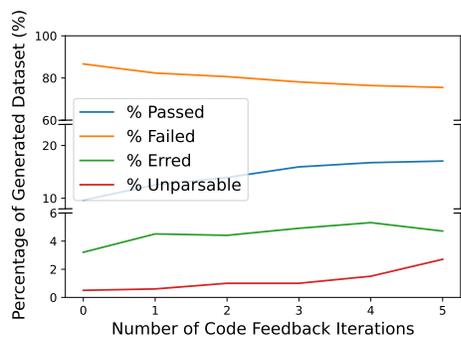


Figure 11: Self-Verification for Leetcode-New

Original Question: {instruction}
New Question:

% hard

Please increase the difficulty of the given programming test question a bit.

Do not provide any hints, solutions or outputs. Only one new instruction is allowed.

Original Question: {instruction}
New Question:

D.2 CROSSOVER PROMPT

I will provide you with a set of coding questions. Please give me a new coding question that combines core concepts from two or more of the given questions.

Please ensure that the new question is novel and does not simply paraphrase

any of the problems I am giving you. Do not include any extra information

that would help a test-taker other than the problem statement itself.

Question 1:

{instruction 1}

Question 2:

{instruction 2}

...

New Question:

D.3 SOLUTION & TEST GENERATION PROMPT

We note that our prompts are similar to those used by Xu et al. (2025).

You are an expert in Python coding.

Task:

Please answer the question and generate unit tests to verify your answer.

```

972 ## Output Format:
973 Your solution and unit tests should be presented in the format
974 within the
975 specified sections below. Ensure your code is within code blocks.
976 For the
977 tests, use pytest style by defining individual test functions
978 (without
979 classes) and using assert statements. Your tests should be
980 implementation
981 independent. Ensure that you include the <|Solution Begin|>,
982 <|Solution End|>, <|Test Begin|>, and <|Test End|> tags as
983 depicted. The
984 solution function must be named solution.

985 <|Solution Begin|>
986 {Solution Code in Python}
987 <|Solution End|>

988 <|Test Begin|>
989 {Unit Test Code in Python}
990 <|Test End|>

991
992 ## Example
993 Below is an example output format implementing a simple a + b
994 function.
995 <|Solution Begin|>
996 def add(a, b):
997     \"\"\"Returns the sum of a and b.\"\"\"
998     return a + b
999 <|Solution End|>

1000 <|Test Begin|>
1001 from solution import add
1002 def test_add_positive_numbers():
1003     assert add(2, 3) == 5
1004 def test_add_with_zero():
1005     assert add(0, 5) == 5
1006     assert add(5, 0) == 5
1007 def test_add_negative_numbers():
1008     assert add(-1, -1) == -2
1009 def test_add_mixed_sign_numbers():
1010     assert add(-1, 3) == 2
1011 <|Test End|>

1012 ## Question:
1013 {problem}
1014
1015
1016 D.4 SOLUTION & TEST GENERATION WITH ITERATIVE FEEDBACK PROMPT
1017
1018 You are an expert in Python coding.
1019 ## Task:
1020 Please answer the question and generate unit tests to verify your
1021 answer. The
1022 entire chat history of your previous attempts to generate
1023 questions and unit
1024 tests is presented below in the "Chat History" section, along with
1025 the output
of running your solution against your tests in a code execution
environment.

```

1026 Please modify only your tests and/or solution to be more correct.
1027

1028 ## Output Format:
1029 <Same as "Output Format" section above>
1030

1031 ## Chat History:
1032 Attempt 1 Solution:
1033 {attempt 1 solution}

1034 Attempt 1 Code Execution Output:
1035 {attempt 1 code output}
1036

1037 Attempt 2 Solution:
1038 {Attempt 2 solution}
1039

1040 Attempt 2 Code Execution Output:
1041 {attempt 2 code output}

1042 ...

1043 ## Question:
1044 {problem}

1045 1046 1047 D.5 POSTPROCESSING PROMPTS

1048 You are an expert in Python coding. Here is a coding problem
1049 with associated test cases. Please rephrase the question so it
1050 describes what the user should output for edge-cases without
1051 changing the essence of the problem. Add as little information
1052 as possible, only describing what the user should output for
1053 edge-cases that cannot be inferred from the problem description.
1054 Do not include anything except for the rewritten problem in your
1055 response and do not include the test cases.

1056 Question:
1057 {question}
1058 Tests:
1059 {tests}

1060 1061 E EXAMPLE OF POST-PROCESSED PROBLEMS

1062 1063 **MBPP-New Dataset Original Problem:**

1064
1065 Write a function that filters a list of usernames stored in a
1066 dictionary, returning only those associated with students who fall
1067 within a specified age range.

1068 **Post-Processed Version:**

1069
1070 Write a function that filters a list of usernames stored in a
1071 dictionary, returning only those associated with students who fall
1072 within a specified age range. Ensure that the function returns an
1073 empty list when there are no students or when the input dictionary
1074 is empty.

1075 **Leetcode-New Dataset Original Problem:**

1076
1077 Alice and Bob are engaged in a strategic game on an infinite 2D
1078 plane with n points provided by their coordinates in two integer
1079 arrays, $xCoord$ and $yCoord$. They take turns, starting with Alice,
attacking a point on the plane to capture it, with the condition

1080 that once a point is attacked, it is removed permanently from the
1081 game, and they have to remove exactly 1 point per turn.
1082

1083 The winner is the player who either removes the point that leaves
1084 no rectangle capable of being formed using the remaining points on
1085 their turn or can force the scenario by optimal play such that the
1086 opponent has no such move left on their subsequent turns.

1087 Given the arrays `xCoord` and `yCoord`, along with knowledge of
1088 optimal strategies for both players, determine if Alice, who
1089 starts the game, can always guarantee a win. Return `true` if Alice
1090 has a winning strategy, or `false` if Bob can always force a win
1091 even with Alice starting first.
1092

1093 **Post-Processed Version:**

1094 Alice and Bob are playing a strategic game on an infinite 2D
1095 plane with `n` points defined by their coordinates in two integer
1096 arrays, `xCoord` and `yCoord`. They alternately take turns, with Alice
1097 starting first, to attack and permanently remove exactly one point
1098 at a time. The objective is for a player to leave no possibility
1099 of forming a rectangle using any four of the remaining points. A
1100 player wins if they achieve this or if they can force a scenario
1101 where the opponent has no such moves left. Given the arrays `xCoord`
1102 and `yCoord`, determine if Alice has a guaranteed winning strategy.
1103 Return `true` if Alice can always win, and `false` if Bob can always
1104 force a win or if there are no points to start with.
1105

1106 **F EXAMPLES OF GENERATED PROBLEMS**

1107
1108 We note that for conciseness, the examples shown in this section are before the postprocessing
1109 step. In this section, we illuminate the different mechanisms through which mutation and crossover
1110 generate interesting problems.
1111

1112 **F.1 MUTATION OPERATION: ADDING OR REMOVING CONSTRAINTS**

1113 This example shows how mutation creates three variants of the question by adding requirements
1114 or removing constraints from the problem. We can see that all three problems require similar
1115 conceptual understanding, but the harder ones simply require more code or bookkeeping.
1116
1117

1118 **Original Question:**

1119 Write a function to filter the height and width of students which
1120 are stored in a dictionary.
1121

1122 **Easy Mutation:**

1123 Write a function that filters a list of usernames stored in a
1124 dictionary, returning only those associated with students who fall
1125 within a specified age range. Ensure that the function returns an
1126 empty list when there are no students or when the input dictionary
1127 is empty.

1128 **Medium Mutation:**

1129 Create a function that filters user profiles based on a
1130 dictionary. The function should return a list of user IDs for
1131 profiles where the age is within a specified range (inclusive)
1132 and the profile only contains lowercase alphabetic characters. If
1133 there are no valid profiles that match the criteria, the function
should return an empty list.

1134 Hard Mutation:
1135 Create a function to filter user passwords from a dictionary,
1136 returning only those that are valid for students whose dimensions
1137 (height, width, and weight) are within a given range. Each valid
1138 password must include at least one uppercase letter, one lowercase
1139 letter, one digit, and one special character. Ensure the solution
1140 appropriately handles and returns an empty dictionary for cases
1141 where no users are provided or all entries are invalid due to
1142 dimension or password criteria.

1143 F.2 MUTATION OPERATION EXAMPLE: CREATIVE MODIFICATION 1144

1145 This example shows how mutation creates three variants of the question by creatively modifying the
1146 central idea of the problem itself. This differs from the previous example; in this case, the harder
1147 questions require a fundamental understanding of new topics.
1148

1149 Original Question:
1150 Write a function to find the perimeter of a rectangle.

1151 Easy Mutation:
1152 Write a function to find the area of a rectangle.

1153 Medium Mutation:
1154 Write a function to calculate the area of a trapezoid given its
1155 base lengths and height.

1156 Hard Mutation:
1157 Write a function to find the area of a rhombus given its diagonals
1158 and verify if the rhombus is also a square by using its side
1159 lengths.

1160 F.3 CROSSOVER OPERATION: COMBINING CONCEPTS 1161

1162 This example shows how crossover creates an interesting, novel question by combining two
1163 unrelated concepts. A key difference between crossover and mutation is that crossover does not
1164 introduce any new concepts or content into the generated problem as it draws from existing ones.
1165 This demonstrates the necessity of combining mutation and crossover into one pipeline; mutation
1166 introduces new concepts into the dataset while crossover takes existing concepts and uses them to
1167 create richer problems.
1168

1169 Seed Problems: 1170

1171 Question 1:
1172 There is a 50 x 50 chessboard with one knight and some pawns on
1173 it. You are given two integers k_x and k_y where (k_x, k_y) denotes
1174 the position of the knight, and a 2D array `positions` where
1175 `positions[i] = [x_i, y_i]` denotes the position of the pawns on the
1176 chessboard.
1177 Alice and Bob play a turn-based game, where Alice goes first. In
1178 each player's turn:
1179

1180 The player selects a pawn that still exists on the board and
1181 captures it with the knight in the fewest possible moves. Note
1182 that the player can select any pawn, it might not be one that can
1183 be captured in the least number of moves.
1184 In the process of capturing the selected pawn, the knight may pass
1185 other pawns without capturing them. Only the selected pawn can be
1186 captured in this turn.
1187

1188
1189 Alice is trying to maximize the sum of the number of moves made by
1190 both players until there are no more pawns on the board, whereas
1191 Bob tries to minimize them.
1192 Return the maximum total number of moves made during the game that
1193 Alice can achieve, assuming both players play optimally.
1194 Note that in one move, a chess knight has eight possible positions
1195 it can move to, as illustrated below. Each move is two cells in a
1196 cardinal direction, then one cell in an orthogonal direction.

1197
1198 Question 2:
1199 You are given an array points where points[i] = [xi, yi]
1200 represents the coordinates of a point on an infinite plane.
1201 Your task is to find the maximum area of a rectangle that:

1202 Can be formed using four of these points as its corners.
1203 Does not contain any other point inside or on its border.
1204 Has its edges parallel to the axes.

1205
1206 Return the maximum area that you can obtain or -1 if no such
1207 rectangle is possible.

1208 **Crossover Generated Question:**

1209
1210 Alice and Bob are playing a strategic game on an infinite 2D
1211 plane with n points defined by their coordinates in two integer
1212 arrays, xCoord and yCoord. They alternately take turns, with Alice
1213 starting first, to attack and permanently remove exactly one point
1214 at a time. The objective is for a player to leave no possibility
1215 of forming a rectangle using any four of the remaining points. A
1216 player wins if they achieve this or if they can force a scenario
1217 where the opponent has no such moves left. Given the arrays xCoord
1218 and yCoord, determine if Alice has a guaranteed winning strategy.
1219 Return true if Alice can always win, and false if Bob can always
1220 force a win or if there are no points to start with.

1221
1222 **F.4 CROSSOVER OPERATION: MULTISTEP REASONING**

1223
1224 This example shows how crossover sometimes creates problems that require multistep reasoning.
1225 The new problem effectively requires a test-taker to solve both the seed problems.

1226
1227 **Seed Problems:**

1228
1229 Question 1:
1230 Write a python function to find the sum of fourth power of first n
1231 odd natural numbers.

1232
1233 Question 2:
1234 Write a python function to find the sum of xor of all pairs of
1235 numbers in the given array.

1236 **Crossover Generated Question:**

1237
1238 Write a python function to compute the combined metric of two
1239 properties for a given array of integers:

1240
1241 1. The sum of the fourth power of the first 'k' odd numbers from
the array.

1242 2. The sum of the bitwise XOR of all pairs of numbers formed by
 1243 selecting one element from the first 'k' even numbers and one
 1244 element from the last 'k' elements in the array.
 1245

1246 The function should take two arguments - the array of integers
 1247 and an integer 'k'. Ensure that 'k' is less than or equal to half
 1248 the length of the array and that the array contains a sufficient
 1249 number of odd and even numbers.
 1250

1251 G EXAMPLE OF CHAIN-OF-THOUGHT THINKING

1252
 1253 This example shows how the model is able to reason about why its solution and tests are erroring
 1254 due to receiving feedback from the code execution environment.
 1255

1256 **Question:** Write a function that takes a list of strings and returns the second shortest string that
 1257 contains all vowels, sorted alphabetically. If no such string exists, return an empty string.
 1258

1259 Initial Solution:

```
1260 def solution(strings):
1261     """
1262     Returns the second shortest string containing all vowels,
1263     sorted alphabetically.
1264     """
1265     all_vowels = set('aeiou')
1266
1267     # Filter strings that contain all vowels
1268     filtered_strings = [s for s in strings if all_vowels <=
1269 set(s)]
1270
1271     # Sort them by length first and then alphabetically
1272     filtered_strings.sort(key=lambda x: (len(x), x))
1273
1274     # Return the second shortest string or empty string if it
1275     doesn't exist
1276     return filtered_strings[1] if len(filtered_strings) > 1 else
1277     ''
1278 {Tests Omitted for Conciseness}
```

1279 Solution after First Feedback iteration:

1280 Let's revise the test cases to ensure they correctly test the
 1281 function's behavior and check the solution implementation
 1282 and its handling of edge cases. First, we'll make sure the
 1283 solution properly handles strings with varying cases (uppercase
 1284 and lowercase) and that it verifies the full vowel set in a
 1285 case-insensitive manner. Additionally, let's adjust the tests to
 1286 cover more scenarios and better reflect the conditions described
 1287 in the question. Here's the revised version of the solution and
 1288 unit tests:

```
1289 {Solutions Omitted for Conciseness}
1290 {Tests Omitted for Conciseness}
```

1291 Solution after Second Feedback iteration:

1292
 1293 We need to ensure that our solution and tests align correctly with
 1294 the task. Let's refine both the solution and the tests so that
 1295 they pass correctly according to the given problem statement. The
 solution needs to ensure that we filter strings containing all

1296 vowels in any case, sort them, and return the second shortest. The
1297 test cases should reflect this logic.
1298

1299 Here's the refined solution and tests:

1300 {Solutions Omitted for Conciseness}
1301 {Tests Omitted for Conciseness}

1304 H TOPIC-LABELING PROMPT

1306 You are a precise tagger for coding interview problems.
1307

1308 Given a problem statement and solution, choose up to THREE topics
1309 from the bank below that best describe the core techniques/data
1310 structures needed to solve it. Return JSON ONLY, no extra text,
1311 as:

1312 {"topics": ["Topic1", "Topic2"]}

1314
1315 Rules:

- 1316 - Use EXACT spellings from the bank (case/spacing must match).
- 1317 - Prefer the most specific tag available (e.g., "Binary Tree" over
1318 "Tree", "Shortest Path" over "Graph" when appropriate).
- 1319 - If the solution critically relies on a data structure (e.g.,
1320 "Heap (Priority Queue)"), include it.
- 1321 - If multiple techniques are essential (e.g., DP + Bitmask),
include both.
- 1322 - Do NOT exceed 3 topics; order them by importance.
- 1323 - If nothing fits, choose the closest general tag (e.g., "Graph",
1324 "Array", "Math")|never invent tags.

1325
1326 Topic Bank (allowed values only):

1327 Array; String; Hash Table; Dynamic Programming; Math; Sorting;
1328 Greedy; Depth-First Search; Binary Search; Database; Matrix; Tree;
1329 Breadth-First Search; Bit Manipulation; Two Pointers; Prefix Sum;
1330 Heap (Priority Queue); Simulation; Binary Tree; Graph; Stack;
1331 Counting; Sliding Window; Design; Enumeration; Backtracking; Union
1332 Find; Linked List; Number Theory; Ordered Set; Monotonic Stack;
1333 Segment Tree; Trie; Combinatorics; Bitmask; Divide and Conquer;
1334 Queue; Recursion; Geometry; Binary Indexed Tree; Memoization; Hash
1335 Function; Binary Search Tree; Shortest Path; String Matching;
1336 Topological Sort; Rolling Hash; Game Theory; Interactive;
1337 Data Stream; Monotonic Queue; Brainteaser; Doubly-Linked List;
1338 Randomized; Merge Sort; Counting Sort; Iterator; Concurrency;
1339 Probability and Statistics; Quickselect; Suffix Array; Line Sweep;
1340 Minimum Spanning Tree; Bucket Sort; Shell; Reservoir Sampling;
1341 Strongly Connected Component; Eulerian Circuit; Radix Sort;
1342 Rejection Sampling; Biconnected Component

1343
1344 Input:

1345 [Problem]
1346 {problem}

1347
1348 [Solution]
1349 {solution}

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

Output:
JSON with key "topics" and UP TO 3 strings from the bank. No
prose, no explanations.
"""