# PLUM: Improving Code LMs with Execution-Guided On-Policy Preference Learning Driven By Synthetic Test Cases

**Dylan Zhang**[1], **Shizhe Diao**[2], **Xueyan Zou**[3] and **Hao Peng**[1]

[1]University of Illinois Urbana-Champaign, [2]The Hong Kong University of Science and Technology, [3]University of California San Diego, [1]University of Illinois Urbana-Champaign

Preference learning provides a promising solution to address the limitations of supervised fine-tuning (SFT) for code language models, where the model is not explicitly trained to differentiate between correct and incorrect code. Recent findings demonstrate that on-policy data is the key to successful preference learning, where the preference data is collected using the same policy LM being trained. Inspired by this, we propose PLUM, an on-policy **P**reference **L**earning framework A**u**gmented with test cases for code L**M**s. The framework operates in three key stages: (1) automatic generation of test cases from natural language instructions, (2) creation of a preference data by evaluating candidate code solutions sampled from the policy, which can then be used to (3) train the policy LM. PLUM levitates the need to train reward models, allowing for large scale on-policy and online preference data collation. PLUM is evaluated on both standard benchmarks (HumanEval, MBPP) and more challenging ones (LiveCodeBench), delivering substantial improvements over original SFT'ed models and other execution-feedback-driven approaches. We show PLUM's benefits are consistent across various widely-used code LMs even they have been well-trained with SFT. For example, PLUM increases pass rates by up to 4.8% on average on standard benchmarks and 11.8% on LiveCodeBench, demonstrating its effectiveness and generalizability. We also demonstrate the benefits of on-policy and online preference learning by comprehensive experimentation.

## 1. Introduction

Language models pre-trained on code corpora have excelled at code generation [40, 25]. Supervised Fine-Tuning (SFT) enhances their ability to follow natural language prompts but focuses on reproducing patterns from training data rather than ensuring code correctness [49, 58]. This leads to models that generate syntactically correct but functionally flawed code, unable to meet real-world requirements like edge cases or algorithmic accuracy [6].

Works like AlphaCode [26] and LeTI [48] have introduced test outcomes as a means to define functional correctness in code generation. Building on the insights from these efforts, we propose leveraging preference learning for refining model behavior. Preference learning trains models to prefer certain solutions (e.g., factual, helpful, or harmless) over undesirable ones (e.g., inaccurate, unhelpful, or harmful). Despite its success in aligning models with human values and improving reasoning in other domains [10, 19, 54, 47, 36], the application of preference learning as a principled and efficient approach in code generation remains under-explored, largely due to the lack of high-quality training data.

Recent research shows that reducing the likelihood of incorrect outputs is more effective for improving model performance than simply maximizing correct responses [41, 44]. Mode-seeking objectives, which

prioritize minimizing errors, have been found to outperform maximum likelihood methods by more efficiently redistributing probability mass across potential outputs. This underscores the importance of applying on-policy and online approaches to enhance preference learning algorithms [44, 20, 41, 29]. Unlike offline preference data, on-policy data remains in-distribution with the model, reducing the risk of misalignment [20, 57, 45, 16]. The main challenge now is how to efficiently obtain preference labels for on-policy data at scale [53]. In programming tasks, test cases present as a native and powerful candidate solution to address this issue. Being able to automatically produce high-quality test cases unlocks the possibility of collecting preference data over programming questions at any scale.

---

**Prompt for Test Case Generation**

You are a teaching assistant helping to write reference solutions and tests for programming questions. Given a programming question, you need to first analyze the problem, then write a reference solution (code), followed by assertions that test student solutions. The test code must be runnable when concatenated at the end of student solutions to check the correctness.

**Programming Question:**
{Question}

Follow the format below:
[Analysis]
{{Natural language analysis of the problem.}}
[Solution]
{{Your solution to the problem}}
[Start Code]
{{Start code for students so that they can follow the I/O protocol.  E.g.
Function signatures, class names etc.}}
[Test Code]
{{Test code that is immediately runnable if concatenated with student code to
check the correctness.}}

---

To this end, we propose **p**reference **l**earning framework a**u**gmented with test cases for training code language **m**odels (**PLUM**), which integrates the process of deriving test cases from natural language specifications into the training process to obtain preference labels for model's on-policy candidate solutions.

PLUM utilizes natural language instructions from well-established datasets such as OSS-Instruct [49], Evol-Instruct-Code [31], and ShareGPT [8]. For each instruction, high-quality test cases are constructed, and multiple solutions are sampled from the model. These solutions are evaluated using the generated test cases, with preference labels assigned based on the results: solutions passing the tests are preferred, while failures are dis-preferred. This dataset then trains the policy using established preference learning algorithms [39, 14, 2]. By relying solely on policy model's self-generated solutions, PLUM eliminates the need for external, synthetic off-policy data, reducing the risk of distributional shifts and poor generalization commonly observed for synthetic off-policy data, enhancing the model's robustness and improving its ability to differentiate between correct and incorrect solutions. In addition, by showcasing the effectiveness of our framework, we demonstrated the feasibility of bypassing tedious (and potentially unstable) reward model training [27, 24] and manual labeling, by automating the process of test case synthesis. These simplicity advantages of PLUM makes online preference training of language models possible.

We evaluate PLUM on a diverse set of state-of-the-art code language models under different set-ups, on commonly used evaluation benchmarks: HumanEval(+) and MBPP(+) [7, 1, 28] as well as more challenging

code generation datasets like LiveCodeBench and LeetCode [22, 19]. We demonstrate that our approach seamlessly integrates with various models in a plug-and-play manner, relying solely on coding instructions to enhance models' code generation capabilities. Furthermore, we show that online training, facilitated by automated test case generation, further boosts model performance particularly on difficult coding benchmarks, echoing findings from other domains [51].

## 2. Preference Learning Augmented with Test Cases for Code LMs (PLUM)

The core of PLUM lies in leveraging recent advancements in on-policy and online preference learning, which have proven effective across various domains [51, 50, 33]. In the context of code generation, PLUM simplifies and scales the process of collecting preference data by using test cases. These test cases act as a lightweight yet robust mechanism for evaluating model outputs.

Algorithm 1 outlines the core mechanism of PLUM, demonstrating how test case generation is embedded into the preference learning loop. This process allows model-generated outputs to be evaluated in real-time using automatically generated test cases, which serve as direct feedback mechanisms for the learning process. By using test cases rather than complex reward models, PLUM simplifies the collection of preference data, maintaining high feedback quality while reducing the complexity associated with reward model training.

### 2.1. The PLUM

As illustrated in Algorithm 1 and Figure 1, our approach takes in a base policy model, a set of natural language programming instructions, and a test case generator. For each iteration, it will produce multiple test cases for the batch of instructions. Then we sample solutions from policy $\pi_\theta$, and execute them against the generated test suite to obtain preference labels. We then update the policy $\pi_\theta := \pi'_\theta$.

### 2.2. Generating Test Cases

A crucial factor in making PLUM successful is the ability to synthesize high-quality test cases for programming questions. In the following subsections, we provide a detailed explanation of the test case generation process, outlining how it contributes to the overall effectiveness of PLUM.

The test cases in PLUM are generated with a test-case generator model over natural instructions from established code generation datasets.[1] In automated testing, ensuring the correctness and completeness of test cases is a persistent challenge due to the lack of reliable oracles to validate test outputs. We adopt two strategic principles: 1) employing self-consistency as an approximate oracle, and 2) generating diverse test suites to minimize overfitting to any particular test instance and mitigate under-specification.

**Collecting instructions from established datasets**   We collect natural language instructions from established datasets including OSS-Instruct [49], Evol-Instruct-Code [31], and ShareGPT [8]. [2]These datasets provide a diverse range of programming tasks and instructions. Although they come with gold/silver solutions in the training splits, these solutions are *never* used in PLUM. Instead, they allow us to directly compare PLUM's performance against SFT, which relies on gold solutions, as we investigate in our experiments. Not

---

[1]We use GPT-4-1106 as the generator model.
[2]We focus on Python due to its wide use and the availability of well-established training and evaluation resources.

requiring gold solutions for training broadens the applicability of PLUM to a wide variety of real-world coding tasks and user requirements.

---

**Algorithm 1** PLUM.

---

**Input:** Natural language instructions $\mathcal{I} = \{q_i\}$, policy model to be trained $\pi_\theta$, generator model $G$, update frequency $T$, chunk size $M$            ▷ Unified for both offline and online alignment
**Output:** Trained policy model $\pi'_\theta$
 1: Initialize preference datasets $\mathcal{D}^+$ and $\mathcal{D}^-$
 2: **for** each chunk $\mathcal{I}_M \subset \mathcal{I}$, where $\mathcal{I}_M$ contains $M$ instructions **do**
 3:      **for** each $q_i \in \mathcal{I}_M$ **do**
 4:          Generate $n$ pairs of reference code and test case $\{(r_{ij}, t_{ij})\}_{j=1}^n$ using $G$      ▷ Test collection
 5:          **for** each pair $(r_{ij}, t_{ij})$ **do**
 6:              **if** $r_{ij}$ passes $t_{ij}$ **then**
 7:                  Add $(q_i, t_{ij})$ to $\mathcal{D}$      ▷ Self-consistency filtering
 8:              **end if**
 9:          **end for**
10:          $\mathcal{S}_{ik} \sim \pi_\theta$ for $k = 1$ to $K$ (sample $K$ solutions for $q_i$)      ▷ On-policy sampling
11:          **for** each solution $s_{ik} \in \mathcal{S}_{ik}$ **do**
12:              **for** each test case $t_{ij}$ in $\mathcal{D}$ **do**
13:                  **if** $s_{ik}$ passes $t_{ij}$ **then**
14:                      Add $(q_i, s_{ik})$ to $\mathcal{D}^+$      ▷ Positive case
15:                  **else**
16:                      Add $(q_i, s_{ik})$ to $\mathcal{D}^-$      ▷ Negative case
17:                  **end if**
18:              **end for**
19:          **end for**
20:      **end for**
21:      Filter out instances with no correct solutions from $\mathcal{D}^+$ and $\mathcal{D}^-$
22:      **if** iteration count $\%T = 0$ **then**      ▷ Policy Update
23:          Train the policy model $\pi_\theta$ using $\mathcal{D}^+$ and $\mathcal{D}^-$ with preference learning to get $\pi'_\theta$
24:          Update policy model $\pi_\theta = \pi'_\theta$
25:      **end if**
26: **end for**
27: **return** $\pi'_\theta$

---

**Generating high-quality test cases**      Given a training instruction in natural language, we prompt a generator model to produce a reference solution, a starter code snippet specifying the function signature, and a suite of test cases using the prompt in Figure 1. The generated test cases are critical for ensuring that the solutions meet the functional requirements specified in the instructions. The correctness of the test cases is central to the success of preference learning. We adopt a consistency-based approach inspired by [5] and [40] for quality control.

We check for consistency between the generated reference solution and the test cases. Pairs where the test cases do not accurately reflect the solution, or the solution does not pass the test cases, are filtered out. This process helps minimize potential noise and enhances the quality of the test cases used in the following stages. The generated reference solutions serve only to control the quality of the test cases and are *never* used in training. Similarly, the solutions provided with the instruction data are *never* used in PLUM. On average, each instruction is paired with 3–5 test cases depending on the dataset.

## 2.3. Sampling Solutions from the Policy to Create the Preference Data

Many preference learning algorithms assume that the preference data is in-distribution for the policy, i.e., the solutions are sampled from the policy model to be trained [39, 14, 2]. In practice, however, preference data often contains solutions sampled from different models than the policy, leaving the data out of distribution [4, 54]. A common workaround is to first perform supervised fine-tuning (SFT) on the same instructions before applying preference learning [39, 54]. This ensures that the policy has a similar distribution to that from which the preference data are sampled.

One of the research questions we aim to answer through PLUM is the standalone effect of preference learning on LMs' coding capability, with or without first performing SFT. To this end, we sample solutions from the policy to be trained and run them against the test cases to create the preference data. For each instruction, we sample $K$ solutions from the policy and evaluate them against the generated test cases. $K$ is set to 20 based on the findings from our preliminary experiments.

| Dataset | Self-Consistency Pass Rate(%) |
|---|---|
| OSS-INSTRUCT | 63.76 |
| EVOL-INSTRUCT | 42.38 |
| SHAREGPT | 45.69 |

Table 1: Self-Consistency Pass Rate Using GPT-4-1106.

With static and execution checks,[3] we identify and filter out solutions that contain syntactic errors and fail to execute, as our focus is on functional correctness.

Moreover, as a recent work points out, training with code snippets containing syntax errors may hurt the model's performance [48]. Solutions passing all test cases are used as the chosen solutions, and those failing at least one the rejected solutions.

An instruction is filtered out if it has no chosen solution after this process.

This aims to ensure that the learned policy does not drift too far from the original one as drastic changes might cause the model to forget previously learned information or to perform poorly on tasks it was previously adept at [39].

## 2.4. Preference Learning

We then proceed to train the model on the on-policy sampled candidate solutions using preference learning algorithm. In this process, we do not need golden solutions paired in the original dataset or GPT-4 during test-generation process. We mainly consider two popular preference learning algorithms - Direct Preference Optimization [39] and Kahneman-Tversky Optimization [14] that have been shown to bring improvements for reasoning tasks [33, 54, 12]. For DPO, we subsample redundant classes and randomly pair positive and negative responses for each programming question. In contrast, we use all available responses when training with KTO.

## 3. Experiments

To demonstrate the effectiveness of PLUM, we evaluate it on established benchmarks: HumanEval [7], MBPP [1], and EvalPlus (a widely-adopted augmented version [28] of them). We also use the more challenging LiveCodeBench [22].

---

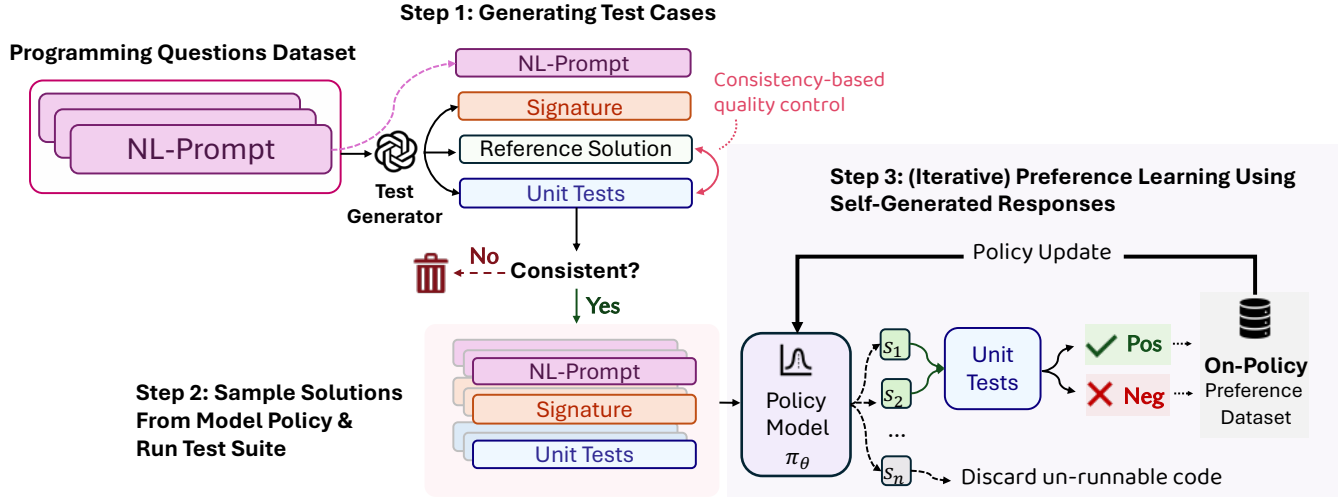[3]We use mypy for the static check: https://mypy-lang.org/.

**Figure** 1: Overview of PLUM. It involves three steps: (1) Generating the test cases; (2) Sampling solutions from the policy and evaluating them against the test cases to collect the preference data for (3) performing preference learning.
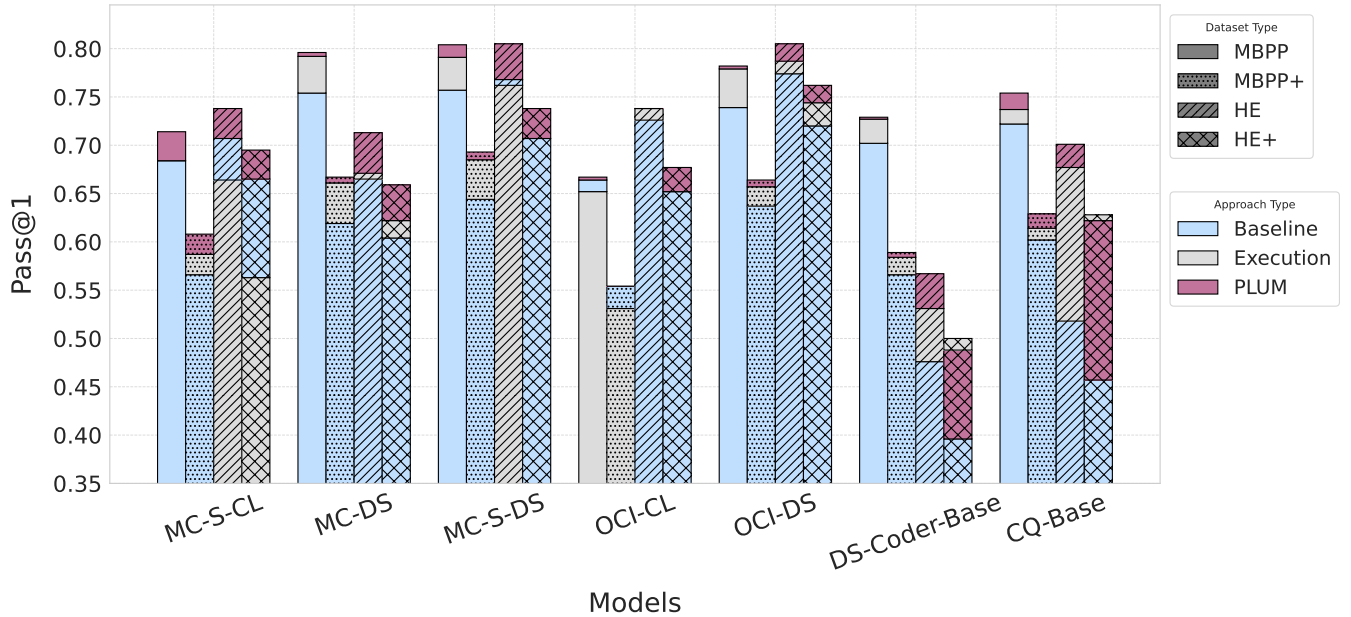


**Figure** 2: Ablations on preference training signals reveal that using un-runnable code as negative instances does not consistently lead to performance improvements. The baseline results refer to the SFT-model without PLUM.

**Datasets** To demonstrate the generality of our approach across different datasets, we evaluated it on three distinct collections of open datasets: OSS-Instruct (GPT-3.5 generated), EvolInstruct, and a Python code generation subset of ShareGPT. We showed that PLUM can significantly enhance model performance in various settings with high data efficiency, even when using only a small, randomly selected subset of SFT datasets.

**Preference Data Collection** We use GPT-4 [34] to generate test cases for each programming question. For the OSS-Instruct dataset and ShareGPT dataset, we query GPT-4 for 3 responses for a randomly chosen subset of 1500 questions, and due to the comparatively more complex nature of the natural language instruction, we generate 6 for each of the EvolInstruct instances for a subset of 1000.

We then sample 20 outputs from the policy using temperature $T = 1$ for the former two and 50 outputs for the latter. This yields around $\sim 60,000$ examples for ShareGPT and OSS-Instruct, and around $\sim 120,000$ for EvolInstruct before any filtering. We present the statistics on the pass ratio of sampled solutions over OSS-Instruct in Figure 3. We included the self-consistency pass rate of the test-generation process with GPT-4-1106 in Table 1.

**Models** We consider a diverse set of strong open language models: MagiCoder [49], OpenCodeInterpreter [58], CodeQwen [3], DeepSeek Coder [19] and StarCoder2 [25]. MagiCoder and OpenCodeIntepreter contain instruction-tuned checkpoints from DeepSeek Coder and CodeLlama [40] base models. In the main text, we focus on instruction-tuned language models, while differing results from directly training base models to Appendix A.5

**Baselines** We evaluate our approach against a variety of baselines, including both prompting-based and fine-tuning techniques. To compare methods that incorporate program correctness through execution feedback, we benchmark our approach against *Reflexion* [42], using instruction-tuned (SFT) models and value-conditioning techniques [26, 48]. Additionally, to contrast preference-learning techniques with the SFT approach, we perform experiments using rejection-sampling-based SFT, utilizing the same set of positive examples as used in KTO. In these experiments, the solutions are also generated on-policy.

### 3.1. Training

In order to demonstrate the generality of the approach when applied to various models and code instruction tuning data distributions, we experimented with different data-model pairs. We followed the procedure described earlier in the paper, and used all positive and negative responses when training with KTO objective.

| Model | Item | MBPP | MBPP+ | HE | HE+ | Avg. |
|-------|------|------|-------|-----|-----|------|
| CODEQWEN-1.5-CHAT [3] | Baseline | 77.7 | 67.2 | 83.5 | 78.7 | 76.8 |
| | Cond.Token | 71.9 | 57.1 | 68.3 | 57.9 | 63.8 |
| | RFT | 81.2 | 67.9 | 84.1 | 81.1 | 78.6 |
| | Cond.Err.Msg | 79.4 | 68.7 | 84.1 | 79.3 | 77.9 |
| | PLUM-DPO | 81.2 | 70.2 | 86.0 | 81.1 | **79.6** |
| | **PLUM-KTO** | **81.0** | **69.0** | **86.0** | **81.1** | **79.3** |
| | *Rel. +* | *4.3* | *2.7* | *3.0* | *3.1* | *3.3* |
| DS-CODER-INSTRUCT [19] | Baseline | 74.9 | 65.6 | 75.4 | 71.3 | 71.8 |
| | Cond.Token | 73.4 | 62.9 | 76.8 | 70.7 | 71.0 |
| | RFT | 74.7 | 64.9 | 74.4 | 66.5 | 70.1 |
| | Cond.Err.Msg | 74.7 | 64.2 | 80.5 | 75.0 | 73.6 |
| | PLUM-DPO | 76.4 | 65.9 | 80.5 | 76.8 | 77.4 |
| | **PLUM-KTO** | **78.2** | **67.9** | **81.7** | **76.8** | **76.2** |
| | *Rel. +* | *4.4* | *3.5* | *8.4* | *7.7* | *6.0* |
| MAGICODER-DS [49] | Baseline | 75.4 | 61.9 | 66.5 | 60.4 | 66.1 |
| | Cond.Token | 75.9 | 62.4 | 68.3 | 62.2 | 67.2 |
| | RFT | 76.2 | 62.2 | 67.7 | 62.2 | 67.1 |
| | Cond.Err.Msg | 74.2 | 62.2 | 66.5 | 59.8 | 65.7 |
| | PLUM-DPO | 75.9 | 63.7 | 67.7 | 61.6 | 67.2 |
| | **PLUM-KTO** | **79.6** | **66.7** | **71.3** | **65.9** | **70.9** |
| | *Rel. +* | *5.6* | *7.8* | *7.2* | *9.1* | *7.4* |
| MAGICODER-S-DS [49] | Baseline | 75.7 | 64.4 | 76.8 | 70.7 | 71.9 |
| | Cond.Token | 73.9 | 63.7 | 75.0 | 71.3 | 71.0 |
| | RFT | 75.4 | 64.4 | 73.2 | 69.5 | 70.6 |
| | Cond.Err.Msg | 75.2 | 65.4 | 75.6 | 70.7 | 71.7 |
| | PLUM-DPO | 76.2 | 64.7 | 78.7 | 73.8 | 73.4 |
| | **PLUM-KTO** | **80.4** | **69.3** | **80.5** | **73.8** | **76.0** |
| | *Rel. +* | *4.4* | *7.4* | *4.4* | *4.5* | *5.2* |
| OCI-DS [58] | Baseline | 73.9 | 63.7 | 77.4 | 72.0 | 71.8 |
| | Cond.Token | 73.9 | 62.9 | 75.6 | 71.3 | 70.9 |
| | RFT | 74.2 | 63.7 | 76.8 | 72.0 | 71.7 |
| | Cond.Err.Msg | 75.0 | 70.7 | 74.4 | 64.7 | 71.2 |
| | PLUM-DPO | 76.4 | 66.4 | 80.5 | 76.2 | 74.9 |
| | **PLUM-KTO** | **78.2** | **66.4** | **80.5** | **76.2** | **75.3** |
| | *Rel. +* | *5.8* | *4.2* | *4.0* | *5.8* | *5.0* |
| OCI-CL [58] | Baseline | 66.4 | 55.4 | 72.6 | 65.2 | 64.9 |
| | Cond.Token | 59.6 | 48.4 | 23.2 | 21.3 | 38.1 |
| | RFT | 63.2 | 52.6 | 60.4 | 56.7 | 58.2 |
| | Cond.Err.Msg | 67.9 | 55.9 | 68.9 | 65.2 | 64.5 |
| | PLUM-DPO | 66.4 | 55.9 | 71.3 | 65.2 | 64.7 |
| | **PLUM-KTO** | **66.7** | **55.4** | **73.8** | **67.7** | **65.9** |
| | *Rel. +* | *0.5* | *0.0* | *1.7* | *3.8* | *1.5* |

Table 2: %Pass@1 on HumanEval (HE) and MBPP, and their enhanced versions (HE+ and MBPP+) when PLUM is applied to OSS-Instruct. The *Rel. +* is computed as the relative percentage increase of PLUM-KTO over baseline. PLUM brings consistent improvements over SFT-ed baseline and outperforms other methods that leverage execution feedback when applied to the same SFT-ed models.

## 3.2. Results

**HumanEval(+) and MBPP(+)**   Table 2 presents the results of PLUM when applied to a subset of 1K instances of the OSS-Instruct-75K dataset.

MagiCoder models (-DS, -S-CL, and -S-DS) and OpenCodeIntepreter models (-CL and -DS) have already seen these instructions during supervised fine-tuning, while DeepSeekCoder-Instruct has not, as it was released earlier than the dataset. CodeQwen chat model uses proprietary data. PLUM-ShareGPT data for preference learning is generated with the same setting. Similarly, Table 4 corresponds to the results when we apply PLUM to EvolInstruct [31] dataset. Since the instructions are comparatively less clear than the OSS-Instruct dataset, we control the number of initial samples to be the same by generating 50 samples for each problem and use about 400 instances in total.

| Model | Item | MBPP/(+) | HE/(+) |
|---|---|---|---|
| DS-CODER -INSTRUCT | Base | 75/66 | 75/71 |
| | Reflexion | 34/- | 43/- |
| | **PLUM** | **78/68** | **82/77** |
| CODEQWEN -7B-CHAT | Base | 78/67 | 84/79 |
| | Reflexion | 74/- | 83/- |
| | **PLUM** | **81/69** | **86/81** |

Table 3: Comparison with Reflexion [42]

PLUM consistently improves the performance of a wide range of code language models across all three settings, regardless of the base models' performance. Remarkably, PLUM can even improve the state-of-the-art 7B model, CodeQwen-7B-Chat, relatively by 3% on average, using either OSS-Instruct or ShareGPT data. These results demonstrate that PLUM is broadly applicable in different datasets and settings.

| Model Families | Data | Item | MBPP | MBPP+ | HE | HE+ | Avg. |
|---|---|---|---|---|---|---|---|
| WIZARDCODER [31] | EVOLINSTRUCT | Baseline | 48.2 | 40.9 | 56.6 | 47.1 | 48.2 |
| | | **PLUM-KTO** | **54.3** | **48.8** | **65.9** | **52.9** | **55.5** |
| | | *Rel. +* | *16.4* | *12.3* | *12.7* | *19.3* | *15.2* |
| DS-CODER-INSTRUCT | | Baseline | 74.9 | 65.6 | 75.4 | 71.3 | 71.8 |
| | | **PLUM-KTO** | **77.7** | **67.7** | **81.7** | **76.8** | **76.0** |
| | | *Rel. +* | *3.7* | *3.2* | *8.4* | *7.7* | *5.8* |
| DS-CODER-INSTRUCT | SHAREGPT -PYTHON | Baseline | 74.9 | 65.6 | 75.4 | 71.3 | 71.8 |
| | | **PLUM-KTO** | **79.2** | **67.9** | **77.4** | **73.8** | **74.6** |
| | | *Rel. +* | *5.7* | *3.5* | *2.8* | *3.5* | *3.9* |
| OCI-DS | | Baseline | 73.9 | 63.7 | 77.4 | 72.0 | 71.8 |
| | | **PLUM-KTO** | **77.7** | **64.4** | **79.9** | **75.6** | **74.4** |
| | | *Rel. +* | *5.1* | *1.1* | *3.2* | *5.0* | *3.6* |
| CODEQWEN-1.5-CHAT | | Baseline | 77.7 | 67.2 | 83.5 | 78.7 | 76.8 |
| | | **PLUM-KTO** | **81.2** | **69.7** | **85.4** | **79.3** | **78.9** |
| | | *Rel. +* | *4.5* | *3.7* | *2.3* | *0.8* | *2.8* |

Table 4: PLUM on other datasets.

We noticed that PLUM-KTO consistently out-performs the baseline, and that PLUM-DPO sometimes under-perform PLUM-KTO. Prior works [33, 54] noticed the phenomenon where DPO can exhibit instability due to reducing reward for the positive class.

**LiveCodeBench**   We further evaluate PLUM using strong instruction-tuned models on the more challenging LiveCodeBench dataset. As shown in Table 5, the models demonstrate overall performance improvements over their respective baselines across the board. Despite the increased difficulty and reasoning required, we show that PLUM can enhance the base models' overall coding performance on interview-level coding

problems from LiveCodeBench.

PLUM proves particularly beneficial for medium-level interview questions, which are often quite challenging for models, especially those with around 7B parameters. This demonstrates that PLUM does more than simply fitting to commonly tested benchmarks; it enhances the models' general coding capabilities in more complex and diverse coding scenarios.

| Model | Item | Easy | Medium | Hard | Overall |
|---|---|---|---|---|---|
| | Baseline | 29.9 | 1.0 | 0 | 11.4 |
| MAGICODER-S-CL | **PLUM-KTO** | **38.1** | **1.8** | **0** | **14.3** |
| | *Rel. +* | *27.4* | *80* | *0* | *25.4* |
| | Baseline | 35.2 | 3.6 | 0.0 | 14.0 |
| MAGICODER-DS | **PLUM-KTO** | **55.6** | **13.1** | **2.2** | **25.8** |
| | *Rel. +* | *58.0* | *266.7* | *-* | *83.9* |
| | Baseline | 48.6 | 12.1 | 0.1 | 22.6 |
| MAGICODER-S-DS | **PLUM-KTO** | **52.1** | **15.5** | **0.1** | **25.0** |
| | *Rel. +* | *7.2* | *28.1* | *0.1* | *10.6* |
| | Baseline | 49.6 | 9.9 | 0.4 | 21.9 |
| OCI-DS | **PLUM-KTO** | **45.8** | **13.7** | **1.2** | **22.3** |
| | *Rel. +* | *-7.7* | *38.4* | *200* | *1.8* |
| | Baseline | 42.7 | 18.8 | 0.9 | 23.2 |
| CODEQWEN-1.5-CHAT | **PLUM-KTO** | **43** | **23.2** | **3** | **25.8** |
| | *Rel. +* | *0.7* | *20* | *230* | *11.2* |

Table 5: %Pass@1 on LiveCodeBench.

**Comparison Against Baselines**   As shown in Table 3, verbal reinforcement learning like Reflexion [42], does not perform well on code language models fine-tuned for code at the scale relevant to our work. This is partly due to the limitations of smaller LLMs in handling various types of instructions effectively.

Approaches like LeTI [48] implicitly optimizes the model for generating correct programs solely through input prompt, without directly enforcing such distinction with its training objective. Unlike preference learning algorithms ,approaches like LeTI [48] optimize models to generate correct programs based solely on the input prompt, without explicitly enforcing correctness through the training objective. As a result, we observe inconsistent outcomes when applying these methods to our tested SFT models, as shown in Table 2. Additionally, our results show that using value-token-conditioned approaches often leads to reduced performance, likely due to the introduction of new value tokens and the absence of clear labels distinguishing good from bad outputs.

**On-Policy vs Off-Policy**   We study how important on-policy training is to the performance. The result is presented in 6. To test this, we apply the same method but use preference data sampled from other models (i.e., off-policy). We selected models of varying sizes, including DeepSeek-Coder-1.3B-Instruct, Qwen-2.5-Coder-1.5B-Instruct, CodeStral-22B [32], and DeepSeek-Coder-33B-Instruct. We followed the exact same data collection and training processes. Our results show that while off-policy training still improves model performance, it generally underperforms compared to on-policy training. These findings are consistent with prior research  [50, 11, 44, 52].

Further, we investigate the effect of synthetic negatives. To this effect, we use a mutation-based approach

| | On-Policy | Off-Policy | | | | |
|---|---|---|---|---|---|---|
| | **PLUM** | **DS-Coder -1.3B -Instruct** | **Qwen2.5- Coder-1.5B -Instruct** | **Codestral -22B** | **DeepSeek -Coder-33B -Instruct** | **Syn. -Neg.** |
| CodeQwen1.5-Chat | **79.3** | 78.4 | 78.1 | 77.3 | 79.0 | 73.9 |
| DS-Coder-Instruct | **76.2** | 76.0 | 75.2 | 74.6 | 74.0 | 69.4 |
| MagiCoder-DS | **70.9** | 68.4 | 68.6 | 68.3 | 67.0 | 61.2 |
| Magicoder-S-DS | **76.0** | 75.8 | 75.3 | 74.1 | 73.8 | 71.6 |

Table 6: Comparison between preference learning with on-/off-policy data using KTO. The numbers are Macro Avg. across HE, MBPP's **base** and **(+)** tests. **Syn.Neg.** is the result by using synthetic negative data.

for synthetically introducing errors into Python code while maintaining its **syntactic correctness**, as detailed in Appendix A.4. This method uses Abstract Syntax Tree (AST) manipulation to apply mutations like argument swapping, operator replacement, control flow changes, off-by-one errors, and return value modification. By injecting these errors, it generates valid but behaviorally altered code. We apply this to on-policy positive examples, creating off-policy negatives for model training under the same setup. Observing that the synthetic negatives could even potentially harm the performance, we confirmed the importance of preference training with more natural, and ideally on-policy negative samples.

This highlights our contribution in demonstrating a method empowered by automated test cases and efficient on-policy preference learning. This approach can be easily adopted to scale the collection of test cases, providing a robust supervision signal for model training.

**Importance of Test Case-Based Preference Learning**   We experiment with including only non-executable samples as rejected solutions, while using the same set of chosen solutions. As shown in Figure 2, we observe that this is worse than PLUM in most cases. More importantly, it does not always improve the model's performance and may even hurt. This has also been noted in previous studies [48]. Although the positive examples used are the same as our oracle-based preference learning, lower-quality negative examples do not necessarily help the model improve due to the additional noise in the preference signal.

| | Easy | Medium | Hard | Average | Rel. Gain |
|---|---|---|---|---|---|
| CodeQwen1.5-Chat | 66.7 | 27.5 | 13.6 | 33.9 | - |
| +PLUM-DPO | 66.7 | 31.9 | 15.9 | 36.7 | 8.3 |
| **+PLUM-DPO-Iter** | **66.7** | **31.9** | **18.2** | **37.2** | **9.7** |
| +PLUM-KTO | 68.9 | 30.8 | 11.4 | 35.6 | 5.0 |
| **+PLUM-KTO-Iter** | **66.7** | **31.9** | **18.2** | **37.2** | **9.7** |
| Magicoder-DS | 33.3 | 17.6 | 9.1 | 19.4 | - |
| +PLUM-DPO | 44.4 | 15.4 | 13.6 | 22.2 | 14.4 |
| **+PLUM-DPO-Iter** | **46.7** | **17.6** | **11.4** | **23.3** | **20.1** |
| +PLUM-KTO | 48.9 | 16.5 | 9.1 | 22.8 | 17.5 |
| **+PLUM-KTO-Iter** | **44.4** | **17.6** | **11.4** | **22.8** | **17.5** |

Table 7: Results on iterative PLUM following Algorithm 1.

**Test Cases Allow By-passing Reward Model Training For Iterative Alignment**   We demonstrate in Table 7 that PLUM enables iterative on-policy alignment while providing accurate preference signals without

requiring reward model training. Notably, iterative preference learning, facilitated by the online feedback loops generated through the test case collection procedure, outperforms offline methods on the challenging LeetCode benchmark, using the same data. This finding aligns with results from other domains [50, 51], further confirming the advantages of online learning. This demonstrates the further potential of PLUM in advancing code language models especially in more challenging problems by its support for efficient online policy improvement.

## 4. Related Works

**Reinforcement Learning and Preference Learning For Reasoning-Related Tasks**   Preference learning algorithms like Direct Preference Optimization (DPO) [39] and Kahneman & Tversky's Optimization(KTO) [14] are popular for their cost efficiency and training stability. Beyond controlling model-user interactions, these methods are now applied to more complex reasoning tasks. Ocra-Math [33] uses iterative preference learning to improve math reasoning in SFT-ed language models, while Eurus [54] leverages preference trees for solving complex problems through multi-step interactions with external feedback.

**Code Generation with Large Language Models**   Code generation has become a key application of generative language models.  Pre-training on code corpora has led to strong performance in models like StarCoder [25], StarCoder2 [30], and DeepSeek-Coder [19], while others like CodeQwen [3] and CodeLlama [40] benefit from continued pre-training on additional code data. To enhance these models, supervised fine-tuning on instruction-response pairs has been employed [31, 49, 58].  Reinforcement learning techniques, like those in CodeRL [24, 43, 27], and reward models used in DeepSeek-Coder-V2 [9] also improve performance using test feedback.

**Test Case Generation with Language Models**   Automated test case generation [35, 17, 37] is crucial for ensuring software quality and safety and has long been a key topic in software engineering. The advent of LLMs has inspired works using transformer models for test generation, either by training models [46, 26] or prompting them [5]. Test cases also help clarify user intent, aligning model-generated programs with user requirements [15, 13].

**The synergy between test cases and code generation**   Programming-by-examples [18] and test-driven programming [38] focus on using test cases to automatically refine programs to meet specifications. This concept has been adapted to enhance deep learning approaches to code synthesis [23, 5, 55].  Recent methods, like CodeT [5] and Parsel [55], use test cases to reduce the search space during inference, while ALGO [56] employs brute-force solutions as oracles to generate test outputs for competitive programming. Our approach, similar to [21], leverages test cases during training to improve models' inherent programming capabilities.

## 5. Conclusion

In this paper, we introduced PLUM, a novel preference learning framework designed to improve the ability of code language models (LMs) to distinguish between correct and incorrect code by leveraging test cases. Our framework tackles the limitations of traditional supervised fine-tuning approaches by embedding on-policy learning directly into the training process. Through the automatic generation and evaluation of test cases,

PLUM enables models to learn from their own outputs without requiring separate reward models or manual labeling, offering a scalable and flexible solution.

The results from our experiments demonstrate the effectiveness and generalizability of PLUM. Furthermore, we performed careful experiments and showed that on-policy preference learning outperforms various off-policy methods, highlighting the crucial role played by on-policy training. Further, we demonstrated PLUM allows for effective online preference learning that further pushes the performance on challenging coding benchmarks.

# References

[1] Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, and Charles Sutton. Program synthesis with large language models, 2021.

[2] Mohammad Gheshlaghi Azar, Mark Rowland, Bilal Piot, Daniel Guo, Daniele Calandriello, Michal Valko, and Rémi Munos. A general theoretical paradigm to understand learning from human preferences, 2023.

[3] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. Qwen technical report, 2023.

[4] Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom Brown, Jack Clark, Sam McCandlish, Chris Olah, Ben Mann, and Jared Kaplan. Training a helpful and harmless assistant with reinforcement learning from human feedback, 2022.

[5] 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.

[6] Jie Chen, Xintian Han, Yu Ma, Xun Zhou, and Liang Xiang. Unlock the correlation between supervised fine-tuning and reinforcement learning in training code large language models, 2024. URL https://arxiv.org/abs/2406.10305.

[7] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin,

Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large language models trained on code, 2021.

[8] Dome Cleston. Sharegpt. https://github.com/domeccleston/sharegpt, 2023.

[9] DeepSeek-AI, Qihao Zhu, Daya Guo, Zhihong Shao, Dejian Yang, Peiyi Wang, Runxin Xu, Y. Wu, Yukun Li, Huazuo Gao, Shirong Ma, Wangding Zeng, Xiao Bi, Zihui Gu, Hanwei Xu, Damai Dai, Kai Dong, Liyue Zhang, Yishi Piao, Zhibin Gou, Zhenda Xie, Zhewen Hao, Bingxuan Wang, Junxiao Song, Deli Chen, Xin Xie, Kang Guan, Yuxiang You, Aixin Liu, Qiushi Du, Wenjun Gao, Xuan Lu, Qinyu Chen, Yaohui Wang, Chengqi Deng, Jiashi Li, Chenggang Zhao, Chong Ruan, Fuli Luo, and Wenfeng Liang. Deepseek-coder-v2: Breaking the barrier of closed-source models in code intelligence, 2024. URL https://arxiv.org/abs/2406.11931.

[10] Hanze Dong, Wei Xiong, Deepanshu Goyal, Yihan Zhang, Winnie Chow, Rui Pan, Shizhe Diao, Jipeng Zhang, Kashun Shum, and Tong Zhang. Raft: Reward ranked finetuning for generative foundation model alignment. *arXiv preprint arXiv:2304.06767*, 2023.

[11] Hanze Dong, Wei Xiong, Bo Pang, Haoxiang Wang, Han Zhao, Yingbo Zhou, Nan Jiang, Doyen Sahoo, Caiming Xiong, and Tong Zhang. Rlhf workflow: From reward modeling to online rlhf, 2024. URL https://arxiv.org/abs/2405.07863.

[12] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim,

Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aaron Grattafiori, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kun Huang, Kunal Chawla, Kushal Lakhotia, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong, Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve

Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vítor Albiero, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaofang Wang, Xiaojian Wu, Xiaolan Wang, Xide Xia, Xilun Wu, Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. The llama 3 herd of models, 2024. URL https://arxiv.org/abs/2407.21783.

[13] Madeline Endres, Sarah Fakhoury, Saikat Chakraborty, and Shuvendu K. Lahiri. Can large language models transform natural language intent into formal method postconditions?, 2024.

[14] Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. Kto: Model alignment as prospect theoretic optimization, 2024.

[15] Sarah Fakhoury, Aaditya Naik, Georgios Sakkas, Saikat Chakraborty, and Shuvendu K Lahiri. Llm-based test-driven interactive code generation: User study and empirical evaluation. *arXiv preprint arXiv:2404.10100*, 2024.

[16] Adam Fisch, Jacob Eisenstein, Vicky Zayats, Alekh Agarwal, Ahmad Beirami, Chirag Nagpal, Pete Shaw, and Jonathan Berant. Robust preference optimization through reward model distillation. *arXiv preprint arXiv:2405.19316*, 2024.

[17] Gordon Fraser and Andrea Arcuri. EvoSuite: automatic test suite generation for object-oriented software. In *Proceedings of the 19th ACM SIGSOFT symposium and the 13th European conference on Foundations of software engineering*, pages 416–419, 2011.

[18] Sumit Gulwani. Programming by examples - and its applications in data wrangling. In *Dependable Software Systems Engineering*, 2016. URL https://api.semanticscholar.org/CorpusID:7866845.

[19] Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, Y. Wu, Y. K. Li, Fuli Luo, Yingfei Xiong, and Wenfeng Liang. Deepseek-coder: When the large language model meets programming – the rise of code intelligence, 2024.

[20] Shangmin Guo, Biao Zhang, Tianlin Liu, Tianqi Liu, Misha Khalman, Felipe Llinares, Alexandre Rame, Thomas Mesnard, Yao Zhao, Bilal Piot, Johan Ferret, and Mathieu Blondel. Direct language model alignment from online ai feedback, 2024. URL https://arxiv.org/abs/2402.04792.

[21] Patrick Haluptzok, Matthew Bowers, and Adam Tauman Kalai. Language models can teach themselves to program better. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=SaRj2ka1XZ3.

[22] 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, 2024.

[23] Sumith Kulal, Panupong Pasupat, Kartik Chandra, Mina Lee, Oded Padon, Alex Aiken, and Percy S Liang. Spoc: Search-based pseudocode to code. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran

Associates, Inc., 2019. URL https://proceedings.neurips.cc/paper_files/paper/2019/file/7298332f04ac004a0ca44cc69ecf6f6b-Paper.pdf.

[24] Hung Le, Yue Wang, Akhilesh Deepak Gotmare, Silvio Savarese, and Steven Chu Hong Hoi. Coderl: Mastering code generation through pretrained models and deep reinforcement learning. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages 21314–21328. Curran Associates, Inc., 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/8636419dea1aa9fbd25fc4248e702da4-Paper-Conference.pdf.

[25] Raymond Li, Loubna Ben allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou, Marc Marone, Christopher Akiki, Jia LI, Jenny Chim, Qian Liu, Evgenii Zheltonozhskii, Terry Yue Zhuo, Thomas Wang, Olivier Dehaene, Joel Lamy-Poirier, Joao Monteiro, Nicolas Gontier, Ming-Ho Yee, Logesh Kumar Umapathi, Jian Zhu, Ben Lipkin, Muhtasham Oblokulov, Zhiruo Wang, Rudra Murthy, Jason T Stillerman, Siva Sankalp Patel, Dmitry Abulkhanov, Marco Zocca, Manan Dey, Zhihan Zhang, Urvashi Bhattacharyya, Wenhao Yu, Sasha Luccioni, Paulo Villegas, Fedor Zhdanov, Tony Lee, Nadav Timor, Jennifer Ding, Claire S Schlesinger, Hailey Schoelkopf, Jan Ebert, Tri Dao, Mayank Mishra, Alex Gu, Carolyn Jane Anderson, Brendan Dolan-Gavitt, Danish Contractor, Siva Reddy, Daniel Fried, Dzmitry Bahdanau, Yacine Jernite, Carlos Muñoz Ferrandis, Sean Hughes, Thomas Wolf, Arjun Guha, Leandro Von Werra, and Harm de Vries. Starcoder: may the source be with you! *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL https://openreview.net/forum?id=KoFOg41haE. Reproducibility Certification.

[26] Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, Thomas Hubert, Peter Choy, Cyprien de Masson d'Autume, Igor Babuschkin, Xinyun Chen, Po-Sen Huang, Johannes Welbl, Sven Gowal, Alexey Cherepanov, James Molloy, Daniel J. Mankowitz, Esme Sutherland Robson, Pushmeet Kohli, Nando de Freitas, Koray Kavukcuoglu, and Oriol Vinyals. Competition-level code generation with alphacode. *Science*, 378(6624):1092–1097, December 2022. ISSN 1095-9203. doi: 10.1126/science.abq1158. URL http://dx.doi.org/10.1126/science.abq1158.

[27] Jiate Liu, Yiqin Zhu, Kaiwen Xiao, QIANG FU, Xiao Han, Yang Wei, and Deheng Ye. RLTF: Reinforcement learning from unit test feedback. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL https://openreview.net/forum?id=hjYmsV6nXZ.

[28] Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. Is your code generated by chatgpt really correct? rigorous evaluation of large language models for code generation, 2023.

[29] Zhihan Liu, Miao Lu, Shenao Zhang, Boyi Liu, Hongyi Guo, Yingxiang Yang, Jose Blanchet, and Zhaoran Wang. Provably mitigating overoptimization in rlhf: Your sft loss is implicitly an adversarial regularizer, 2024. URL https://arxiv.org/abs/2405.16436.

[30] Anton Lozhkov, Raymond Li, Loubna Ben Allal, Federico Cassano, Joel Lamy-Poirier, Nouamane Tazi, Ao Tang, Dmytro Pykhtar, Jiawei Liu, Yuxiang Wei, Tianyang Liu, Max Tian, Denis Kocetkov, Arthur Zucker, Younes Belkada, Zijian Wang, Qian Liu, Dmitry Abulkhanov, Indraneil Paul, Zhuang Li, Wen-Ding Li, Megan Risdal, Jia Li, Jian Zhu, Terry Yue Zhuo, Evgenii Zheltonozhskii, Nii Osae Osae Dade, Wenhao Yu, Lucas Krauß, Naman Jain, Yixuan Su, Xuanli He, Manan Dey, Edoardo Abati, Yekun Chai, Niklas Muennighoff, Xiangru Tang, Muhtasham Oblokulov, Christopher Akiki, Marc Marone, Chenghao Mou, Mayank Mishra, Alex Gu, Binyuan Hui, Tri Dao, Armel Zebaze, Olivier Dehaene,

Nicolas Patry, Canwen Xu, Julian McAuley, Han Hu, Torsten Scholak, Sebastien Paquet, Jennifer Robinson, Carolyn Jane Anderson, Nicolas Chapados, Mostofa Patwary, Nima Tajbakhsh, Yacine Jernite, Carlos Muñoz Ferrandis, Lingming Zhang, Sean Hughes, Thomas Wolf, Arjun Guha, Leandro von Werra, and Harm de Vries. Starcoder 2 and the stack v2: The next generation, 2024.

[31] Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xiubo Geng, Wenxiang Hu, Chongyang Tao, Jing Ma, Qingwei Lin, and Daxin Jiang. Wizardcoder: Empowering code large language models with evol-instruct, 2023.

[32] MistralAI. Codestral-22b-v0.1, 2024. URL https://huggingface.co/mistralai/Codestral-22B-v0.1. Accessed: 2024-09-28.

[33] Arindam Mitra, Hamed Khanpour, Corby Rosset, and Ahmed Awadallah. Orca-math: Unlocking the potential of slms in grade school math, 2024.

[34] OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantini-dis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O'Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Shep-pard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie

Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. Gpt-4 technical report, 2024.

[35] Carlos Pacheco, Shuvendu K Lahiri, Michael D Ernst, and Thomas Ball. Feedback-directed random test generation. In *29th International Conference on Software Engineering (ICSE'07)*, pages 75–84. IEEE, 2007.

[36] Richard Yuanzhe Pang, Weizhe Yuan, Kyunghyun Cho, He He, Sainbayar Sukhbaatar, and Jason Weston. Iterative reasoning preference optimization, 2024. URL https://arxiv.org/abs/2404.19733.

[37] Annibale Panichella, Fitsum Meshesha Kifetew, and Paolo Tonella. Reformulating branch coverage as a many-objective optimization problem. In *2015 IEEE 8th international conference on software testing, verification and validation (ICST)*, pages 1–10. IEEE, 2015.

[38] Daniel Perelman, Sumit Gulwani, Dan Grossman, and Peter Provost. Test-driven synthesis. In *PLDI '14, June 09-11, 2014, Edinburgh, United Kingdom*, June 2014. URL https://www.microsoft.com/en-us/research/publication/test-driven-synthesis/.

[39] Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model, 2023.

[40] Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, Ivan Evtimov, Joanna Bitton, Manish Bhatt, Cristian Canton Ferrer, Aaron Grattafiori, Wenhan Xiong, Alexandre Défossez, Jade Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier, Thomas Scialom, and Gabriel Synnaeve. Code llama: Open foundation models for code, 2024.

[41] Amrith Setlur, Saurabh Garg, Xinyang Geng, Naman Garg, Virginia Smith, and Aviral Kumar. Rl on incorrect synthetic data scales the efficiency of llm math reasoning by eight-fold, 2024. URL https://arxiv.org/abs/2406.14532.

[42] Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning, 2023. URL https://arxiv.org/abs/2303.11366.

[43] Parshin Shojaee, Aneesh Jain, Sindhu Tipirneni, and Chandan K. Reddy. Execution-based code generation using deep reinforcement learning. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL https://openreview.net/forum?id=0XBuaxqEcG.

[44] Fahim Tajwar, Anikait Singh, Archit Sharma, Rafael Rafailov, Jeff Schneider, Tengyang Xie, Stefano Ermon, Chelsea Finn, and Aviral Kumar. Preference fine-tuning of llms should leverage suboptimal, on-policy data. *arXiv preprint arXiv:2404.14367*, 2024.

[45] Yunhao Tang, Daniel Zhaohan Guo, Zeyu Zheng, Daniele Calandriello, Yuan Cao, Eugene Tarassov, Rémi Munos, Bernardo Ávila Pires, Michal Valko, Yong Cheng, and Will Dabney. Understanding the performance gap between online and offline alignment algorithms, 2024. URL https://arxiv.org/abs/2405.08448.

[46] Michele Tufano, Dawn Drain, Alexey Svyatkovskiy, Shao Kun Deng, and Neel Sundaresan. Unit test case generation with transformers and focal context, 2021.

[47] Peiyi Wang, Lei Li, Zhihong Shao, R. X. Xu, Damai Dai, Yifei Li, Deli Chen, Y. Wu, and Zhifang Sui. Math-shepherd: Verify and reinforce llms step-by-step without human annotations, 2024.

[48] Xingyao Wang, Hao Peng, Reyhaneh Jabbarvand, and Heng Ji. Leti: Learning to generate from textual interactions, 2024.

[49] Yuxiang Wei, Zhe Wang, Jiawei Liu, Yifeng Ding, and Lingming Zhang. Magicoder: Source code is all you need, 2023.

[50] Wei Xiong, Hanze Dong, Chenlu Ye, Ziqi Wang, Han Zhong, Heng Ji, Nan Jiang, and Tong Zhang. Iterative preference learning from human feedback: Bridging theory and practice for RLHF under KL-constraint. In *Forty-first International Conference on Machine Learning*, 2024. URL https://openreview.net/forum?id=c1AKcA6ry1.

[51] Wei Xiong, Chengshuai Shi, Jiaming Shen, Aviv Rosenberg, Zhen Qin, Daniele Calandriello, Misha Khalman, Rishabh Joshi, Bilal Piot, Mohammad Saleh, Chi Jin, Tong Zhang, and Tianqi Liu. Building math agents with multi-turn iterative preference learning, 2024. URL https://arxiv.org/abs/2409.02392.

[52] Wenda Xu, Jiachen Li, William Yang Wang, and Lei Li. Bpo: Supercharging online preference learning by adhering to the proximity of behavior llm. *arXiv preprint arXiv:2406.12168*, 2024.

[53] Sen Yang, Leyang Cui, Deng Cai, Xinting Huang, Shuming Shi, and Wai Lam. Not all preference pairs are created equal: A recipe for annotation-efficient iterative preference learning, 2024. URL https://arxiv.org/abs/2406.17312.

[54] Lifan Yuan, Ganqu Cui, Hanbin Wang, Ning Ding, Xingyao Wang, Jia Deng, Boji Shan, Huimin Chen, Ruobing Xie, Yankai Lin, Zhenghao Liu, Bowen Zhou, Hao Peng, Zhiyuan Liu, and Maosong Sun. Advancing llm reasoning generalists with preference trees, 2024.

[55] Eric Zelikman, Qian Huang, Gabriel Poesia, Noah D. Goodman, and Nick Haber. Parsel: Algorithmic reasoning with language models by composing decompositions, 2023.

[56] Kexun Zhang, Danqing Wang, Jingtao Xia, William Yang Wang, and Lei Li. Algo: Synthesizing algorithmic programs with llm-generated oracle verifiers, 2023.

[57] Shenao Zhang, Donghan Yu, Hiteshi Sharma, Ziyi Yang, Shuohang Wang, Hany Hassan, and Zhaoran Wang. Self-exploring language models: Active preference elicitation for online alignment. *arXiv preprint arXiv:2405.19332*, 2024.

[58] Tianyu Zheng, Ge Zhang, Tianhao Shen, Xueling Liu, Bill Yuchen Lin, Jie Fu, Wenhu Chen, and Xiang Yue. Opencodeinterpreter: Integrating code generation with execution and refinement, 2024.

# A. Appendix

## A.1. Test case generation

To produce test cases for each programming instruction, we queried OpenAI GPT-4 with temperature 0 and max response token 4096.

## A.2. Training Details

We trained the model using the KTO objective, with a learning rate $5 \times 10^{-7}$, linear scheduler, $\beta = 0.1$, and maintained the desirable-to-undesirable ratio to be 1. We train each model using 8-bit quantized LoRA for 3 epochs on a Nvidia A-100 GPU with 40 GB memory.

## A.3. Additional Results: What if directly using competitive coding datasets?

There are existing competitive coding datasets like CODE CONTESTS [26] which are equipped with test cases. Using more general datasets like OSS-INSTRUCT [49], paired with synthetic test cases, is more advantageous for preference learning in code language models than using competitive coding datasets like Code Contests. Competitive coding datasets present complex problems with intricate edge cases, which can overwhelm the model and obscure fundamental instruction-following and preference-learning goals. In contrast, OSS-Instruct provides more accessible, more uniform instructions that allow for cleaner and more straightforward alignment. This helps models learn functional correctness more effectively without being distracted by the nuances of competitive coding. Additionally, OSS-Instruct, sourced from real-world open-source projects, avoids domain-specific biases that can arise from competitive coding, making it more generalizable and applicable across diverse programming environments.

We also conducted experiments to repeat the same process using CODECONTESTS dataset. Due to the challenging nature of this dataset, we sampled more to get the same number of positive and negative cases as in OSS-INSTRUCT and SHAREGPT etc. As shown in Table 8, training on this dataset seems ineffective.

| | Item | MBPP | MBPP+ | HE | HE+ | Avg. |
|---|---|---|---|---|---|---|
| | Baseline | 75.7 | 64.4 | 76.8 | 70.7 | 71.9 |
| MAGICODER-S-DS | KTO | 74.9 | 64.7 | 75.4 | 72.6 | 71.9 |
| | DPO | 74.9 | 64.9 | 76.8 | 71.3 | 72.0 |
| | Baseline | 75.4 | 61.9 | 66.5 | 60.4 | 66.1 |
| MAGICODER-DS | KTO | 75.7 | 63.2 | 65.2 | 59.8 | 66.0 |
| | DPO | 75.7 | 63.2 | 65.2 | 59.8 | 66.0 |

Table 8: Training on Code Contests

## A.4. Generation of Synthetic Negatives

We present the algorithm we used to generate synthetic negatives below in Algorithm A.4.

## A.5. Results On Base Models

Below we present the results of directly applying PLUM on base models without performing supervised fine-tuning and the comparison with training using SFT in Tables 9 and 10.

---

**Algorithm 2** MutateCode Algorithm

---

**Require:** *source_code* as a string, mutation probability $P$
**Ensure:** *mutated_code* as a string
 1: **Parse** the source code into an AST: *tree* ← ParseAST(*source_code*)
 2: **Initialize** mutation rules:
 3:     Swap function arguments
 4:     Change arithmetic/logical operators
 5:     Modify control flow (negate conditions, swap if-else blocks)
 6:     Introduce off-by-one errors in loops
 7:     Remove exception handling blocks
 8:     Alter return values
 9: **Define** *Mutator* class:
10:     **Function** visit_FunctionDef(node):
11:         Store function signatures, recursively traverse AST
12:     **Function** visit_Assign(node):
13:         Track variable types, recursively traverse AST
14:     **Function** visit_Call(node):
15:         With probability $P$, swap arguments if types match
16:         With probability $P$, replace function call with another compatible one
17:     **Function** visit_If(node):
18:         With probability $P$, negate the condition or swap if-else blocks
19:     **Function** visit_For(node):
20:         With probability $P$, introduce off-by-one error in loop range
21:     **Function** visit_Try(node):
22:         With probability $P$, remove exception handling block
23:     **Function** visit_Return(node):
24:         With probability $P$, alter the return value
25: **Apply** the Mutator to the AST: *mutated_tree* ← Mutator().visit(*tree*)
26: **Perform** syntactic validation: *is_valid* ← SyntaxCheck(*mutated_tree*)
27: **if** *is_valid* = **False then**
28:      **return** original source code or error
29: **end if**
30: **Convert** the mutated AST back to code: *mutated_code* ← ASTtoSource(*mutated_tree*)
31: **return** *mutated_code*

---

| Model Families | Item | MBPP | MBPP+ | HE | HE+ |
|---|---|---|---|---|---|
| CODEQWEN-BASE | Base | 72.2 | 60.2 | 51.8 | 45.7 |
| | **OSS-Instruct** | **75.4** | **62.9** | **70.1** | **62.2** |
| | Rel. + | 4.4 | 4.5 | 35.3 | 36.1 |
| | **ShareGPT-Python** | **76.4** | **64.9** | **73.2** | **67.1** |
| | *Rel. +* | *5.8* | *7.8* | *41.3* | *46.8* |
| DS-CODER-BASE | Base | 70.2 | 56.6 | 47.6 | 39.6 |
| | OSS-Instruct | 72.9 | 58.9 | 56.7 | 48.8 |
| | Rel. + | 3.9 | 4.1 | 19.1 | 23.2 |
| | ShareGPT-Python | 75.4 | 60.7 | 64 | 53.7 |
| | Rel. + | 6.4 | 7.2 | 34.5 | 35.6 |
| STARCODER2-BASE | Base | 54.4 | 45.6 | 35.4 | 29.9 |
| | **OSS-Instruct** | **60.4** | **49.1** | **46.3** | **39.6** |
| | Rel. + | 11 | 7.7 | 30.8 | 32.4 |
| | **ShareGPT-Python** | **63.9** | **51.9** | **50** | **42.1** |
| | Rel. + | 17.5 | 13.8 | 41.2 | 40.8 |

Table 9: Results on base model training.

| Model | Type | MBPP | MBPP+ | HE | HE+ | Avg. (Base) | Avg. (+) | Avg. (All) |
|---|---|---|---|---|---|---|---|---|
| STARCODER2-BASE | **Baseline** | 54.4 | 45.6 | 35.4 | 29.9 | 44.9 | 37.8 | 41.3 |
| | **SFT** | 62.2 | 49.4 | 41.5 | 35.4 | 51.9 | 50.6 | 47.1 |
| | **PLUM-KTO** | 60.4 | 49.1 | 46.3 | 39.6 | **53.4** | **51.2** | **62.2** |
| CODEQWEN-BASE | **Baseline** | 72.2 | 60.2 | 51.8 | 45.7 | 62.0 | 53.0 | 57.5 |
| | **SFT** | 73.4 | 62.4 | 67.7 | 59.1 | 70.6 | 66.5 | 65.7 |
| | **PLUM-KTO** | 75.4 | 62.9 | 70.1 | 62.2 | **72.8** | **67.8** | **67.7** |
| DS-CODER-BASE | **Baseline** | 70.2 | 56.6 | 47.6 | 39.6 | 58.9 | 57.8 | 53.5 |
| | **SFT** | 71.7 | 57.1 | 56.1 | 48.8 | 63.9 | 60.5 | 58.4 |
| | **PLUM-KTO** | 72.9 | 58.9 | 56.7 | 48.8 | **64.8** | **61.9** | **59.3** |

Table 10: PLUM vs. SFT for base models.

## A.6. Distribution of policy model correctness

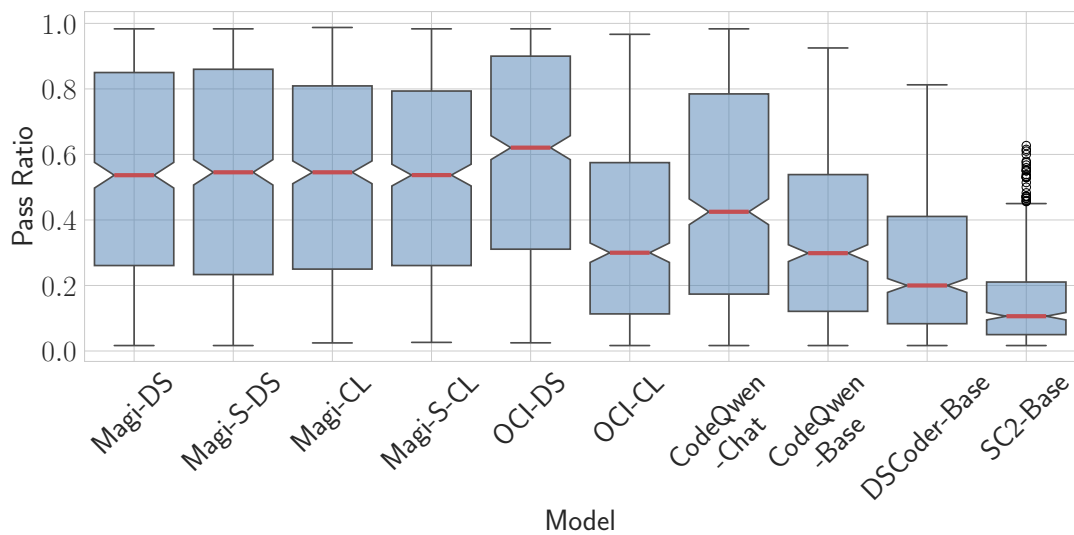Figure 3 shows the pass ratio on OSS-Instruct dataset of models we consider in this study.

**Figure** 3: Distribution of policy model correctness ratio on OSS-Instruct dataset.