PLUM: IMPROVING CODE LMS WITH EXECUTION GUIDED ON-POLICY PREFERENCE LEARNING DRIVEN BY SYNTHETIC TEST CASES

Anonymous authors

Paper under double-blind review

Abstract

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 Preference Learning framework Augmented with test cases for code LMs. 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.

030 1 INTRODUCTION

031

006

012 013

014

015

016

017

018

019

021

025

026

027

028 029

032 Language models pre-trained on code corpora have excelled at code generation (Rozière et al., 033 2024; Li et al., 2023). Supervised Fine-Tuning (SFT) enhances their ability to follow natural 034 language prompts but focuses on reproducing patterns from training data rather than ensuring code correctness (Wei et al., 2023; Zheng et al., 2024). This leads to models that generate syntactically 035 correct but functionally flawed code, unable to meet real-world requirements like edge cases or 036 algorithmic accuracy (Chen et al., 2024). Works like AlphaCode (Li et al., 2022) and LeTI (Wang 037 et al., 2024b) 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, 040 helpful, or harmless) over undesirable ones (e.g., inaccurate, unhelpful, or harmful). Despite its 041 success in aligning models with human values and improving reasoning in other domains (Dong et al., 042 2023; Guo et al., 2024a; Yuan et al., 2024; Wang et al., 2024a; Pang et al., 2024), the application of 043 preference learning as a principled and efficient approach in code generation remains under-explored, 044 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 (Setlur et al., 2024; Tajwar et al., 2024). 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 (Tajwar et al., 2024; Guo et al., 2024b; Setlur et al., 2024; Liu et al., 2024). Unlike offline preference data, on-policy data remains in-distribution with the model, reducing the risk of misalignment (Guo et al., 2024b; Zhang et al., 2024; Tang et al., 2024; Fisch et al., 2024).
The main challenge now is how to efficiently obtain preference labels for on-policy data at scale (Yang et al., 2024b). 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.

057 058	Prompt for Test Case Generation
059	You are a teaching assistant helping to write reference solutions and tests for programming questions.
060	Given a programming question, you need to first analyze the problem, then write a reference solution
061	(code), followed by assertions that test student solutions. The test code must be runnable when
062	concatenated at the end of student solutions to check the correctness.
063	Programming Question:
064	{Question}
065	
066	Follow the format below:
067	[Analysis] {{Natural language analysis of the problem }}
068	[Solution]
069	{{Your solution to the problem}}
070	[Start Code]
071	{{Start code for students so that they can follow the 1/0 protocol.
072	[Test Code]
073	{{Test code that is immediately runnable if concatenated with
074	student code to check the correctness.}}
075	

077 To this end, we propose preference learning framework augmented with test cases for training code 078 language models (PLUM), which integrates the process of deriving test cases from natural language 079 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 (Wei et al., 2023), Evol-Instruct-Code (Luo et al., 2023), and ShareGPT (Cleston, 2023). For 081 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 083 based on the results: solutions passing the tests are preferred, while failures are dis-preferred. This 084 dataset then trains the policy using established preference learning algorithms (Rafailov et al., 2023; 085 Ethayarajh et al., 2024; Azar et al., 2023). 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 087 shifts and poor generalization commonly observed for synthetic off-policy data, enhancing the 088 model's robustness and improving its ability to differentiate between correct and incorrect solutions. 089 In addition, by showcasing the effectiveness of our framework, we demonstrated the feasibility of bypassing tedious (and potentially unstable) reward model training (Liu et al., 2023a; Le et al., 2022) and manual labeling, by automating the process of test case synthesis. These simplicity advantages 091 of PLUM makes online preference training of language models possible. 092

We evaluate PLUM on a diverse set of state-of-the-art code language models under different set-ups, on commonly used benchmarks: HumanEval(+) and MBPP(+) (Chen et al., 2021; Austin et al., 2021; Liu et al., 2023b) as well as more challenging code generation datasets like LiveCodeBench and LeetCode (Jain et al., 2024; Guo et al., 2024a). 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 (Xiong et al., 2024b).

101 102

103

076

2 PREFERENCE LEARNING AUGMENTED WITH TEST CASES FOR CODE LMS (PLUM)

104 105

The core of PLUM lies in leveraging recent advancements in on-policy and online preference learning,
 which have proven effective across various domains (Xiong et al., 2024b;a; Mitra et al., 2024). In
 the context of code generation, PLUM simplifies and scales the preference data collection process

by using test cases. These test cases act as a lightweight yet robust mechanism to evaluate 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.

Algorithm 1 PLUM.

117

Inpu f	It: Natural language instructions $\mathcal{I} = \{q_i\}$, policy model to be trained π_{θ} , greateners T , chunk size M \triangleright Unified for both of	enerator model G, update fline and online alignment
Outr	put: Trained policy model π'_{θ}	
1: Î	nitialize preference datasets \mathcal{D}^+ and \mathcal{D}^-	
2: f	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 $\{(r_{ij}, t_{ij})\}_{i=1}^n$ 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:	$S_{ik} \sim \pi_{\theta}$ for $k = 1$ to K (sample K solutions for q_i)	On-policy sampling
11:	for each solution $s_{ik} \in S_{ik}$ do	
12:	for each test case t_{ij} in \mathcal{D} do	
13:	if s_{ik} fails t_{ij} then	
14:	Add (q_i, s_{ik}) to \mathcal{D}^-	▷ Negative case
15:	end if	
16:	Add (q_i,s_{ik}) to \mathcal{D}^+	Positive case
17:	end for	
18:	end for	
19:	end for	
20:	Filter out instances with no correct solutions from \mathcal{D}^+ and \mathcal{D}^-	
21:	if iteration count $\%T = 0$ then	▷ Policy Update
22:	Train the policy model π_{θ} using \mathcal{D}^+ and \mathcal{D}^- with preference learning to	$\sigma \operatorname{get} \pi'_{\theta}$
23:	Update policy model $\pi_{\theta} = \pi'_{\theta}$	
24:	end if	
25: 0	end for	
26: 1	return $\pi_{ heta}$	

143 144

145

146

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 produces multiple test cases for the batch of instructions. Then we sample solutions from policy π_{θ} , and execute them against the generated test suite to obtain preference labels. We then update the policy $\pi_{\theta} := \pi'_{\theta}$.

151 152

153

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.

157 The test cases in PLUM are generated with a test-case generator model over natural instructions 158 from established code generation datasets.¹ In automated testing, ensuring the correctness and 159 completeness of test cases is a persistent challenge due to the lack of reliable oracles to validate 160 test outputs. We adopt two strategic principles: 1) employing self-consistency as an approximate

¹⁶¹

¹We use GPT-4-1106 as the generator model.

oracle, and 2) generating diverse test suites to minimize overfitting to any particular test instance and
 mitigate under-specification.

165 **Collecting instructions from established datasets** We collect natural language instructions from 166 established datasets including OSS-Instruct (Wei et al., 2023), Evol-Instruct-Code (Luo et al., 2023), 167 and ShareGPT (Cleston, 2023). ²These datasets provide a diverse range of programming tasks and 168 instructions. Although they come with gold/silver solutions in the training splits, these solutions are 169 never used in PLUM. Instead, they allow us to directly compare PLUM's performance against SFT, 170 which relies on gold solutions, as we investigate in our experiments. Not requiring gold solutions for training broadens the applicability of PLUM to a wide variety of real-world coding tasks and user 171 requirements. 172

173

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 consistencybased approach inspired by Chen et al. (2023) and Rozière et al. (2024) 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.

- 186
- 187 188

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 Rafailov et al. (2023); Ethayarajh et al. (2024); Azar et al. (2023). In practice, however, preference data often contains solutions sampled from different models than the policy, leaving the data out of distribution Bai et al. (2022); Yuan et al. (2024). A common workaround is to first perform supervised fine-tuning (SFT) on the same instructions before applying preference learning (Rafailov et al., 2023; Yuan et al., 2024). 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 Ksolutions from the policy and evaluate them against the

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

solutions from the policy and evaluate them against the Table 1: Self-Consistency Pass Rate Usgenerated test cases. *K* is set to 20 based on the findings ing GPT-4-1106.

from our preliminary experiments. 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 (Wang et al., 2024b). Solutions passing all test cases are used as the chosen solutions, and those failing at least one the rejected solutions.

209 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 Rafailov et al. (2023).

213

 ²We focus on Python due to its wide use and the availability of well-established training and evaluation
 resources.

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

Model	Item	MBPP	MBPP+	HE	HE+	Avg
	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	78.0
CODEQWEN-1.5-CHAT (Bai et al., 2023)	Cond.Err.Msg	79.4	68.7	84.1	79.3	77.9
	PLUM-DPO	81.2	70.2	86.0	81.1	79.
	PLUM-KTO	81.0	69.0	86.0	81.1	79.
	Rel. +	4.3	2.7	3.0	3.1	3.3
	Baseline	74.9	65.6	75.4	71.3	71.
	Cond.Token	73.4	62.9	76.8	70.7	71.
	RFT	74.7	64.9	74.4	66.5	70.
DS-CODER-INSTRUCT (Guo et al., 2024a)	Cond.Err.Msg	74.7	64.2	80.5	75.0	73.0
	PLUM-DPO	76.4	65.9	80.5	76.8	77.4
	PLUM-KTO	78.2	67.9	81.7	76.8	76.2
	<i>Rel</i> . +	4.4	3.5	8.4	7.7	6.0
	Baseline	75.4	61.9	66.5	60.4	66.
	Cond.Token	75.9	62.4	68.3	62.2	67.
MAGICODER-DS (Wei et al., 2023)	RFT	76.2	62.2	67.7	62.2	67.
	Cond.Err.Msg	74.2	62.2	66.5	59.8	65.
	PLUM-DPO	75.9	63.7	67.7	61.6	67.
	PLUM-KTO	79.6	66.7	71.3	65.9	70.
	<i>Rel.</i> +	5.6	7.8	7.2	9.1	7.4
	Baseline	75.7	64.4	76.8	70.7	71.9
	Cond.Token	73.9	63.7	75.0	71.3	71.
	RFT	75.4	64.4	73.2	69.5	70.
MAGICODER-S-DS (Wei et al., 2023)	Cond.Err.Msg	75.2	65.4	75.6	70.7	71.
	PLUM-DPO	76.2	64.7	78.7	73.8	73.
	PLUM-KTO	80.4	69.3	80.5	73.8	76.
	<i>Rel.</i> +	4.4	7.4	4.4	4.5	5.2
	Baseline	73.9	63.7	77.4	72.0	71.
	Cond.Token	73.9	62.9	75.6	71.3	70.
	RFT	74.2	63.7	76.8	72.0	71.
OCI-DS (Zheng et al., 2024)	Cond.Err.Msg	75.0	70.7	74.4	64.7	71.
	PLUM-DPO	76.4	66.4	80.5	76.2	74.
	PLUM-KTO	78.2	66.4	80.5	76.2	75.
	Rel. +	5.8	4.2	4.0	5.8	5.0
	Baseline	66.4	55.4	72.6	65.2	64.
	Cond.Token	59.6	48.4	23.2	21.3	38.
	RFT	63.2	52.6	60.4	56.7	58.
OCI-CL (Zheng et al., 2024)	Cond.Err.Msg	67.9	55.9	68.9	65.2	64.
	PLUM-DPO	66.4	55.9	71.3	65.2	64
	PLUM-KTO	66.7	55.4	73.8	67.7	65.
	Ral +	0.5	0.0	17	3.8	15

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.

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 Rafailov et al. (2023) and Kahneman-Tversky Optimization (Ethayarajh et al., 2024) that have been shown to bring improvements for reasoning tasks (Mitra
et al., 2024; Yuan et al., 2024; Dubey et al., 2024). 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.



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) preference learning.

3 EXPERIMENTS

To demonstrate the effectiveness of PLUM, we evaluate it on established benchmarks: HumanEval (Chen et al., 2021), MBPP (Austin et al., 2021), and EvalPlus (a widely-adopted augmented version (Liu et al., 2023b) of them). We also use the more challenging LiveCodeBench (Jain et al., 2024).

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 powerful language models to generate test cases for each programming question. The results reported in the main text use test cases generated by GPT-4 OpenAI et al. (2024). An ablation of test case generators is in Appendix A.5. 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 us

We then sample 20 outputs from the policy us-306 ing temperature T = 1 for the former two and 307 50 outputs for the latter. This yields around $\sim 60,000$ examples for ShareGPT and OSS-308 Instruct, and around $\sim 120,000$ for EvolIn-309 struct before any filtering. We present the statis-310 tics on the pass ratio of sampled solutions over 311 OSS-Instruct in Figure 3. We included the self-312 consistency pass rate of the test-generation pro-313 cess with GPT-4-1106 in Table 1.

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

Table 3: Comparison with Reflexion (Shinn et al., 2023)

314

284

285

287

288 289

290

291

292

293

³¹⁵ **Models** We consider a diverse set of strong

open language models: MagiCoder (Wei et al., 2023), OpenCodeIntepreter (Zheng et al., 2024),
CodeQwen (Bai et al., 2023), DeepSeek Coder (Guo et al., 2024a) and StarCoder2 (Li et al., 2023).
MagiCoder and OpenCodeIntepreter contain instruction-tuned checkpoints from DeepSeek Coder
and CodeLlama (Rozière et al., 2024) base models. In the main text, we focus on instruction-tuned
language models, while differing results from directly training base models to Appendix A.4

321

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* (Shinn et al., 2023), using

Model Families	Data	Item	MBPP	MBPP+	HE	HE+	Avg
WIZARDCODER		Baseline	48.2	40.9	56.6	47.1	48.2
(Luc et al. 2022)		PLUM-KTO	54.3	48.8	65.9	52.9	55.5
(Luo et al., 2023)	EVOL INSTRUCT	Rel. +	16.4	12.3	12.7	19.3	15.2
	- EVOLINSIKUUI	Baseline	74.9	65.6	75.4	71.3	71.8
DS-CODER-INSTRUCT		PLUM-KTO	77.7	67.7	81.7	76.8	76.0
		Rel. +	3.7	3.2	8.4	7.7	5.8
		Baseline	74.9	65.6	75.4	71.3	71.8
DS-CODER-INSTRUCT		PLUM-KTO	79.2	67.9	77.4	73.8	74.6
		<i>Rel.</i> +	5.7	3.5	2.8	3.5	3.9
	- SHADECDT	Baseline	73.9	63.7	77.4	72.0	71.8
OCI-DS	DATION	PLUM-KTO	77.7	64.4	79.9	75.6	74. 4
	-PYTHON	Rel. +	5.1	1.1	3.2	5.0	3.6
	-	Baseline	77.7	67.2	83.5	78.7	76.8
CODEQWEN-1.5-CHAT		PLUM-KTO	81.2	69.7	85.4	79.3	78.9
		<i>Rel.</i> +	4.5	3.7	2.3	0.8	2.8

Table 4: PLUM on other datasets.

instruction-tuned (SFT) models and value-conditioning techniques (Li et al., 2022; Wang et al., 2024b). 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.

352 353

354

338

339 340 341

342

343

344

345 346

347

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.

357 MagiCoder models (-DS, -S-CL, and -S-DS) and OpenCodeIntepreter models (-CL and -DS) have 358 already seen these instructions during supervised fine-tuning, while DeepSeekCoder-Instruct has 359 not, as it was released earlier than the dataset. CodeQwen chat model uses proprietary data. PLUM-360 ShareGPT data for preference learning is generated with the same setting. Similarly, Table 4 corresponds to the results when we apply PLUM to EvolInstruct (Luo et al., 2023) dataset. Since 361 the instructions are comparatively less clear than the OSS-Instruct dataset, we control the number of 362 initial samples to be the same by generating 50 samples for each problem and use about 400 instances 363 in total. 364

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 OSSInstruct or ShareGPT data. These results demonstrate that PLUM is broadly applicable in different
datasets and settings.

We noticed that PLUM-KTO consistently out-performs the baseline, and that PLUM-DPO sometimes under-perform PLUM-KTO. Prior works (Mitra et al., 2024; Yuan et al., 2024) noticed the
phenomenon where DPO can exhibit instability due to reducing reward for the positive class.

373

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 per formance 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.



Figure 2: Ablation studies on preference training signals show that merely using un-runnable code as negative instances does not consistently enhance performance. In contrast, PLUM effectively improves the model by introducing functional correctness signals. Baseline results refer to the SFT model without PLUM.

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
	<i>Rel.</i> +	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
	<i>Rel.</i> +	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
	<i>Rel.</i> +	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
	<i>Rel.</i> +	-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
	<i>Rel.</i> +	0.7	20	230	11.2

Table 5: %Pass@1 on 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.

421 Comparison Against Baselines As shown in Table 3, verbal reinforcement learning like Reflex 422 ion (Shinn et al., 2023), does not perform well on code language models fine-tuned for code at the
 423 scale relevant to our work. This is partly due to the limitations of smaller LLMs in handling various
 424 types of instructions effectively.

Approaches like LeTI (Wang et al., 2024b) 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 Wang et al. (2024b) 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 lead to reduced performance, likely due to the
 token addition to tokenizer and the absence of clear labels distinguishing good from bad outputs.

	•			Off-Policy			
	PLUM	DS-C -1. -Inst	oder Qv 3B Cod ruct -Ii	ven2.5- ler-1.5B ıstruct	Codestral -22B	DeepSeek -Coder-33B -Instruct	Syn -Neg
CodeQwen1.5-Ci	нат 79.3	78	.4	78.1	77.3	79.0	73.9
S-CODER-INSTRU	UCT 76.2	76	.0	75.2	74.6	74.0	69.4
MAGICODER-D	S 70.9	68	.4	68.6 75.2	68.3	67.0	61.2
MAGICODER-S-L	70.0	15	.0	13.3	/4.1	/3.8	/1.
Table 6: Compariso	on between on-/off	-policy p	preference les	arning us	ing KTO on (CodeQwen-1.5	-Сна
		Easy	Medium	Hard	Average	Rel. Gain	
CODE	QWEN1.5-CHAT	66.7	27.5	13.6	33.9	-	
	+PLUM-DPO	66.7	31.9	15.9	36.7	8.3	
+P	LUM-DPO-Iter	66.7	31.9	18.2	37.2	9.7	
P	+PLUM-KTO	68.9	30.8	11.4	35.6	5.0	
+P	LUM-KTO-Iter	66.7	31.9	18.2	37.2	9.7	
MA	DI LIM DDO	33.5 11 1	1/.0	9.1 12.4	19.4	-	
ı D		44.4 46.7	13.4 17.6	13.0	22.2	14.4 20 1	
TI	+PLUM-KTO	48.9	16.5	91	23.5	17.5	
+P	LUM-KTO-Iter	44.4	17.6	11.4	22.8	17.5	
			1710		22.0	1/10	
nce. To test this, we ff-policy). We sele	e apply the same cted models of v	method	whether "o but use pre sizes, inclu	n-policy ference of ding Dec	" training m lata sampled epSeek-Cod	hakes a significated from other model of the second	ant di odels (ct, Qv
ence. To test this, we off-policy). We sele 2.5-Coder-1.5B-Inst followed the exact sa raining still improve These findings are of the al 2024. Xu et al	e apply the same cted models of v ruct, CodeStral-2 ame data collections so model perform consistent with p	method arying 22B (M on and t ance, it prior res	whether "o but use pre sizes, inclu- istralAI, 20 raining pro- generally u earch (Xic	n-policy ference of ding Deo (24), and cesses. C nderperf ong et al	" training m data sampled epSeek-Cod DeepSeek- Our results sl orms compa ., 2024a; D	hakes a significant of from other model ler-1.3B-Instruct Coder-33B-Instruct coder-34B-Instruct coder-34B-In	ant di odels (ct, Qv truct. off-po y trair 4; Taj
nce. To test this, we ff-policy). We sele .5-Coder-1.5B-Inst ollowed the exact sa raining still improve these findings are of t al., 2024; Xu et al	e apply the same cted models of v ruct, CodeStral-2 ame data collection es model perform consistent with p l., 2024).	method arying 22B (M on and t ance, it orior res	whether "o but use pre sizes, inclu- istralAI, 20 raining proo generally u earch (Xio	n-policy ference of ding Dec (24), and cesses. C nderperf ong et al	" training m data sampled epSeek-Cod DeepSeek- Dur results sl orms compa ., 2024a; D	hakes a significated from other model of the signal from the signal content of the signal from	ant di odels ct, Qv truct. off-pc v train 4; Taj
nce. To test this, we ff-policy). We sele .5-Coder-1.5B-Inst blowed the exact sa aining still improve hese findings are of t al., 2024; Xu et al urther, we investig	e apply the same cted models of v rruct, CodeStral-2 ame data collection es model perform consistent with p l., 2024).	address method arying (22B (M on and t ance, it orior res	whether "o but use pre sizes, inclu- istralAI, 20 raining proo generally u earch (Xio tic negative pressinte Press	n-policy ference of ding Dea (24), and cesses. Conderperf ong et al	" training m data sampled epSeek-Cod DeepSeek- Dur results sl orms compa ., 2024a; D nis effect, w	hakes a significated from other modeler-1.3B-Instruct Coder-33B-Instruct Coder-33B-Instruct now that while our on-policy ong et al., 2024 we use a mutative aintaining its of	ant di odels ct, Qv truct. off-pc y train 4; Taj on-ba
nce. To test this, we ff-policy). We sele .5-Coder-1.5B-Inst ollowed the exact sa aining still improve 'hese findings are of t al., 2024; Xu et al further, we investig pproach for synthe orrectness, as detai	e apply the same cted models of v ruct, CodeStral-2 ame data collection es model perform consistent with p l., 2024). gate the effect of etically introduc iled in Appendix	address method arying 22B (M on and t ance, it rior res synthe ing erro A.8. Th	whether "o but use pre sizes, inclu- istralAI, 20 generally u earch (Xic tic negative ors into Py is method u	n-policy ference of ding Deo (24), and cesses. Of nderperf ong et al es. To the thon coo ses Abst	" training m data sampled epSeek-Cod DeepSeek- Our results sl orms compa ., 2024a; D nis effect, w de while m ract Syntax '	hakes a significand d from other model ler-1.3B-Instruct Coder-33B-Instruct coder-33B-Instruct on that while our model to on-policy ong et al., 2024 we use a mutative aintaining its so Tree (AST) man	ant di odels ct, Qv truct. off-pc / train 4; Taj on-ba synta
nce. To test this, we ff-policy). We sele .5-Coder-1.5B-Inst blowed the exact sa aining still improve hese findings are of t al., 2024; Xu et al urther, we investig pproach for synthe prectness , as detai	e apply the same cted models of v ruct, CodeStral-2 ame data collection es model perform consistent with p l., 2024). gate the effect of etically introduc iled in Appendix 2 ike argument swa	address method arying 22B (M on and t ance, it rior res synthe ing erro A.8. Th apping	whether "o but use pre sizes, inclu- istralAI, 20 raining pro- generally u earch (Xic tic negative ors into Py is method u operator re	n-policy ference of ding Dee (24), and cesses. Of nderperf ong et al es. To the thon coordinates ses Abstriplaceme	" training m data sampled epSeek-Cod DeepSeek- Our results sl orms compa ., 2024a; D his effect, w de while m ract Syntax ' nt. control f	hakes a significant of from other model ler-1.3B-Instruct Coder-33B-Instruct coder-33B-Instruct on that while of the to on-policy ong et al., 2024 we use a mutative aintaining its so Tree (AST) man low changes, of	ant di odels ct, Qv truct. off-pc y train 4; Taj on-ba synta ipula ff-bv-
nce. To test this, we ff-policy). We sele 5-Coder-1.5B-Inst billowed the exact sa aining still improve hese findings are of al., 2024; Xu et al urther, we investig pproach for synthe prectness , as detai apply mutations 1 rors, and return v	e apply the same cted models of v ruct, CodeStral-2 ame data collection consistent with p l., 2024). gate the effect of etically introduct led in Appendix 2 ike argument swa	address method arying 22B (M on and t ance, it rior res synthe ing erro A.8. Th apping, b By in	whether "o but use pre sizes, inclu- istralAI, 20 raining pro- generally u earch (Xic tic negative ors into Py is method u operator re jecting the	on-policy ference of ding Dec (24), and cesses. Conderperf ong et al es. To the thon coordinates ses Abstriplaceme se errors	" training m data sampled epSeek-Cod DeepSeek- Our results sl orms compa ., 2024a; D his effect, w de while m ract Syntax ' nt, control f	hakes a significated from other model of the from other model of the from t	ant di odels ct, Qv truct. off-pc y train 4; Taj on-ba synta iipula ff-by- avior
nce. To test this, we ff-policy). We sele 5-Coder-1.5B-Inst billowed the exact sa aining still improve hese findings are of al., 2024; Xu et al urther, we investig pproach for synthe prectness, as detail apply mutations 1 trors, and return va- tered code. We app	e apply the same cted models of v ruct, CodeStral-2 ame data collection es model perform consistent with p l., 2024). gate the effect of etically introduc led in Appendix 2 ike argument swa alue modification ply this to on-po	address method arying 22B (M on and t ance, it rior res synthe ing erro A.8. Th apping, 1. By in licy pos	whether "o but use pre sizes, inclu- istralAI, 20 raining proo generally u earch (Xic tic negative ors into Py is method u operator re jecting these jecting these jecting these	n-policy ference of ding Dec (24), and cesses. Of nderperf ong et al es. To the thon coordinates ses Abstriplaceme se errors	" training m data sampled epSeek-Cod DeepSeek- Our results sl orms compa ., 2024a; D nis effect, w de while m ract Syntax ' nt, control f , it generate ating off-po	hakes a significated from other model of the significated from other model of the signal of the sign	ant di odels ct, Qv truct. off-po v train 4; Taj on-ba synta iipula ff-by- avior
nce. To test this, we ff-policy). We sele .5-Coder-1.5B-Inst blowed the exact sa aining still improve hese findings are of t al., 2024; Xu et al urther, we investig pproach for synthe orrectness , as detail o apply mutations 1 rrors, and return va- ltered code. We ap- aining under the sa	e apply the same cted models of v ruct, CodeStral-2 ame data collection es model perform consistent with p l., 2024). gate the effect of etically introduc led in Appendix 2 ike argument swa alue modification ply this to on-po ame setup. Obse	address method arying 22B (M on and t ance, it rior res synthe ing erro A.8. Th apping, b. By in licy pos rying th	whether "o but use pre sizes, inclu- istralAI, 20 raining proo generally u earch (Xio tic negative ors into Py is method u operator re jecting these itive example at the synt	n-policy ference of ding Dec (24), and cesses. C nderperf ong et al es. To the thon co- ses Absti- placeme se errors ples, crea- hetic new	" training m data sampled epSeek-Cod DeepSeek- Our results sl orms compa ., 2024a; D nis effect, w de while m ract Syntax ' nt, control f , it generate ating off-po gatives coul	hakes a significated from other model of the significated from other model of the signification of the significati	ant di odels ct, Qv truct. off-pc / train 4; Taj on-ba synta iipula ff-by- avior or mo
nce. To test this, we ff-policy). We sele .5-Coder-1.5B-Inst collowed the exact sa raining still improve these findings are of t al., 2024; Xu et al ourther, we investig pproach for synthe orrectness , as detain to apply mutations la rrors, and return van latered code. We appraining under the sa re performance, we	e apply the same cted models of v ruct, CodeStral-2 ame data collection es model perform consistent with p l., 2024). gate the effect of etically introduc led in Appendix 2 ike argument swa alue modification ply this to on-po ame setup. Obse e confirmed the in	address method arying (22B (M on and t ance, it prior res synthe ing erro A.8. Th apping, h. By in licy pos rving th pportar	whether "o but use pre sizes, inclu- istralAI, 20 raining proo generally u earch (Xio tic negative ors into Py is method u operator re jecting the sitive examp at the synt ce of prefe	n-policy ference of ding Dec (24), and cesses. C nderperf ong et al es. To the thon co- ses Absti- placeme se errors ples, crea- hetic neg- rence tra	" training m data sampled epSeek-Cod DeepSeek- Our results sl orms compa ., 2024a; D nis effect, w de while m ract Syntax ' nt, control f , it generate ating off-po gatives coul ining with	hakes a significated from other model from other modeler-1.3B-Instruct Coder-33B-Instruct Coder-33B-Instruct ow that while our of the on-policy ong et al., 2024 we use a mutative aintaining its so Tree (AST) man low changes, of the so valid but beh licy negatives for d even potentiated and the other of the other nore natural, and the other other other other and the other other other other other other other so the other oth	ant di dels ct, Qv truct. off-pc / train i; Taj on-ba synta synta i pula ff-by- avior for mo or
nce. To test this, we ff-policy). We sele .5-Coder-1.5B-Inst oblowed the exact sa raining still improve these findings are of t al., 2024; Xu et al further, we investig pproach for synthe orrectness , as detail to apply mutations 1 rrors, and return val ltered code. We ap raining under the sa ne performance, we n-policy negative s	e apply the same cted models of v rruct, CodeStral-2 ame data collection es model perform consistent with p l., 2024). gate the effect of etically introduce led in Appendix 2 ike argument swa alue modification ply this to on-po ame setup. Obse e confirmed the in amples.	Synthe apping, acception on and t ance, it rior res synthe ing erro A.8. Th apping, h. By in licy pos rving th mportan	whether "o but use pre sizes, inclu- istralAI, 200 raining proo generally u earch (Xio tic negative ors into Py is method u operator re jecting the sitive examp- nat the synt ace of prefe	n-policy ference of ding Dec (24), and cesses. C nderperf ong et al es. To the thon coordinates ses Abstriplaceme se errors placeme se errors ples, creates hetic negrence tra	" training m data sampled epSeek-Cod DeepSeek- our results sl orms compa ., 2024a; D his effect, w de while m ract Syntax ' nt, control f , it generate ating off-po gatives coul ining with r	hakes a significated from other model from other modeler-1.3B-Instruct Coder-33B-Instruct Coder-33B-Instruct ow that while our of the on-policy ong et al., 2024 we use a mutativation of the output of the output we use a mutativation of the output of the	ant di dels (ct, Qv truct. off-pc / train i; Taj on-ba yynta avior or mo or mo or mo illy h id ide
nce. To test this, we ff-policy). We sele .5-Coder-1.5B-Inst ollowed the exact sa aining still improve hese findings are of t al., 2024; Xu et al ourther, we investig pproach for synthe orrectness , as detai to apply mutations 1 rrors, and return val letered code. We ap raining under the sa ne performance, we n-policy negative s	e apply the same cted models of v rruct, CodeStral-2 ame data collection es model perform consistent with p l., 2024). gate the effect of etically introduc led in Appendix 2 alue modification ply this to on-po ame setup. Obse e confirmed the in amples.	address method arying 22B (M on and t ance, it rior res synthe ing erro A.8. Th apping, h. By in licy pos rving th mportan	whether "o but use pre sizes, inclu- istralAI, 20 raining proo generally u earch (Xio tic negative ors into Py is method u operator re jecting thes tive examp at the synt ce of prefe-	n-policy ference of ding Dec (24), and cesses. Conderperf ong et al es. To the thon coordinates ses Abstriplaceme se errors ples, creates hetic negrence trans	" training m data sampled epSeek-Cod DeepSeek-Cod DeepSeek- Our results sl orms compa ., 2024a; D his effect, w de while m ract Syntax ' nt, control f , it generate ating off-po gatives coul ining with r	hakes a significated from other model of the significated from other model of the signification of the significati	ant di dels ct, Qv truct. off-pc / train trai train trai trai train train trai train train train train train train train train train train train train trai train trai train trai trai train trai trai trai train trai trai trai trai trai trai trai trai
nce. To test this, we off-policy). We sele 2.5-Coder-1.5B-Inst ollowed the exact sa raining still improve These findings are of t al., 2024; Xu et al Further, we investig pproach for synthe orrectness , as detail to apply mutations 1 errors, and return val ltered code. We ap raining under the sa he performance, we on-policy negative s This highlights our of fficient on-policy pro-	e apply the same cted models of v ruct, CodeStral-2 ame data collection es model perform consistent with p l., 2024). gate the effect of etically introduc- led in Appendix 2 ike argument swa alue modification ply this to on-po- ame setup. Obse e confirmed the in amples.	address method arying 22B (M on and t ance, it rrior res synthe ing erro A.8. Th apping, 1. By in licy pos rving th nportan emonstr g, This	whether "o but use pre sizes, inclu- istralAI, 200 raining proo generally u earch (Xio tic negative ors into Py is method u operator re jecting these tive example the synt is configure at the synt at the synt at a met	n-policy ference of ding Dec (24), and cesses. C nderperf ong et al es. To the thon coordinate ses Abstriplaceme se errors ples, creation hetic neg rence trans	" training m data sampled epSeek-Cod DeepSeek-Cod DeepSeek-Our results sl orms compa ., 2024a; D nis effect, w de while m ract Syntax ' nt, control f , it generate ating off-po gatives coul ining with r	hakes a significated from other model of the significated from other model of the signification of the significati	ant di dels ct, Qv truct. off-pc / train 1; Taj on-ba synta ff-by- avion or me ally h ad ide cases
nce. To test this, we off-policy). We sele 2.5-Coder-1.5B-Inst collowed the exact sa raining still improve These findings are of t al., 2024; Xu et al Further, we investig pproach for synthe orrectness , as detain o apply mutations 1 errors, and return va ltered code. We ap raining under the sa he performance, we on-policy negative s This highlights our of fficient on-policy prest est cases, providing	e apply the same cted models of v ruct, CodeStral-2 ame data collection es model perform consistent with p l., 2024). gate the effect of etically introduc led in Appendix 2 ike argument swa alue modification ply this to on-po ame setup. Obse e confirmed the in amples.	address method arying (22B (M on and t ance, it rior res synthe ing erro A.8. Th apping, h. By in licy pos rving th mportan emonstr g. This sion sis	whether "o but use pre sizes, inclu- istralAI, 200 raining proo generally u earch (Xio tic negative ors into Py is method u operator re jecting thes titive example the synt is conformed at the synt at the synt at the synt at the synt at general proof the synt the synt at the synt at the synt the synt the synt at the synt the synt at the synt the synt synt the synt sy	n-policy ference of ding Dec (24), and cesses. C nderperf ong et al es. To the thon coordinate ses Abstriplaceme se errors ples, create hetic neg rence trans hod emp an be eas del traini	" training m data sampled epSeek-Cod DeepSeek-Cod DeepSeek-Our results sl orms compa ., 2024a; D nis effect, w de while m ract Syntax ' nt, control f , it generate ating off-po gatives coul ining with n oowered by a ily adopted ing.	hakes a significated from other model from other modeler-1.3B-Instruct Coder-33B-Instruct Coder-33B-Instruct how that while our of the on-policy ong et al., 2024 we use a mutative aintaining its as Tree (AST) man how changes, of the other of the other of the licy negatives for d even potentian nore natural, are automated test of to scale the coll	ant di dels ct, Qv truct. off-pc / train 4; Taj on-ba synta iipula iff-by- avior or mod illy h id ide cases lectio
nce. To test this, we ff-policy). We sele .5-Coder-1.5B-Inst blowed the exact sa aining still improve these findings are of t al., 2024; Xu et al urther, we investig pproach for synthe orrectness , as detail o apply mutations 1 rrors, and return va- ltered code. We ap 'aining under the sa- ne performance, we n-policy negative s his highlights our of fficient on-policy pre- vst cases, providing	e apply the same cted models of v ruct, CodeStral-2 ame data collection es model perform consistent with p l., 2024). gate the effect of etically introduc led in Appendix 2 ike argument swa alue modification ply this to on-po ame setup. Obse e confirmed the in amples. contribution in de reference learning a robust supervise	address method arying 5 22B (M on and t ance, it prior res synthe ing erro A.8. Th apping, h. By in licy pos rving th mportan emonstr g. This sion sig	whether "o but use pre sizes, inclu- istralAI, 20 raining proo generally u earch (Xio tic negative ors into Py is method u operator re jecting these itive examp- nat the synt ice of prefe- ating a met approach ca- gnal for moo	n-policy ference of ding Dec (24), and cesses. Conderperf ong et al es. To the thon coordinates ses Abstriplaceme se errors placeme se errors placeme transplaceme se errors placeme se errors placeme transplaceme se errors placeme se errors placeme se errors placeme transplaceme se errors placeme transplaceme se errors placeme transplaceme tr	" training m data sampled epSeek-Cod DeepSeek-Cod DeepSeek-Our results sl orms compa ., 2024a; D nis effect, w de while m ract Syntax ' nt, control f , it generate ating off-po gatives coul ining with n owered by a ily adopted ing.	hakes a significated from other model from other modeler-1.3B-Instruct Coder-33B-Instruct Coder-33B-Instruct how that while our of the on-policy ong et al., 2024 we use a mutative aintaining its so Tree (AST) man how changes, of the other of the other of the licy negatives for d even potentian nore natural, are automated test of to scale the coll	ant di dels (ct, Qv truct. off-pc / train 4; Taj on-ba synta avior fr-by- avior fr-by- avior fr hy- di de cases lectio
ice. To test this, we f-policy). We sele 5-Coder-1.5B-Inst llowed the exact sa aining still improve nese findings are of al., 2024; Xu et al urther, we investig poroach for synthe proach for synthe protectness, as detail apply mutations l rors, and return va tered code. We ap aining under the sa e performance, we h-policy negative s his highlights our of ficient on-policy p st cases, providing	e apply the same cted models of v ruct, CodeStral-2 ame data collection es model perform consistent with p l., 2024). gate the effect of etically introduc led in Appendix 2 ike argument swa alue modification ply this to on-po ame setup. Obse e confirmed the in amples. contribution in de reference learning a robust supervi	address method arying ; 22B (M on and t ance, it rior res 'synthe ing erro A.8. Th apping, h. By in licy pos rving th nportan emonstr g. This sion sig	whether "o but use pre sizes, inclu- istralAI, 20 raining proo generally u earch (Xio tic negative ors into Py is method u operator re jecting the stitive example the synt is of prefe- ating a met approach ca gnal for mo	n-policy ference of ding Dec (24), and cesses. Conderperf ong et al es. To the thon coordinates placeme se errors placeme se errors placeme transplaceme se errors placeme transplaceme se errors placeme se errors placeme se errors placeme se errors placeme se errors placeme transplaceme se errors placeme se errors placeme transplaceme se errors placeme se errors placeme transplaceme se errors placeme transplaceme se errors placeme transplaceme transplaceme se errors placeme transplaceme se errors placeme transplaceme tra	" training m data sampled epSeek-Cod DeepSeek-Cod DeepSeek-Cod DeepSeek-Cod orms compa ., 2024a; D nis effect, w de while m ract Syntax ' nt, control f , it generate ating off-po gatives coul ining with n powered by a ily adopted ing.	hakes a significated from other model from other modeler-1.3B-Instruct Coder-33B-Instruct Coder-33B-Instruct now that while our of the on-policy ong et al., 2024 we use a mutative aintaining its served (AST) man low changes, of the ven gotentiated by the licy negatives features for the other anore natural, are automated test of to scale the coll with including	ant di dels ct, Qv truct. off-pc v train 4; Taj on-ba synta iipula ff-by- avior for mo- ully h id ide cases lectio

studies (Wang et al., 2024b). 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.
Test Cases Allow By-passing Reward Model Training For Iterative Alignment We demonstrate

always improve the model's performance and may even hurt. This has also been noted in previous

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 (Xiong et al., 2024a;b), 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) (Rafailov et al., 2023) and Kahneman & Tversky's Optimization(KTO) (Ethayarajh et al., 2024) 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 (Mitra et al., 2024) uses iterative preference learning to improve math reasoning in SFT-ed language models, while Eurus (Yuan et al., 2024) 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 (Li et al., 2023), StarCoder2 (Lozhkov et al., 2024), and DeepSeek-Coder (Guo et al., 2024a), while others like CodeQwen (Bai et al., 2023) and CodeLlama (Rozière et al., 2024) benefit from continued pre-training on additional code data. To enhance these models, supervised fine-tuning on instruction-response pairs has been employed (Luo et al., 2023; Wei et al., 2023; Zheng et al., 2024). Reinforcement learning techniques, like those in CodeRL (Le et al., 2022; Shojaee et al., 2023; ?), and reward models used in DeepSeek-Coder-V2 (DeepSeek-AI et al., 2024) also improve performance using test feedback.

Test Case Generation with Language Models Automated test case generation (Pacheco et al., 2007; Fraser & Arcuri, 2011; Panichella et al., 2015) 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 (Tufano et al., 2021; Li et al., 2022) or prompting them (Chen et al., 2023). Test cases also help clarify user intent, aligning model-generated programs with user requirements (Fakhoury et al., 2024; Endres et al., 2024).

The synergy between test cases and code generation Programming-by-examples (Gulwani, 2016) and test-driven programming (Perelman et al., 2014) 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 (Kulal et al., 2019; Chen et al., 2023; Zelikman et al., 2023). Recent methods, like CodeT (Chen et al., 2023) and Parsel (Zelikman et al., 2023), use test cases to reduce the search space during inference, while ALGO (Zhang et al., 2023) employs brute-force solutions as oracles to generate test outputs for competitive programming. Our approach, similar to Haluptzok et al. (2023), 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.

540 REFERENCES

- 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.
- 545 Mohammad Gheshlaghi Azar, Mark Rowland, Bilal Piot, Daniel Guo, Daniele Calandriello, Michal
 546 Valko, and Rémi Munos. A general theoretical paradigm to understand learning from human
 547 preferences, 2023.
- 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.
- 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.
- 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.
- 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.
- 569 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared 570 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, 571 Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, 572 Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios 573 Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, 574 Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, 575 Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, 576 Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob 577 McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating 578 large language models trained on code, 2021. 579
- 580 Dome Cleston. Sharegpt. https://github.com/domeccleston/sharegpt, 2023.
- 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.
- 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.
- 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.

594 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha 595 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, 596 Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston 597 Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, 598 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 600 Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, 601 Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip 602 Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme 603 Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, 604 Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, 605 Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, 606 Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu 607 Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph 608 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 610 Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas 611 Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, 612 Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, 613 Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, 614 Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur Celebi, Patrick Alrassy, Pengchuan 615 Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, 616 Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, 617 Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohit Girdhar, Rohit 618 Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, 619 Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia 620 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, 621 Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek 622 Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, 623 Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent 624 Gonguet, Virginie Do, Vish Vogeti, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, 625 Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, 626 Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen 627 Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe 628 Papakipos, Aaditya Singh, Aaron Grattafiori, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya 629 Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alex 630 Vaughan, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Anam Yunus, Andrei 631 Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita Saraf, Arkabandhu Chowdhury, Ashley 632 Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin 633 Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, 634 Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt 635 Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Changhan Wang, Changkyu Kim, Chao 636 Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Damon 637 Civin, Dana Beaty, Daniel Kreymer, Daniel Li, Danny Wyatt, David Adkins, David Xu, Davide 638 Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, 639 Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily 640 Hahn, Emily Wood, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix 641 Kreuk, Feng Tian, Firat Ozgenel, Francesco Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, 642 Govind Thattai, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen 644 Suk, Henry Aspegren, Hunter Goldman, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Irina-645 Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Japhet Asher, Jean-Baptiste 646 Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, 647 Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie,

648 Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik 649 Prasad, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly 650 Michelena, Keqian Li, Kun Huang, Kunal Chawla, Kushal Lakhotia, Kyle Huang, Lailin Chen, 651 Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, 652 Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, 653 Maxim Naumov, Maya Lathi, Meghan Keneally, Michael L. Seltzer, Michal Valko, Michelle 654 Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, 655 Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, 656 Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, 657 Nikolay Pavlovich Laptev, Ning Dong, Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia 658 Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro 659 Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, 660 Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, 661 Raghu Nayani, Rahul Mitra, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan 662 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 663 Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, 664 Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, 665 Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan 666 Govindaprasad, Sumit Gupta, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, 667 Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe 668 Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, 669 Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vítor Albiero, Vlad Ionescu, 670 Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, 671 Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaofang Wang, Xiaojian Wu, Xiaolan Wang, 672 Xide Xia, Xilun Wu, Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, 673 Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, 674 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. 675

- Madeline Endres, Sarah Fakhoury, Saikat Chakraborty, and Shuvendu K. Lahiri. Can large language
 models transform natural language intent into formal method postconditions?, 2024.
- Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. Kto: Model
 alignment as prospect theoretic optimization, 2024.
- 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.

685

686

687 688

689

690

- 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.
- 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*, pp. 416–419, 2011.
- Sumit Gulwani. Programming by examples and its applications in data wrangling. In Depend *able Software Systems Engineering*, 2016. URL https://api.semanticscholar.org/
 CorpusID:7866845.
- 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, 2024a.
- 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, 2024b. URL https://arxiv.org/abs/2402.04792.

 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.

- Binyuan Hui, Jian Yang, Zeyu Cui, Jiaxi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Keming Lu, Kai Dang, Yang Fan, Yichang Zhang, An Yang, Rui Men, Fei Huang, Bo Zheng, Yibo Miao, Shanghaoran Quan, Yunlong Feng, Xingzhang Ren, Xuancheng Ren, Jingren Zhou, and Junyang Lin. Qwen2.5-coder technical report, 2024. URL https://arxiv.org/abs/2409.12186.
- 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.
- 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 (eds.), 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.
- 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 (eds.), Advances in Neural Information Processing Systems, volume 35, pp. 21314–21328. Curran Associates, Inc., 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/8636419dealaa9fbd25fc4248e702da4-Paper-Conference.pdf.
- Raymond Li, Loubna Ben allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou, 727 Marc Marone, Christopher Akiki, Jia LI, Jenny Chim, Qian Liu, Evgenii Zheltonozhskii, Terry Yue 728 Zhuo, Thomas Wang, Olivier Dehaene, Joel Lamy-Poirier, Joao Monteiro, Nicolas Gontier, Ming-729 Ho Yee, Logesh Kumar Umapathi, Jian Zhu, Ben Lipkin, Muhtasham Oblokulov, Zhiruo Wang, 730 Rudra Murthy, Jason T Stillerman, Siva Sankalp Patel, Dmitry Abulkhanov, Marco Zocca, Manan 731 Dey, Zhihan Zhang, Urvashi Bhattacharyya, Wenhao Yu, Sasha Luccioni, Paulo Villegas, Fedor 732 Zhdanov, Tony Lee, Nadav Timor, Jennifer Ding, Claire S Schlesinger, Hailey Schoelkopf, Jan 733 Ebert, Tri Dao, Mayank Mishra, Alex Gu, Carolyn Jane Anderson, Brendan Dolan-Gavitt, Danish 734 Contractor, Siva Reddy, Daniel Fried, Dzmitry Bahdanau, Yacine Jernite, Carlos Muñoz Ferrandis, 735 Sean Hughes, Thomas Wolf, Arjun Guha, Leandro Von Werra, and Harm de Vries. Starcoder: may 736 the source be with you! Transactions on Machine Learning Research, 2023. ISSN 2835-8856. URL https://openreview.net/forum?id=KoFOg41haE. Reproducibility Certification. 737
- 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.
- 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*, 2023a. ISSN 2835-8856. URL https://openreview.net/forum?id=hjYmsV6nXZ.
- 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, 2023b.
- 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.
- 755 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,

756 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, 758 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, 760 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 761 Scholak, Sebastien Paquet, Jennifer Robinson, Carolyn Jane Anderson, Nicolas Chapados, Mostofa 762 Patwary, Nima Tajbakhsh, Yacine Jernite, Carlos Muñoz Ferrandis, Lingming Zhang, Sean Hughes, 763 Thomas Wolf, Arjun Guha, Leandro von Werra, and Harm de Vries. Starcoder 2 and the stack v2: 764 The next generation, 2024. 765

- 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.
- 770 MistralAI. Codestral-22b-v0.1, 2024. URL https://huggingface.co/mistralai/ Codestral-22B-v0.1. Accessed: 2024-09-28.
- Arindam Mitra, Hamed Khanpour, Corby Rosset, and Ahmed Awadallah. Orca-math: Unlocking the potential of slms in grade school math, 2024.
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni 775 Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor 776 Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, 777 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 780 Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, 781 Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, 782 Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, 783 Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty 784 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 785 Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua 786 Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike 787 Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon 788 Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne 789 Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo 790 Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, 791 Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, 793 Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy 794 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 798 Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, 799 Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo 800 Noh, Long Ouyang, Cullen O'Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, 801 Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, 802 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, 804 Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis 805 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 Sheppard, 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

810	Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya,
811	Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason
812	Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff,
813	Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu,
814	Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba,
815	Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang,
816	William Zhuk, and Barret Zoph. Gpt-4 technical report, 2024.

- 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)*, pp. 75–84. IEEE, 2007.
- 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.
- 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), pp. 1–10. IEEE, 2015.
- 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/.
- 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.
- 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 Ilama: Open foundation models for code, 2024.
- 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.
- 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.
- 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.
- 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.
- 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.
- Michele Tufano, Dawn Drain, Alexey Svyatkovskiy, Shao Kun Deng, and Neel Sundaresan. Unit test case generation with transformers and focal context, 2021.
- 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, 2024a.
- 863 Xingyao Wang, Hao Peng, Reyhaneh Jabbarvand, and Heng Ji. Leti: Learning to generate from textual interactions, 2024b.

- Yuxiang Wei, Zhe Wang, Jiawei Liu, Yifeng Ding, and Lingming Zhang. Magicoder: Source code is all you need, 2023.
- 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*, 2024a. URL
 https://openreview.net/forum?id=clAKcA6ry1.
- 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, 2024b. URL https://
 arxiv.org/abs/2409.02392.
- 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.

An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, 878 Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong 879 Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jianxin Yang, Jin Xu, 880 Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, 882 Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin 883 Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng 884 Ren, Xuejing Liu, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, 885 Zeyu Cui, Zhenru Zhang, Zhifang Guo, and Zhihao Fan. Qwen2 technical report, 2024a. URL https://arxiv.org/abs/2407.10671. 887

- 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, 2024b. URL https://arxiv.org/abs/2406.17312.
- 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.
- Eric Zelikman, Qian Huang, Gabriel Poesia, Noah D. Goodman, and Nick Haber. Parsel: Algorithmic reasoning with language models by composing decompositions, 2023.
- Kexun Zhang, Danqing Wang, Jingtao Xia, William Yang Wang, and Lei Li. Algo: Synthesizing
 algorithmic programs with llm-generated oracle verifiers, 2023.
- 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.
 - 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.
- 906 907

894

899

903

904

905

- 908
- 909 910

- 912
- 913
- 914
- 915
- 916 917

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 5×10^{-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 Li et al. (2022) which are equipped with test cases. Using more general datasets like OSS-INSTRUCT (Wei et al., 2023), 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

A.4 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.

Model	Туре	MBPP	MBPP+	HE	HE+	Avg. (Base)	Avg. (+)	Avg. (All)
	Baseline	54.4	45.6	35.4	29.9	44.9	37.8	41.3
STARCODER2-BASE	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
	Baseline	72.2	60.2	51.8	45.7	62.0	53.0	57.5
CODEQWEN-BASE	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
	Baseline	70.2	56.6	47.6	39.6	58.9	57.8	53.5
DS-CODER-BASE	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.

Model Families	Item	MBPP	MBPP+	HE	HE+
	Base	72.2	60.2	51.8	45.7
	OSS-Instruct	75.4	62.9	70.1	62.2
CODEQWEN-BASE	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
	Base	54.4	45.6	35.4	29.9
	OSS-Instruct	60.4	49.1	46.3	39.6
STARCODER2-BASE	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.

A.5 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.

Test Generator	Algorithm	MBPP	MBPP+	HE	HE+	Avg.	LeetCode	LCB
-	Baseine	77.7	67.2	83.5	78.7	76.8	33.9	23.2
CDT /	KTO	81.0	69.0	86.0	81.1	79.3	35.2	25.8
Ur I-4	DPO	81.2	70.2	86.0	81.1	79.6	36.7	25.8
LLAMAS 70P	KTO	79.4	69.9	84.8	80	78.5	36.7	24.5
LLAMAJ-/UD	DPO	79.4	66.2	84.1	79.3	77.3	36.1	24.5
LLAMA 2 405P	KTO	79.4	67.7	84.1	79.9	77.8	36.1	23.8
LLAMAJ-40JD	DPO	79.2	66.9	85.4	80.5	78.0	36.6	25.5
CDT 2 5 TUDDO	KTO	80.2	67.9	84.8	79.9	78.2	36.1	24.0
GF 1-3.3-10KB0	DPO	79.7	67.7	84.8	79.9	78.0	36.1	23.8
CPT/O MINI	KTO	81.2	69.2	85.4	80.5	79.1	36.7	24.0
OF 140-MINI	DPO	80.5	67.6	85.4	81.1	78.7	36.7	23.8
	KTO	79.7	67.2	85.4	81.7	78.5	36	24.0
CLAUDE-3-HAIKU	DPO	79.9	67.2	86	81.7	78.7	36	25.5

Table 11: Ablation of test case generator models. We used CODEQWEN-1.5-CHAT as the policy model. LCB stands for LiveCodeBench.

Model	Item	LeetCode	LiveCodeBench
QWEN-2.5-INSTRUCT-14B	Baseline	55.0	46.0
	PLUM-DPO	58.3	47.0
QWEN-2.5-CODER-14B	Baseline	58.3	32.2
	PLUM-DPO	61.7	35.0

Table 12: PLUM on more powerful policy models.

1051 A.6 ABLATION ON TEST CASE GENERATOR

To demonstrate the robustness of PLUM across different models as test case generators, and more specifically, to showcase its ability to boost the policy model's performance with more efficient options for test case generator thus proving its scalability, we conduct extensive experiments with multiple other test case generator models. We considered proprietary models with much more affordable API access and presumably less powerful than GPT-4 (GPT-3.5-Turbo, GPT4o-mini, Claude-3-Haiku). and Open-weight models (Llama 3.1 70B and 405B).

Table 11 displays the results of the ablation study.Consistent performance gains were observed acrossthe experiments.

Importantly, the use of more cost-efficient test case generators does not compromise PLUM's effectiveness. This demonstrates the scalability of our approach, enabling its practical application across a wide range of test generators and resource constraints.

A.7 IMPROVING STRONGER LANGUAGE MODELS

To validate the generalizability of the approach to more powerful post-trained models with larger parameter size and more sophisticated training, we conducted experiments with QWEN-2.5-INSTRUCT-14B Yang et al. (2024a) and QWEN-2.5-CODER-14B Hui et al. (2024) models. These models are fine-tuned from larger and stronger pre-trained models have undergone more sophisticated post-training including reinforcement learning and preference alignment.

As presented in Table 12, PLUM can further improve these models' performance especially on
 challenging programming benchmarks like LeetCode and LiveCodeBench. This not only validates
 the effectiveness of our approach, but highlights the potential of PLUM to be applied to complement
 other post-training techniques.

- A.8 GENERATION OF SYNTHETIC NEGATIVES
- 1079 We present the algorithm we used to generate synthetic negatives below in Algorithm A.8.

Alg	orithm 2 MutateCode Algorithm
Rec	uire: source code as a string mutation probability P
Ens	une: mutated code as a string
1.	Parse the source code into an AST: $tree \leftarrow ParseAST(source code)$
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
10:	Function visit If (node).
17:	Function VISIL_II($IIOUC$). With probability D , pagata the condition or swap if also blocks
10. 10.	Function visit For(node):
19. 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: <i>mutated_tree</i> \leftarrow Mutator().visit(<i>tree</i>)
26:	Perform syntactic validation: $is_valid \leftarrow SyntaxCheck(mutated_tree)$
27:	if <i>is_valid</i> = False then
28:	return original source code or error
29:	end if
20.	Convert the mutated AST back to code: $mutated_code \leftarrow ASTtoSource(mutated_tree)$
50.	