

# PREFERENCE OPTIMIZATION FOR REASONING WITH PSEUDO FEEDBACK

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Preference optimization techniques, such as Direct Preference Optimization (DPO), are frequently employed to enhance the reasoning capabilities of large language models (LLMs) in domains like mathematical reasoning and coding, typically following supervised fine-tuning. These methods rely on high-quality labels for reasoning tasks to generate preference pairs; however, the availability of reasoning datasets with human-verified labels is limited. In this study, we introduce a novel approach to generate pseudo feedback for reasoning tasks by framing the labeling of solutions to reason problems as an evaluation against associated *test cases*. We explore two forms of pseudo feedback based on test cases: one generated by frontier LLMs and the other by extending self-consistency to multi-test-case. We conduct experiments on both mathematical reasoning and coding tasks using pseudo feedback for preference optimization, and observe improvements across both tasks. Specifically, using Mathstral-7B as our base model, we improve MATH results from 58.3 to 68.6, surpassing both NuminaMath-72B and GPT-4-Turbo-1106-preview. In GSM8K and College Math, our scores increase from 85.6 to 90.3 and from 34.3 to 42.3, respectively. Building on Deepseek-coder-7B-v1.5, we achieve a score of 24.3 on LiveCodeBench (from 21.1), surpassing Claude-3-Haiku.

## 1 INTRODUCTION

Large language models (LLMs) have demonstrated exceptional capabilities in reasoning tasks such as math reasoning and coding (Roziere et al., 2023; Dubey et al., 2024; Guo et al., 2024). A *de facto* pipeline for enhancing the reasoning capabilities of LLMs involves further exposing them to reasoning specific data through continued pre-training or supervised fine-tuning (Roziere et al., 2023; Dubey et al., 2024; Yu et al., 2023; Tang et al., 2024; Dong et al., 2023), followed by preference learning techniques such as direct preference optimization (DPO; Rafailov et al. (2023)) or proximal policy optimization (PPO; Schulman et al. (2017)). Both DPO and PPO depend on reliable labels for reasoning problems to generate preference pairs and train reward models (Lightman et al., 2024; Uesato et al., 2022). Unfortunately, reasoning datasets with large-scale, human-verified labels remain limited, and scaling them through domain experts is becoming increasingly time-consuming and expensive, particularly as LLMs continue to evolve in capabilities (Burns et al., 2024; Bowman et al., 2022), which greatly limits the potential of preference learning methods such as DPO and PPO.

Scalable oversight (Bowman et al., 2022) demonstrates that the annotation effort of human experts can be significantly reduced with the assistance of non-expert LLMs. However, complete elimination of human annotation remains unattainable. Building on this, Khan et al. (2024a) further reduced labeling costs by incorporating a debating mechanism, though this approach is constrained to reasoning tasks with a finite answer space (e.g., multiple-choice questions). Other works have employed self-consistency-based answers or their variants as pseudo-labels to filter self-generated solutions (Huang et al., 2022; Yang et al., 2024c), but these methods struggle to generalize to reasoning tasks that lack explicit answer labels (e.g., coding).

To address these challenges, we frame the labeling of solutions to reasoning problems as the evaluation of these solutions against the *test cases* of the problems. For tasks with explicit answer labels (e.g., mathematical reasoning and multiple-choice questions), we treat them as cases with a single test pair, where the input is empty, and the output is the answer label. In contrast, for tasks without

054 explicit answer labels (e.g., coding), we consider them as problems with multiple test case pairs. A  
055 solution to a reasoning problem is deemed correct if and only if it passes all associated test cases.  
056 Sample solutions generated by an LLM for the same problem can be validated using the test case  
057 suite, with correct and incorrect solutions used to construct preference pairs for DPO training or to  
058 train a reward model for PPO. In this paper, we propose two types of pseudo feedback (i.e., pseudo  
059 test cases) for reasoning problems, both of which eliminate the need for human experts and can be  
060 applied at scale. First, we explore pseudo feedback from frontier LLMs, where we decompose the  
061 process of creating pseudo test cases into multiple steps to ensure that each step is manageable for  
062 frontier LLMs. Intuitively, if an LLM can pass test cases carefully curated by a stronger LLM, it  
063 is likely to provide a correct solution. Previous work Wang et al. (2022); Snell et al. (2024) has  
064 demonstrated that self-consistency improves the reasoning performance of LLMs. Based on this  
065 insight, we introduce a second form of pseudo feedback, utilizing self-consistency from our policy  
066 LLM, which is of vital importance when frontier LLMs are no longer available. Unlike the method  
067 in Wang et al. (2022), which is limited to single-test-case problems, our self-consistency feedback  
068 is designed to generalize to problems with multiple test cases. We also find that these two types of  
069 pseudo feedback complement each other and can be applied iteratively in a pipeline. We conducted  
070 experiments on both mathematical reasoning and coding using pseudo feedback for preference op-  
071 timization and we observe improvements across both tasks. Specifically, using Mathstral-7B as our  
072 base model, we improved our MATH results from 58.3 to 68.6, surpassing both NuminaMath-72B  
073 and GPT-4-Turbo-1106-preview. In GSM8K and College Math, our results increased from  
074 85.6 to 90.3 and from 34.3 to 42.3, respectively. Building on Deepseek-coder-7B-v1.5, we achieved  
a score of 24.3 on LiveCodeBench (from 21.1), surpassing Claude-3-Haiku.

075 In a nutshell, our contribution in this paper can be summarized as follows:

- 077 • We formulate the labeling of solutions to reasoning problems as the process of evaluating  
078 them against the associated *test cases*, which facilitates preference optimization.
- 079 • We explore two types of pseudo feedback based on *test cases*: one created from frontier  
080 LLMs and the other derived from generalized self-consistency w.r.t. multiple test cases.
- 081 • Experiments on mathematical reasoning and coding demonstrate the superiority of these  
082 two types of feedback. We also find they can be applied in a pipeline and iteratively to  
083 further improve the reasoning performance.

## 087 2 RELATED WORK

088 LLMs exhibit remarkable capabilities by tuning on high-quality data annotated by experts or more  
089 advanced models (Achiam et al., 2023). However, these external annotations can be costly, posing  
090 a challenge to further enhance model’s performance. Inspired by the natural evolution process of  
091 human intelligence, researchers explore self-evolution methods (Tao et al., 2024) that enable models  
092 to autonomously acquire, refine, and learn from their own knowledge. Some works (Wang et al.,  
093 2023b; Ding et al., 2024) reformulate the training objective to directly model performance improve-  
094 ment. Others tune the model with its own responses. They first filter the model’s outputs relying  
095 on ground truth labels (Zelikman et al., 2022; Wang et al., 2024b), expert annotations (Dubey et al.,  
096 2024), or more advanced models (Yang et al., 2024a; Kirchner et al., 2024), and then use the result-  
097 ing refined examples for supervised or contrastive learning (Chen et al., 2024b; Yuan et al., 2024).  
098 However, they still depend on external supervision and cannot extend to larger unlabeled datasets.  
099 Recent work (Huang et al., 2022) constructs pseudo labels via self-consistency, but the improvement  
100 is limited, possibly due to model collapse (Shumailov et al., 2023; Alemohammad et al., 2023).

102 For related work about mathematical reasoning (Wei et al., 2022; He-Yueya et al., 2023; Chen et al.,  
103 2021b; 2022; Lightman et al., 2024; Wang et al., 2024a; Jiao et al., 2024; Lai et al., 2024; Cobbe  
104 et al., 2021b; Li et al., 2022a; Weng et al., 2022; Yu et al., 2023; Luo et al., 2023; Mitra et al.,  
105 2024; Yue et al., 2023) and code generation (Guo et al., 2024; DeepSeek-AI et al., 2024; Nijkamp  
106 et al., 2023b;a; Zelikman et al., 2022; Li et al., 2023a; 2022b; Wei et al., 2024; Le et al., 2022; Liu  
107 et al., 2023; Dou et al., 2024; Weyssow et al., 2024), we will discuss them in Appendix F due to the  
limitation of space.

### 3 METHOD

In reasoning tasks such as mathematical reasoning and coding, the solution to a problem can be verified using a *standard* answer or a set of test cases. This property makes it possible to automatically create preference pairs for an LLM solving reasoning tasks and further improve reasoning capabilities of the LLM with preference optimization. However, annotating reasoning problems with answers or test cases manually is expensive and time consuming. As a result, this process is difficult to executed in large scale. Therefore, we propose PFPO (Pseudo-Feedback Preference Optimization), a method to automatically create *pseudo* answers or test cases to facilitate preference learning. In this section, we first introduce preference optimization for reasoning in Section 3.1 (assuming gold answers or test cases are available). Then we will go to details of PFPO, which creates pseudo answers or test cases.

#### 3.1 PREFERENCE OPTIMIZATION FOR REASONING

Suppose we have a set of reasoning problems  $x$  with their test cases  $T$ :  $\mathcal{D} = \{(x_i, T_i)\}_{i=1}^{|\mathcal{D}|}$ , where  $T = \{\langle i_1, o_1 \rangle, \langle i_2, o_2 \rangle, \dots, \langle i_{|T|}, o_{|T|} \rangle\}$  and  $\langle i_k, o_k \rangle$  is the input-output pair of a test case. Note that  $T$  is a generalized representation for either a collection of test cases or the gold answer for problem  $x$ . If  $x$  is a coding problem,  $T$  is a set of test cases to verify the correctness of the corresponding solution of  $x$ . While if  $x$  is one of the other reasoning problems such as mathematical reasoning or multi-choice science questions, there is only one test case in  $T = \{\langle i, o \rangle\}$  and the input  $i$  is empty. For example, “compute  $1 + 1$ ” is a math question with  $i = \emptyset$  as its test case input and  $o = 2$  as its test case output.

Given a reasoning problem  $x$  and its test cases  $T$ , we are ready to evaluate the correctness of a solution  $y$  produced by an LLM  $\pi_\theta$  as follows:

$$r = \frac{1}{|T|} \left( \sum_{k=1}^{|T|} \mathbb{1}(g(y, i_k) = o_k) \right) \quad (1)$$

where  $g(\cdot, \cdot)$  is a function to either execute the solution  $y$  or extract the answer from  $y$ . In the most strict form,  $y$  is a correct solution to problem  $x$  when  $r = 1$ . Otherwise (i.e.,  $r < 1$ ),  $y$  is an incorrect solution. Note that in mathematical reasoning, there is only one test case and  $r \in \{0, 1\}$ .

Note that given a problem  $x$  and its corresponding test cases  $T$ , the process of verifying an arbitrary solution  $y$  does not need any human labeling effort. We can construct preference pairs for an LLM automatically as follows. First, we use an LLM  $\pi_\theta$  to sample  $N$  solutions  $Y = \{y_1, y_2, \dots, y_N\}$  for problem  $x$  and obtain their verification results  $R = \{r_1, r_2, \dots, r_N\}$ . To further improve  $\pi_\theta$ , we can use PPO (Schulman et al., 2017) to optimize these feedback online or use DPO (Rafailov et al., 2023) to do preference optimization offline. In this work, we employ DPO due to its simplicity. Then, we create preference pairs from  $R$  and valid pairs  $(y_w, y_l)$  requires  $r_w = 1$  and  $r_l < 1$ .

$$\mathcal{P} = \{(y_w, y_l) | r_w = 1, r_l < 1, r_w \in R, r_l \in R\} \quad (2)$$

Given these valid preference pairs, We optimize our LLM  $\pi_\theta$  using the following objective:

$$\mathcal{L}_{\text{DPO}}(\pi_\theta; \pi_{\text{ref}}; \mathcal{D}) = -\mathbb{E}_{x \in \mathcal{D}, y_w, y_l \sim \pi_\theta(\cdot|x)} \left[ \log \sigma \left( \beta \log \frac{\pi_\theta(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \beta \log \frac{\pi_\theta(y_l|x)}{\pi_{\text{ref}}(y_l|x)} \right) \right] \quad (3)$$

where  $\pi_{\text{ref}}$  is the reference model before the DPO stage (usually it is the model of the supervised fine-tuning stage).  $\beta$  is a hyper-parameter to control the distance between  $\pi_\theta$  and  $\pi_{\text{ref}}$ .

#### 3.2 PSEUDO FEEDBACK PREFERENCE OPTIMIZATION FOR REASONING

In this section, we introduce how to obtain pseudo feedback for reasoning problems and the method of leverage them for preference learning.

##### 3.2.1 PSEUDO FEEDBACK FROM FRONTIER LLM

**Single-Test-Case Feedback** For single-test-case reasoning tasks such as mathematical reasoning, the input is *explicitly* given in the problem itself. Therefore, we do not need to create new test cases.

162  
163  
164  
165  
166  
167  
168  
169  
170  
171  
172  
173  
174  
175  
176  
177  
178  
179  
180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190  
191  
192  
193  
194  
195  
196  
197  
198  
199  
200  
201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215

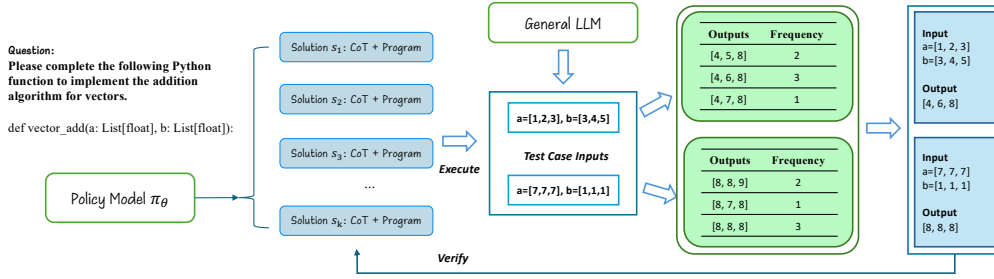


Figure 1: The process of employing self-consistency (*i.e.*, majority voting) to construct pseudo test cases for code generation problem. The outputs owing the highest frequency will be treated as pseudo outputs for verifying generated programs.

Given a problem  $x$ , we can use a frontier LLM to generate a solution  $\tilde{y}$  and extract its *pseudo* answer  $g(\tilde{y}, \cdot)$  as our pseudo feedback (see Equation 1). The solution  $y \sim \pi_\theta(\cdot|x)$  from our model is likely to be correct if  $g(y, \emptyset) = g(\tilde{y}, \emptyset)$ :

$$r = \mathbb{1}(g(y, \emptyset) = g(\tilde{y}, \emptyset)) \quad (4)$$

Since solutions of many reasoning datasets used for LLM supervised fine-tuning are created by frontier LLM (Taori et al., 2023; Tang et al., 2024), we can re-use the SFT datasets and extract the pseudo feedback as a *free lunch*.

**Multi-Test-Case Feedback** For multi-test-case reasoning tasks such as coding, test cases for coding problems are usually not available and manually label them are expensive. We choose to generate pseudo test cases by prompting frontier LLMs. There are three steps to generate test cases as shown in Figure 1:

- Step 1: Given a problem  $x$ , generate input test cases  $\mathcal{I} = \{\tilde{i}_1, \tilde{i}_2, \dots, \tilde{i}_K\}$  by prompting a *general*<sup>1</sup> LLM.
- Step 2: Generate pseudo (code) solutions  $\mathcal{Y} = \{y'_1, y'_2, \dots, y'_{|\mathcal{Y}|}\}$  for problem  $x$  using a frontier LLM.
- Step 3: Generate pseudo output test cases  $\mathcal{O} = \{o'_1, o'_2, \dots, o'_K\}$  using majority voting by executing all solutions in  $\mathcal{Y}$  for each input test case in  $\mathcal{I}$ .

The output test case  $o'_k$  corresponds to the input  $\tilde{i}_k$  is obtained as follows: after executing all pseudo solutions, we obtain a set of candidate pseudo output  $\mathcal{O}'_k = \{g(y'_1, \tilde{i}_k), g(y'_2, \tilde{i}_k), \dots, g(y'_{|\mathcal{Y}|}, \tilde{i}_k)\}$ . The output test case  $o'_k$  is the most frequent element in  $\mathcal{O}'_k$ :

$$o'_k = \arg \max_{o \in \mathcal{O}'_k} f(o) \quad (5)$$

where  $f(o) = |\{x \in \mathcal{O}'_k \mid x = o\}|$  is a frequency function that gives the number of times an element  $o$  appears in  $\mathcal{O}'_k$ . The resulting set of pseudo test cases is  $T' = \{\langle \tilde{i}_1, o'_1 \rangle, \langle \tilde{i}_2, o'_2 \rangle, \dots, \langle \tilde{i}_K, o'_K \rangle\}$ . At this point, we can verify arbitrary solution  $y$  to problem  $x$  as in Equation 1.

Note that we do not choose to generate both input and output test cases in a single step by prompting LLMs, as Gu et al. (2024) have pointed out that, generating the test case output based given input is a challenging task, which requires strong reasoning capabilities of LLMs. Also note that the single-test-case pseudo feedback described earlier is essentially equivalent to multi-test-case method feedback with number of test cases equals to one and the input test case is empty.

### 3.2.2 PSEUDO FEEDBACK FROM SELF-CONSISTENCY

Methods above leverage frontier LLMs to create pseudo feedback. We can alternatively create feedback from our own policy model  $\pi_\theta$  to facilitate self-improvement without external guidance. We

<sup>1</sup>Here we differentiate *general* LLM with the *frontier* one as generating only the inputs is much easier compared with solving the problem itself. Thus this process does not necessarily rely on SOTA LLMs.

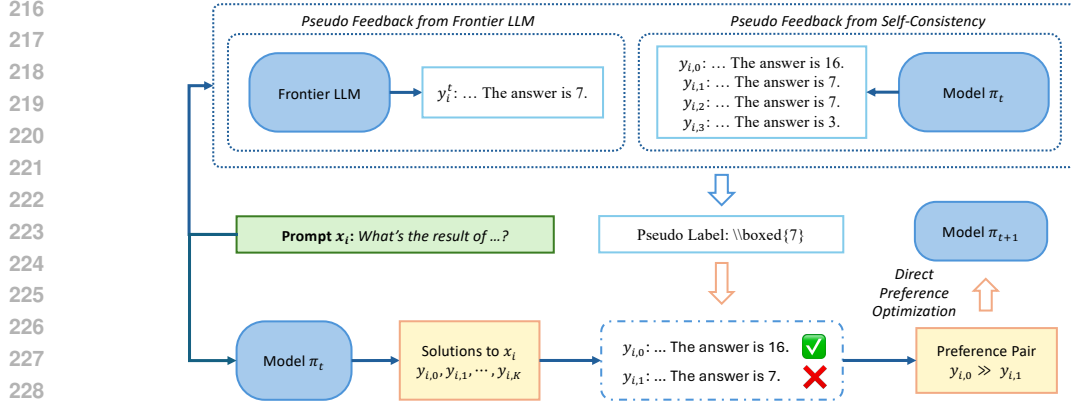


Figure 2: The overall training workflow of our method. For simplicity, we only show single step with outcome feedback for mathematical reasoning. Given an arbitrary prompt, we will sample multiple solutions from the current policy model and construct preference pairs according to pseudo feedback from frontier LLM or self-consistency. Finally, the constructed preference pair will be used to improve the policy model through DPO.

start from the method for the multi-test-case reasoning tasks, since the single-test-case counterpart is a special case of it. Specifically, we re-use the input test cases generated in Step 1 (Section 3.2.1). The main difference starts from Step 2. In the second step, we use our policy model  $\pi_\theta$  to sample pseudo solutions. The pseudo output in the third step is also based on executing all pseudo solutions from our policy model  $\pi_\theta$ . We can apply the same process to single-test-case reasoning tasks such as mathematical reasoning, which is equivalent to using majority voted answer from  $\pi_\theta$  samples as pseudo feedback. We can again use Equation 1 to verify the correctness of solutions from  $\pi_\theta$ .

### 3.2.3 PREFERENCE LEARNING UNDER PSEUDO FEEDBACK

Given the problem  $x$  and the pseudo feedback (i.e., test cases)  $T'$  we have created, the preference optimization process is as follows. We first sample  $N$  solutions  $Y = \{y_1, y_2, \dots, y_N\}$  to problem  $x$  from our policy  $\pi_\theta$ . We then obtain our verification results  $R = \{r_1, r_2, \dots, r_N\}$  using Equation 1 (i.e., executing all solutions on all test cases). We then move to create preference pairs using a different method as in Equation 2:

$$\mathcal{P}_o = \{(y_w, y_l) | r_w \geq \epsilon, r_w - r_l > \sigma, r_w \in R, r_l \in R\} \quad (6)$$

where  $\epsilon$  and  $\sigma$  are two hyper-parameters controlling the quality lower-bound of positive samples and the margin, respectively. Because our pseudo test cases may contains errors and if a solution  $y_k$  is required to pass all test cases, we may end up with no positive solutions for problem  $x$ . As a result, if a solution passes enough tests ( $r_w \geq \epsilon$ ) and significantly more tests than another solution ( $r_w - r_l > \sigma$ ), we treat them  $((y_w, y_l))$  as an eligible preference pair.

The above preference pairs in  $\mathcal{P}_o$  are based on outcome feedback of test cases. Recent studies (Wang et al., 2024a; Jiao et al., 2024) demonstrate that the outcome feedback can be used to estimate the expected returns of intermediate reasoning steps, which can help the model produce better reasoning trajectory. Motivated from this, we also construct the step-level process preference data. Following pDPO (Jiao et al., 2024), given the solution prefix  $\hat{y}$ , we employ the same policy model to sample  $M$  completions following  $\hat{y}$ , and treat the averaged outcome feedback  $\hat{r}$  of the completions as the expected returns of  $\hat{y}$ .

$$\hat{r} = \mathbb{E}_{y \sim \pi_\theta(\cdot | x, \hat{y})} r(\hat{y} \circ y) \quad (7)$$

where  $\circ$  is the concatenation operator. After that, the process preference data can be defined as:

$$\mathcal{P}_s = \{(\hat{y}_w, \hat{y}_l) | \hat{r}_w \geq \epsilon, \hat{r}_w - \hat{r}_l > \sigma, \hat{r}_w \in R, \hat{r}_l \in R\} \quad (8)$$

The final preference data we use is a combination of the outcome and process preference datasets  $\mathcal{P} = \mathcal{P}_o \cup \mathcal{P}_s$ . We use the DPO objective (Equation 3) to optimize the policy model  $\pi_\theta$ .

**Iterative Training** PFPO can be applied after the supervised fine-tuning (SFT) stage and we can train the policy model iteratively with both the feedback from frontier LLMs and from self-consistency. Imperially, we find applying LLM feedback first and then followed by self-consistency feedback rather than the opposite achieves better results. Figure 2 illustrates the process of single step.

## 4 EXPERIMENT

### 4.1 EXPERIMENTAL SETUP

**Prompt Collection** For mathematical reasoning, we followed Tang et al. (2024) to create 800K prompts under the help of GPT-4o. 500K of them are paired with one solution written by GPT-4o to construct pseudo feedback from frontier LLM. The 500K data is named as MathScale-500K, and the other prompts are called MathScale-300K. We also filtered out around 790K prompts from NuminaMath<sup>2</sup> by removing those that we cannot extract the predicted answers from or having appeared in the test set of MATH. For validation, we randomly sampled 2,000 question-solution pairs from the training set of MWPBench (Tang et al., 2024) after removing the questions from GSM8K (Cobbe et al., 2021a).

For code generation, we have collected the problems from the training set of APPs (Hendrycks et al., 2021a), Magicoder (Wei et al., 2024) and xCodeEval (Khan et al., 2024b), which contains 5,000, 9,000 and 6,400 questions, respectively. We remove all prompts where the test case inputs are failed to be synthesized. We randomly sampled 500 questions from the training set of APPs for validation. For APPs and xCodeEval, we use GPT-4o to generate the test case inputs. For Magicoder, we employ Mistral-Large-Instruct-2407<sup>3</sup> for test case inputs generation, because of large size of the original magicoder dataset. The detailed prompt can be found in Appendix C.1.

**Evaluation** For mathematical reasoning, the performance is evaluated on the test set of MATH (Hendrycks et al., 2021b), GSM8K (Cobbe et al., 2021a), and College Math (Tang et al., 2024) by Accuracy. For code generation, we evaluate the models on HumanEval (Chen et al., 2021a), MBPP (sanitized version) (Austin et al., 2021), APPs (Hendrycks et al., 2021a), and Live-CodeBench (Jain et al., 2024) by Pass@1. Without specific clarification, all evaluations are conducted using zero-shot prompting and greedy decoding.

For simplicity, we only highlight the main results our method. For more details include the detailed source of prompts, from which checkpoints are the models initialized, please refer to Appendix A. All hyper-parameters for different experiments can be found in Appendix B.

### 4.2 EXPERIMENTAL RESULTS

#### 4.2.1 MATHEMATICAL REASONING

We took two models for experiments: Llama-3.1-8B-base (Dubey et al., 2024) and Mathstral-7B-v0.1. We first conducted SFT on 500K MathScale data with GPT-4o annotation, followed by our method, with GPT-4o generated labels as pseudo feedback. As shown in Table 1, the pseudo feedback from GPT-4o can achieve consistent improvements on Llama-3.1-8b-base and Mathstral-7B. On MATH and College MATH, PFPO-LLM have made 1.2 and 4.1 averaged absolute improvements compared with Llama-3.1-8b w/ SFT and Mathstral-7B w/ SFT, respectively.

At the second phase, we started iterative pDPO training on unseen prompts with self-consistency-based pseudo feedback. For Llama-3.1-8b, we used the prompts from NuminaMath-790K to synthesize solutions and construct pseudo feedback via self-consistency. The prompts are divided into non-overlapped splits for iterative training. As shown in the table, by employing pseudo feedback, the models achieve continuous improvements across different iterations. Specifically, our method achieves the best results at Iteration 5. Compared with Llama-3.1-8b w/ PFPO-LLM Iter. 0, it achieves consistent improvements with 2.8 on MATH, 3.0 on GSM8K, as well 2.2 on Col-

<sup>2</sup><https://huggingface.co/datasets/AI-MO/NuminaMath-CoT>

<sup>3</sup><https://huggingface.co/mistralai/Mistral-Large-Instruct-2407>

Table 1: Overall results on mathematical reasoning benchmarks. PFPO-LLM refers to the training phase employing the pseudo feedback from frontier model (GPT-4o), while PFPO-Self indicates the phase using pseudo feedback constructed from self-generated solutions. NuminaMath-72B-CoT is built on Qwen2-72B by fine-tuning on NuminaMath. †: Results are from Chan et al. (2024). We employ an evaluation strategy similar to Yang et al. (2024b).

	MATH	GSM8K	College Math
GPT-4o-2024-0512	78.7	95.8	46.7
GPT-4-Turbo-2024-0409	72.8	94.8	44.2
GPT-4-Turbo-1106-preview†	64.3	—	—
GPT-4-0613	55.0	93.5	39.0
NuminaMath-72B-CoT (Beeching et al., 2024)	67.1	91.7	39.8
Llama-3.1-8B-Instruct (Dubey et al., 2024)	47.5	84.5	27.5
Llama-3.1-70B-Instruct (Dubey et al., 2024)	68.1	95.5	41.8
Llama-3.1-8B-base (Dubey et al., 2024)	20.3 (4-shot)	56.7 (8-shot)	20.1 (4-shot)
w/ SFT	53.8	85.1	34.6
w/ PFPO-LLM Iter. 0	55.0	86.6	35.8
w/ PFPO-Self Iter. 1	55.9	87.6	36.6
w/ PFPO-Self Iter. 2	56.6	88.9	37.0
w/ PFPO-Self Iter. 3	57.0	88.8	36.7
w/ PFPO-Self Iter. 4	57.4	89.1	37.6
w/ PFPO-Self Iter. 5	<b>57.8</b>	<b>89.6</b>	<b>38.0</b>
Mathstral-7B-v0.1 (Mistral AI Team, 2024b)	58.3	85.6	34.3
w/ SFT	61.4	87.3	38.4
w/ PFPO-LLM Iter. 0	66.7	90.0	41.3
w/ PFPO-Self Iter. 1	67.8	<b>90.8</b>	42.0
w/ PFPO-Self Iter. 2	<b>68.6</b>	90.3	42.2
w/ PFPO-Self Iter. 3	68.2	90.4	<b>42.3</b>

lege Math, revealing the potential of iterative preference optimization via pseudo feedback from self-consistency.

For Mathstral-7B, we use the prompts in MathScale-300K for iterative training with pseudo feedback, since we did not observe improvements on NuminaMath. The prompts across different iterations are the same. As shown in the table, Mathstral-7B w/ PFPO-Self Iter. 2 achieves 1.9 absolute improvements on MATH, compared with the SFT model. And it can outperform the stronger counterparts like NuminaMath-72B-CoT<sup>4</sup>, Llama-3.1-70B-Instruct, and GPT-4-Turbo-1106-preview, with only 7B parameters, which have demonstrated the effectiveness of pseudo feedback. Besides, we find the performance will saturate after several iterations. We discuss the possible reasons in Appendix A.2.

#### 4.2.2 CODE GENERATION

For code generation, we selected Deepseek-coder-7B-v1.5-Instruct (Guo et al., 2024) for experiments. We first use GPT-4o to generate 11 program solutions for each question in the APPs training set, and use the ground-truth test cases to remove those having failed tests. The left solutions are kept to first fine-tune the base model. The resulted model is referred as *w/ SFT (APPs)*.

**Direct Preference Optimization via Test Case Execution Feedback** As described in Section 3.2.2, we have constructed preference dataset via executing the generated programs over real or synthetic test cases. The evaluation results are shown in Table 2 and 3.

First, we aim to discuss the effectiveness of fully synthetic test cases, a topic that has not yet been extensively explored. We use *w/ DPO* and *w/ pDPO* to denote methods utilizing ground truth test cases to gather execution feedback, while PFPO-Self Iter. 0 (APPs) employs the same prompt set but simulates execution feedback using pseudo test cases. From the results presented, we observe that pseudo test cases outperform ground truth ones in almost all benchmarks, with the exception of HumanEval. In

Table 4: The averaged number of test cases of each problem in the training set of APPs.

Avg. No.	Original	Synthetic
Training	5.16	<b>9.95</b>

<sup>4</sup><https://huggingface.co/AI-MO/NuminaMath-72B-CoT>

Table 2: Overall results (Pass@1) on program generation benchmarks. PFPO-Self refers to our training from pseudo feedback method, and the content in the brackets afterwards indicates the source of prompts. Specifically, *M.C.* refers to the prompt set of Magicoder (Wei et al., 2024), and *xCode*. is the short for xCodeEval (Khan et al., 2024b). *Introductory*, *Interview*, and *Competition* indicate the three difficulty levels of APPs. w/ (p)DPO (APPs) refers to that the execution feedback is synthesized based on the groundtruth test cases annotated in APPs training set.

	APPs				HumanEval	MBPP
	Overall	Introductory	Interview	Competition		
GPT-4-0613	35.1	61.8	34.4	10.6	87.8	82.1
GPT-4o-2024-0513	34.0	56.6	32.2	16.7	93.3	87.2
Llama-3.1-8B-Instruct (Dubey et al., 2024)	11.5	29.4	8.5	2.7	72.6	71.2
Llama-3.1-70B-Instruct (Dubey et al., 2024)	24.9	51.8	21.3	9.1	80.5	83.3
Codestral-22B-V0.1 (Mistral AI Team, 2024a)	20.3	45.2	16.9	5.8	81.1	78.2
CodeQwen1.5-7B-chat (Qwen Team, 2024)	8.6	24.1	16.8	2.0	85.6	80.5
Qwen2.5-Coder-7B-Instruct (Hui et al., 2024)	15.7	37.3	12.3	4.1	85.4	86.0
Deepseek-coder-33B-Instruct (Guo et al., 2024)	18.4	44.2	14.5	4.4	77.4	79.0
Deepseek-coder-v1.5-Instruct	14.3	35.7	10.8	3.2	75.6	73.9
w/ SFT (APPs)	15.4	37.8	11.6	4.1	72.0	72.8
w/ DPO (APPs)	16.3	36.2	13.3	5.3	74.4	74.3
w/ pDPO (APPs)	16.9	37.3	13.8	6.1	73.8	73.2
<b>w/ PFPO-LLM Iter. 0 (APPs)</b>	<b>17.9</b>	<b>38.3</b>	<b>14.7</b>	<b>7.1</b>	<b>73.8</b>	<b>75.9</b>
w/ PFPO-Self Iter. 0 (APPs)	17.4	37.5	14.8	5.4	73.2	75.1
w/ PFPO-Self Iter. 1 (APPs & M.C.)	18.0	39.2	14.9	6.2	<b>79.3</b>	<b>75.5</b>
w/ PFPO-Self Iter. 2 (APPs & M.C. & xCode.)	<b>19.1</b>	<b>40.9</b>	<b>15.9</b>	<b>6.9</b>	73.8	75.1

Table 3: Overall results on LiveCodeBench. We follow the recommended setting by sampling 10 solutions for each problem with temperature as 0.2, and estimating the Pass@1 results. The cutoff date of the test questions is from **2023-09-01** to **2024-09-01**. All results except those of our models are referenced from the official leaderboard (<https://livecodebench.github.io/>).

	Overall	Easy	Medium	Hard
Claude-3.5-Sonnet	51.3	87.2	45.3	11.0
Claude-3-Sonnet	26.9	67.2	7.3	1.4
Claude-3-Haiku	24.0	61.3	5.5	0.9
GPT-3.5-Turbo-0125	24.0	55.0	11.6	0.3
Llama-3.1-70B-Instruct (Dubey et al., 2024)	31.8	67.9	17.3	4.1
Llama-3-70B-Instruct (Dubey et al., 2024)	27.4	59.4	15.6	1.3
CodeQwen1.5-7B-Chat (Qwen Team, 2024)	16.8	35.9	10.9	0.3
DeepSeekCoder-V2-236B (DeepSeek-AI et al., 2024)	41.9	79.9	32.0	4.9
Deepseek-Coder-33B-Instruct (Guo et al., 2024)	23.4	56.1	8.6	0.9
Deepseek-coder-7B-v1.5-Instruct	21.1	51.3	7.4	0.2
w/ SFT (APPs)	22.9	53.0	10.6	0.2
w/ DPO (APPs)	22.9	53.7	9.4	1.0
w/ pDPO (APPs)	22.9	55.0	8.1	1.3
<b>w/ PFPO-LLM Iter. 0 (APPs)</b>	<b>24.0</b>	<b>56.8</b>	<b>9.3</b>	<b>1.4</b>
w/ PFPO-Self Iter. 0 (APPs)	23.4	54.2	10.3	0.7
w/ PFPO-Self Iter. 1 (APPs & M.C.)	23.7	55.8	9.5	1.1
w/ PFPO-Self Iter. 2 (APPs & M.C. & xCode)	<b>24.3</b>	<b>56.8</b>	<b>9.8</b>	<b>1.6</b>

particular, PFPO-Self Iter. 0 leads both 0.5 absolute improvements on APPs and LiveCodeBench compared with the groundtruth pDPO. This improvement is attributed to the potential of increasing the number of synthetic test cases, thereby reducing the false positive rate associated with missing corner cases. As shown in Table 4, the questions in the APPs training set contain only 5.16 test cases on average, whereas we can ensure that each problem has approximately 10 test cases.

Besides, by comparing w/ PFPO-LLM Iter. 0 (APPs) and w/ PFPO-Self Iter. 0 (APPs), we find that the pseudo feedback generated by frontier LLM also demonstrates better results on code generation. This is also reflected through the quality of synthetic test cases, and we have a deep analysis in Appendix A.4.



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

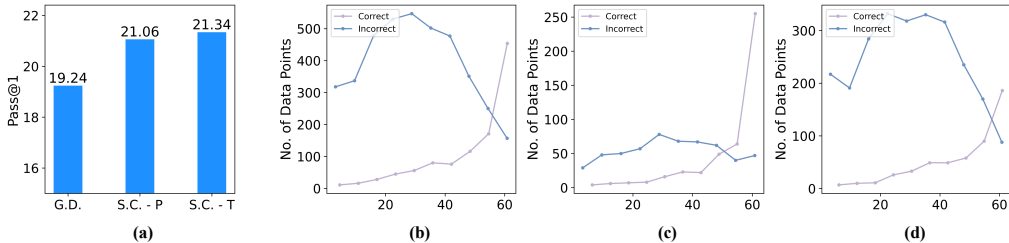


Figure 3: (a) The Pass@1 performance of the test set of APPs of our model variant, w/ PFPO-Self Iter. 2 (APPs & M.C. & xCode, DPO) under different strategies. *G.D.* indicates greedy decoding, *S.C.* is the short for self-consistency. We took the groundtruth test case inputs to execute the programs sampled from our policy model. For each group of test case inputs, we employ self-consistency across different programs to determine the pseudo outputs. For *S.C. - P*, we select one of the programs matching the pseudo outputs for evaluation. For *S.C. - T*, we simply check whether the pseudo outputs are consistent with the ground-truth outputs, which indeed highlights an upper bound of test case self-consistency, since sometimes we fail to select at least one solution matching all the pseudo outputs. (b) - (d) We measure the top-1 frequency among the outputs and average them by the number of test cases. The *x*-axis indicates the averaged top-1 frequency, and the *y*-axis demonstrates the corresponding amount of data samples. The three figures show the relationships on (b) all questions, (c) introductory-level questions, and (d) interview-level questions, respectively.

**Iterative Training** Afterwards, we attempted to introduce more training data in an iterative manner to see if the pseudo feedback can continuously enhance the performance. As shown in Table 2 and 3, after three rounds of training, PFPO-Self Iter. 2 has made 3.7 and 1.4 absolute improvements on APPs and LiveCodeBench, respectively. We also observed more significant improvements on the Easy and Hard level of LiveCodeBench, *i.e.*, 3.8 and 1.4 absolute points, respectively. Besides, on HumanEval and MBPP, we did not observe extra improvements after introducing xCodeEval. This is probably because xCodeEval contains only standard input/output problems, which could lead to distribution shift. The small scale of test cases and problems in HumanEval and MBPP also make the evaluation not really stable.

### 4.3 RELATIONSHIP BETWEEN SELF-CONSISTENCY AND RELIABLE PSEUDO FEEDBACK

In this section, we will explore the impact of test-case-based self-consistency during inference for coding. We sample 64 solutions for each problem in the APPs test set and apply self-consistency over the provided test cases to select the final programs. Two strategies are developed within this process: **Self-Consistency for Programs (S.C. - P)** and **Self-Consistency for Test Cases (S.C. - T)**. Following Section 3.2.2, we execute the sampled programs **on the test case inputs annotated by the dataset** and determine the pseudo outputs for each group of inputs through majority voting.

First of all, we hope to evaluate the correctness of the self-consistency based pseudo outputs by checking the consistency with the ground truth outputs. If this is true, we treat this question is passed. This strategy is called S.C. - T. After that, we can check if there exists at least one program among the sampled solutions such that its outputs can match all the pseudo outputs, which is called S.C. - P. In other words, for S.C. - T, we assume that there is such a program among the candidates that can obtain the expected pseudo outputs, so we only care about if the pseudo outputs are as expected. For S.C. - P, we also consider if such a program is included in the sampled solutions. To this end, S.C. - T can be considered the upper bound of self-consistency, as there are cases where no solution matches the pseudo outputs. As illustrated in Figure 3 (a), the test-case-based self-consistency also makes significant improvements during test time. Besides, by continuously sampling solutions against the pseudo outputs, the pass rate would be further improved.

Besides, we also want to explore the reliance of self-consistency based pseudo feedback. From Figure 3 (b) to (d), we display the distribution of top-1 output frequency averaged by the number of test cases, *i.e.*, the confidence of the policy model, on overall, introductory-level, and interview-level problems. We find that, self-consistency over test cases provide reliable pseudo feedback when different programs achieve high consistency rate. Although the margin between positive and negative solutions is reduced with the increase of difficulty, by continuously improving the sampling budget,

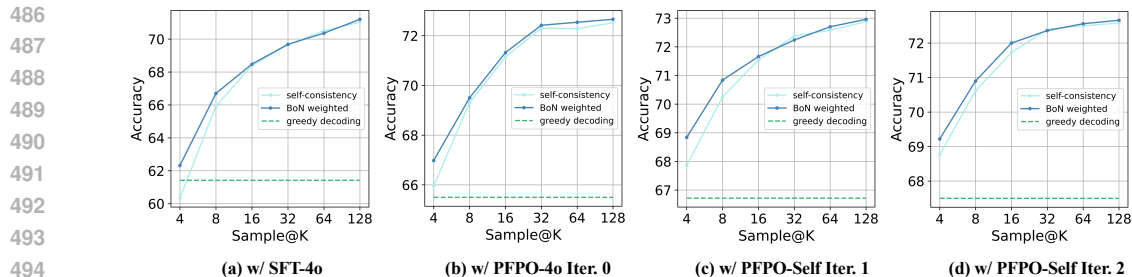


Figure 4: Performance comparison of four Mathstral-7B variants by scaling inference

the uncertainty can also be reduced. This may help to both control the difficulty of training set, and lower down the possible inaccurate feedback. On mathematical reasoning, similar conclusion still holds. Figure 5 demonstrates the accumulated ratio of corrected predictions with the development of self-consistency-based prediction ratio among all candidate answers.

#### 4.4 SCALING INFERENCE AND REWARD MODELING FROM PSEUDO FEEDBACK

In this section, we increase the inference budget (Snell et al., 2024) to assess how it further improves mathematical reasoning. Specifically, we explore two strategies: *self-consistency* and *best-of-N weighted* (Li et al., 2023b). The reward model used for *best-of-N weighted* is optimized on the same training set of Mathstral-7B w/ PFPO-Self Iter. 1, thus also benefiting from the pseudo feedback derived from self-consistency.

The results in Figure 4 indicate that: (i) Scaling inference makes significant improvements, and employing an extra reward model to score responses can bring more benefits, especially when using smaller sampling budget. (ii) Pseudo feedback can also enhance reward modeling. (iii) As highlighted by Snell et al. (2024), combining pseudo reward models with other inference techniques (e.g., weighted best-of-N (Li et al., 2023b), look-ahead search, and step-level beam search) may improve performance. We leave the relevant exploration as future work.

## 5 CONCLUSION

In this paper, we demonstrated the potential of synthesizing pseudo-feedback from both frontier LLMs and the model itself. By incorporating feedback from GPT-4o, we improved Llama-3.1-8b-base-SFT’s performance on the MATH from 53.8 to 55.0 and enhanced Mathstral-7B’s performance from 61.4 to 66.7. Furthermore, by leveraging self-generated feedback based on self-consistency, we increased Llama-3.1-8b’s performance to 57.8 and Mathstral-7B’s to 68.6. For code generation, we introduced a novel test case-level self-consistency strategy. By employing fully self-constructed pseudo test cases, we boosted the SFT model’s performance from 15.4 to 19.1 on APPs and from 22.9 to 24.3 on LiveCodeBench. Additionally, we analyzed the relationship between self-consistency and reliance on synthetic pseudo labels, offering insights for future research. For future work, we hope to combine the pseudo reward models with scaling inference techniques to seek more improvements on more challenge scenarios.

## REFERENCES

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.

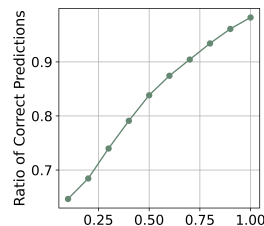


Figure 5: The accumulated ratio of correct predictions over all questions. The  $x$ -axis denotes the ratio of the top-1 frequency among all predicted candidates, *i.e.*, the confidence for self-consistency. Each point in the figure indicating how many predictions are correct when the confidence is lower than the  $x$ -value.

- 540 Sina Alemohammad, Josue Casco-Rodriguez, Lorenzo Luzi, Ahmed Imtiaz Humayun, Hossein  
541 Babaei, Daniel LeJeune, Ali Siahkoochi, and Richard G Baraniuk. Self-consuming generative  
542 models go mad. *arXiv preprint arXiv:2307.01850*, 2023.
- 543  
544 Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan,  
545 Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language  
546 models. *arXiv preprint arXiv:2108.07732*, 2021.
- 547 Edward Beeching, Shengyi Costa Huang, Albert Jiang, Jia Li, Benjamin Lipkin, Zihan Qina, Kashif  
548 Rasul, Ziju Shen, Roman Soletskyi, and Lewis Tunstall. Numinamath 72b cot. [https://](https://huggingface.co/AI-MO/NuminaMath-72B-CoT)  
549 [huggingface.co/AI-MO/NuminaMath-72B-CoT](https://huggingface.co/AI-MO/NuminaMath-72B-CoT), 2024.
- 550  
551 Samuel R. Bowman, Jeeyoon Hyun, Ethan Perez, Edwin Chen, Craig Pettit, Scott Heiner, Kamile  
552 Lukosiute, Amanda Askeell, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron  
553 McKinnon, Christopher Olah, Daniela Amodei, Dario Amodei, Dawn Drain, Dustin Li, Eli Tran-  
554 Johnson, Jackson Kernion, Jamie Kerr, Jared Mueller, Jeffrey Ladish, Joshua Landau, Kamal  
555 Ndousse, Liane Lovitt, Nelson Elhage, Nicholas Schiefer, Nicholas Joseph, Noemí Mercado,  
556 Nova DasSarma, Robin Larson, Sam McCandlish, Sandipan Kundu, Scott Johnston, Shauna  
557 Kravec, Sheer El Showk, Stanislav Fort, Timothy Telleen-Lawton, Tom Brown, Tom Henighan,  
558 Tristan Hume, Yuntao Bai, Zac Hatfield-Dodds, Ben Mann, and Jared Kaplan. Measuring  
559 progress on scalable oversight for large language models. *arXiv preprint arXiv:2211.03540*, 2022.
- 560 Collin Burns, Pavel Izmailov, Jan Hendrik Kirchner, Bowen Baker, Leo Gao, Leopold Aschenbren-  
561 ner, Yining Chen, Adrien Ecoffet, Manas Joglekar, Jan Leike, Ilya Sutskever, and Jeffrey Wu.  
562 Weak-to-strong generalization: Eliciting strong capabilities with weak supervision. In *Forty-first*  
563 *International Conference on Machine Learning*. OpenReview.net, 2024.
- 564 Xin Chan, Xiaoyang Wang, Dian Yu, Haitao Mi, and Dong Yu. Scaling synthetic data creation with  
565 1,000,000,000 personas. *CoRR*, abs/2406.20094, 2024.
- 566  
567 Changyu Chen, Zichen Liu, Chao Du, Tianyu Pang, Qian Liu, Arunesh Sinha, Pradeep Varakan-  
568 tham, and Min Lin. Bootstrapping language models with DPO implicit rewards. *CoRR*,  
569 abs/2406.09760, 2024a.
- 570 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared  
571 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri,  
572 Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan,  
573 Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian,  
574 Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fo-  
575 tios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex  
576 Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders,  
577 Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec  
578 Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob Mc-  
579 Grew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large  
580 language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021a.
- 581 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared  
582 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large  
583 language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021b.
- 584 Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W Cohen. Program of thoughts prompt-  
585 ing: Disentangling computation from reasoning for numerical reasoning tasks. *arXiv preprint*  
586 *arXiv:2211.12588*, 2022.
- 587  
588 Zixiang Chen, Yihe Deng, Huizhuo Yuan, Kaixuan Ji, and Quanquan Gu. Self-play fine-tuning  
589 converts weak language models to strong language models. *arXiv preprint arXiv:2401.01335*,  
590 2024b.
- 591 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,  
592 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John  
593 Schulman. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*,  
2021a.

- 594 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,  
595 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to  
596 solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021b.  
597
- 598 DeepSeek-AI, Qihao Zhu, Daya Guo, Zhihong Shao, Dejian Yang, Peiyi Wang, Runxin Xu, Y. Wu,  
599 Yukun Li, Huazuo Gao, Shirong Ma, Wangding Zeng, Xiao Bi, Zihui Gu, Hanwei Xu, Damai  
600 Dai, Kai Dong, Liyue Zhang, Yishi Piao, Zhibin Gou, Zhenda Xie, Zhewen Hao, Bingxuan Wang,  
601 Junxiao Song, Deli Chen, Xin Xie, Kang Guan, Yuxiang You, Aixin Liu, Qiushi Du, Wenjun Gao,  
602 Xuan Lu, Qinyu Chen, Yaohui Wang, Chengqi Deng, Jiashi Li, Chenggang Zhao, Chong Ruan,  
603 Fuli Luo, and Wenfeng Liang. Deepseek-coder-v2: Breaking the barrier of closed-source models  
604 in code intelligence. *CoRR*, abs/2406.11931, 2024.
- 605 Mucong Ding, Souradip Chakraborty, Vibhu Agrawal, Zora Che, Alec Koppel, Mengdi Wang, Am-  
606 rit Bedi, and Furong Huang. Sail: Self-improving efficient online alignment of large language  
607 models. *arXiv preprint arXiv:2406.15567*, 2024.  
608
- 609 Hanze Dong, Wei Xiong, Deepanshu Goyal, Yihan Zhang, Winnie Chow, Rui Pan, Shizhe Diao,  
610 Jipeng Zhang, Kashun Shum, and Tong Zhang. RAFT: reward ranked finetuning for generative  
611 foundation model alignment. *TMLR*, 2023, 2023.
- 612 Shihan Dou, Yan Liu, Haoxiang Jia, Limao Xiong, Enyu Zhou, Wei Shen, Junjie Shan, Caishuang  
613 Huang, Xiao Wang, Xiaoran Fan, Zhiheng Xi, Yuhao Zhou, Tao Ji, Rui Zheng, Qi Zhang, Xuan-  
614 jing Huang, and Tao Gui. StepCoder: Improve code generation with reinforcement learning from  
615 compiler feedback. *CoRR*, abs/2402.01391, 2024.  
616
- 617 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha  
618 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony  
619 Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark,  
620 Arun Rao, Aston Zhang, Aurélien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Rozière,  
621 Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris  
622 Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong,  
623 Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny  
624 Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino,  
625 Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael  
626 Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson,  
627 Graeme Nail, Grégoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar,  
628 Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel M. Kloumann, Ishan Misra,  
629 Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar,  
630 Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng  
631 Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park,  
632 Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya  
633 Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, and et al. The llama 3 herd of  
634 models. *CoRR*, abs/2407.21783, 2024.
- 635 Alex Gu, Baptiste Rozière, Hugh James Leather, Armando Solar-Lezama, Gabriel Synnaeve, and  
636 Sida Wang. CruxEval: A benchmark for code reasoning, understanding and execution. In *Forty-  
637 first International Conference on Machine Learning, ICML*. OpenReview.net, 2024.
- 638 Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao  
639 Bi, Y. Wu, Y. K. Li, Fuli Luo, Yingfei Xiong, and Wenfeng Liang. Deepseek-coder: When the  
640 large language model meets programming - the rise of code intelligence. *CoRR*, abs/2401.14196,  
641 2024.
- 642 Joy He-Yueya, Gabriel Poesia, Rose E Wang, and Noah D Goodman. Solving math word problems  
643 by combining language models with symbolic solvers. *arXiv preprint arXiv:2304.09102*, 2023.  
644
- 645 Dan Hendrycks, Steven Basart, Saurav Kadavath, Mantas Mazeika, Akul Arora, Ethan Guo, Collin  
646 Burns, Samir Puranik, Horace He, Dawn Song, and Jacob Steinhardt. Measuring coding challenge  
647 competence with APPS. In *Proceedings of the Neural Information Processing Systems Track on  
Datasets and Benchmarks*, 2021a.

- 648 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,  
649 and Jacob Steinhardt. Measuring mathematical problem solving with the MATH dataset. In  
650 *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*,  
651 2021b.
- 652 Jiaxin Huang, Shixiang Shane Gu, Le Hou, Yuexin Wu, Xuezhi Wang, Hongkun Yu, and Jiawei  
653 Han. Large language models can self-improve. *arXiv preprint arXiv:2210.11610*, 2022.
- 654 Binyuan Hui, Jian Yang, Zeyu Cui, Jiayi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang,  
655 Bowen Yu, Kai Dang, et al. Qwen2. 5-coder technical report. *arXiv preprint arXiv:2409.12186*,  
656 2024.
- 657 Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando  
658 Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free  
659 evaluation of large language models for code. *CoRR*, abs/2403.07974, 2024.
- 660 Fangkai Jiao, Chengwei Qin, Zhengyuan Liu, Nancy F. Chen, and Shafiq Joty. Learning planning-  
661 based reasoning by trajectories collection and process reward synthesizing. *arXiv preprint*  
662 *arXiv:2402.00658*, 2024.
- 663 Akbir Khan, John Hughes, Dan Valentine, Laura Ruis, Kshitij Sachan, Ansh Radhakrishnan, Ed-  
664 ward Grefenstette, Samuel R. Bowman, Tim Rocktäschel, and Ethan Perez. Debating with  
665 more persuasive LLMs leads to more truthful answers. In Ruslan Salakhutdinov, Zico Kolter,  
666 Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.),  
667 *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Pro-*  
668 *ceedings of Machine Learning Research*, pp. 23662–23733. PMLR, 21–27 Jul 2024a. URL  
669 <https://proceedings.mlr.press/v235/khan24a.html>.
- 670 Mohammad Abdullah Matin Khan, M. Saiful Bari, Xuan Do Long, Weishi Wang, Md. Rizwan  
671 Parvez, and Shafiq Joty. Xcodeeval: An execution-based large scale multilingual multitask bench-  
672 mark for code understanding, generation, translation and retrieval. In *Proceedings of the 62nd*  
673 *Annual Meeting of the Association for Computational Linguistics*, pp. 6766–6805. Association  
674 for Computational Linguistics, 2024b.
- 675 Jan Hendrik Kirchner, Yining Chen, Harri Edwards, Jan Leike, Nat McAleese, and Yuri Burda.  
676 Prover-verifier games improve legibility of llm outputs. *arXiv preprint arXiv:2407.13692*, 2024.
- 677 Xin Lai, Zhuotao Tian, Yukang Chen, Senqiao Yang, Xiangru Peng, and Jiaya Jia. Step-dpo: Step-  
678 wise preference optimization for long-chain reasoning of llms. *CoRR*, abs/2406.18629, 2024.
- 679 Hung Le, Yue Wang, Akhilesh Deepak Gotmare, Silvio Savarese, and Steven Chu-Hong Hoi.  
680 Coderl: Mastering code generation through pretrained models and deep reinforcement learning. In  
681 *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Informa-*  
682 *tion Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December*  
683 *9, 2022*, 2022.
- 684 Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao  
685 Mou, Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, Qian Liu, Evgenii Zheltonozh-  
686 skii, Terry Yue Zhuo, Thomas Wang, Olivier Dehaene, Mishig Davaadorj, Joel Lamy-Poirier,  
687 João Monteiro, Oleh Shliazhko, Nicolas Gontier, Nicholas Meade, Armel Zebaze, Ming-Ho Yee,  
688 Logesh Kumar Umabathi, Jian Zhu, Benjamin Lipkin, Muhtasham Oblokulov, Zhiruo Wang,  
689 Rudra Murthy V, Jason T. Stillerman, Siva Sankalp Patel, Dmitry Abulkhanov, Marco Zocca,  
690 Manan Dey, Zhihan Zhang, Nour Fahmy, Urvashi Bhattacharyya, Wenhao Yu, Swayam Singh,  
691 Sasha Luccioni, Paulo Villegas, Maxim Kunakov, Fedor Zhdanov, Manuel Romero, Tony Lee,  
692 Nadav Timor, Jennifer Ding, Claire Schlesinger, Hailey Schoelkopf, Jan Ebert, Tri Dao, Mayank  
693 Mishra, Alex Gu, Jennifer Robinson, Carolyn Jane Anderson, Brendan Dolan-Gavitt, Danish  
694 Contractor, Siva Reddy, Daniel Fried, Dzmitry Bahdanau, Yacine Jernite, Carlos Muñoz Ferran-  
695 dis, Sean Hughes, Thomas Wolf, Arjun Guha, Leandro von Werra, and Harm de Vries. Starcoder:  
696 may the source be with you! *Transactions on Machine Learning Research*, 2023, 2023a.
- 697 Yifei Li, Zeqi Lin, Shizhuo Zhang, Qiang Fu, Bei Chen, Jian-Guang Lou, and Weizhu Chen. On the  
698 advance of making language models better reasoners. *arXiv preprint arXiv:2206.02336*, 2022a.
- 699  
700  
701

- 702 Yifei Li, Zeqi Lin, Shizhuo Zhang, Qiang Fu, Bei Chen, Jian-Guang Lou, and Weizhu Chen. Making  
703 language models better reasoners with step-aware verifier. In *Proceedings of the 61st Annual*  
704 *Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 5315–  
705 5333. Association for Computational Linguistics, July 2023b.
- 706 Yujia Li, David H. Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond,  
707 Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, Thomas Hubert, Peter Choy, Cy-  
708 prien de Masson d’Autume, Igor Babuschkin, Xinyun Chen, Po-Sen Huang, Johannes Welbl,  
709 Sven Gowal, Alexey Cherepanov, James Molloy, Daniel J. Mankowitz, Esme Sutherland Robson,  
710 Pushmeet Kohli, Nando de Freitas, Koray Kavukcuoglu, and Oriol Vinyals. Competition-level  
711 code generation with alphacode. *CoRR*, abs/2203.07814, 2022b.
- 712 Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan  
713 Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. In *The Twelfth*  
714 *International Conference on Learning Representations*. OpenReview.net, 2024.
- 715 Jiatae Liu, Yiqin Zhu, Kaiwen Xiao, Qiang Fu, Xiao Han, Wei Yang, and Deheng Ye. RLTF: rein-  
716 forcement learning from unit test feedback. *Trans. Mach. Learn. Res.*, 2023, 2023.
- 717 Haipeng Luo, Qingfeng Sun, Can Xu, Pu Zhao, Jianguang Lou, Chongyang Tao, Xiubo Geng, Qing-  
718 wei Lin, Shifeng Chen, and Dongmei Zhang. Wizardmath: Empowering mathematical reasoning  
719 for large language models via reinforced evol-instruct. *arXiv preprint arXiv:2308.09583*, 2023.
- 720 Mistral AI Team. Codestral. [https://huggingface.co/mistralai/](https://huggingface.co/mistralai/Codestral-22B-v0.1)  
721 [Codestral-22B-v0.1](https://huggingface.co/mistralai/Codestral-22B-v0.1), 2024a.
- 722 Mistral AI Team. Mathstral. [https://huggingface.co/mistralai/](https://huggingface.co/mistralai/Mathstral-7B-v0.1)  
723 [Mathstral-7B-v0.1](https://huggingface.co/mistralai/Mathstral-7B-v0.1), 2024b.
- 724 Arindam Mitra, Hamed Khanpour, Corby Rosset, and Ahmed Awadallah. Orca-math: Unlocking  
725 the potential of slms in grade school math. *arXiv preprint arXiv:2402.14830*, 2024.
- 726 Erik Nijkamp, Hiroaki Hayashi, Caiming Xiong, Silvio Savarese, and Yingbo Zhou. Codegen2:  
727 Lessons for training llms on programming and natural languages. *ICLR*, 2023a.
- 728 Erik Nijkamp, Bo Pang, Hiroaki Hayashi, Lifu Tu, Huan Wang, Yingbo Zhou, Silvio Savarese,  
729 and Caiming Xiong. Codegen: An open large language model for code with multi-turn program  
730 synthesis. *ICLR*, 2023b.
- 731 Qwen Team. Code with codeqwen1.5, April 2024. URL [https://qwenlm.github.io/](https://qwenlm.github.io/blog/codeqwen1.5/)  
732 [blog/codeqwen1.5/](https://qwenlm.github.io/blog/codeqwen1.5/).
- 733 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D. Manning, Stefano Ermon, and  
734 Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model.  
735 In *Advances in Neural Information Processing Systems*, 2023.
- 736 Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi  
737 Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, et al. Code llama: Open foundation models for  
738 code. *arXiv preprint arXiv:2308.12950*, 2023.
- 739 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy  
740 optimization algorithms. *CoRR*, abs/1707.06347, 2017.
- 741 Ilia Shumailov, Zakhar Shumaylov, Yiren Zhao, Yarin Gal, Nicolas Papernot, and Ross Ander-  
742 son. The curse of recursion: Training on generated data makes models forget. *arXiv preprint*  
743 *arXiv:2305.17493*, 2023.
- 744 Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling LLM test-time compute opti-  
745 mally can be more effective than scaling model parameters. *CoRR*, abs/2408.03314, 2024.
- 746 Zhengyang Tang, Xingxing Zhang, Benyou Wang, and Furu Wei. Mathscale: Scaling instruction  
747 tuning for mathematical reasoning. *arXiv preprint arXiv:2403.02884*, 2024.

- 756 Zhengwei Tao, Ting-En Lin, Xiancai Chen, Hangyu Li, Yuchuan Wu, Yongbin Li, Zhi Jin, Fei  
757 Huang, Dacheng Tao, and Jingren Zhou. A survey on self-evolution of large language models.  
758 *arXiv preprint arXiv:2404.14387*, 2024.
- 759  
760 Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy  
761 Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model.  
762 [https://github.com/tatsu-lab/stanford\\_alpaca](https://github.com/tatsu-lab/stanford_alpaca), 2023.
- 763 Jonathan Uesato, Nate Kushman, Ramana Kumar, H. Francis Song, Noah Y. Siegel, Lisa Wang,  
764 Antonia Creswell, Geoffrey Irving, and Irina Higgins. Solving math word problems with process-  
765 and outcome-based feedback. *arXiv preprint arXiv:2211.14275*, 2022.
- 766  
767 Peiyi Wang, Lei Li, Zhihong Shao, Runxin Xu, Damai Dai, Yifei Li, Deli Chen, Yu Wu, and Zhi-  
768 fang Sui. Math-shepherd: Verify and reinforce llms step-by-step without human annotations. In  
769 *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics*, pp.  
770 9426–9439. Association for Computational Linguistics, 2024a.
- 771 Tianduo Wang, Shichen Li, and Wei Lu. Self-training with direct preference optimization improves  
772 chain-of-thought reasoning, 2024b.
- 773  
774 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdh-  
775 ery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models.  
776 *arXiv preprint arXiv:2203.11171*, 2022.
- 777 Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and  
778 Hannaneh Hajishirzi. Self-instruct: Aligning language models with self-generated instructions.  
779 In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics*, pp.  
780 13484–13508. Association for Computational Linguistics, 2023a.
- 781 Ziqi Wang, Le Hou, Tianjian Lu, Yuexin Wu, Yunxuan Li, Hongkun Yu, and Heng Ji. Enable lan-  
782 guage models to implicitly learn self-improvement from data. *arXiv preprint arXiv:2310.00898*,  
783 2023b.
- 784  
785 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny  
786 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in*  
787 *neural information processing systems*, 35:24824–24837, 2022.
- 788 Yuxiang Wei, Zhe Wang, Jiawei Liu, Yifeng Ding, and Lingming Zhang. Magicoder: Empowering  
789 code generation with oss-instruct. In *International Conference on Machine Learning*. PMLR,  
790 2024.
- 791  
792 Yixuan Weng, Minjun Zhu, Shizhu He, Kang Liu, and Jun Zhao. Large language models are rea-  
793 soners with self-verification. *arXiv preprint arXiv:2212.09561*, 2, 2022.
- 794 Martin Weysow, Aton Kamanda, and Houari A. Sahraoui. Codeultrafeedback: An llm-as-a-judge  
795 dataset for aligning large language models to coding preferences. *CoRR*, abs/2403.09032, 2024.
- 796  
797 Wei Xiong, Hanze Dong, Chenlu Ye, Ziqi Wang, Han Zhong, Heng Ji, Nan Jiang, and Tong Zhang.  
798 Iterative preference learning from human feedback: Bridging theory and practice for RLHF under  
799 kl-constraint. In *Forty-first International Conference on Machine Learning*. OpenReview.net,  
800 2024.
- 801 An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li,  
802 Chengyuan Li, Dayiheng Liu, Fei Huang, et al. Qwen2 technical report. *arXiv preprint*  
803 *arXiv:2407.10671*, 2024a.
- 804  
805 An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu,  
806 Jianhong Tu, Jingren Zhou, Junyang Lin, Keming Lu, Mingfeng Xue, Runji Lin, Tianyu Liu,  
807 Xingzhang Ren, and Zhenru Zhang. Qwen2.5-math technical report: Toward mathematical ex-  
808 pert model via self-improvement. *CoRR*, abs/2409.12122, 2024b.
- 809 Yuqing Yang, Yan Ma, and Pengfei Liu. Weak-to-strong reasoning. *arXiv preprint*  
*arXiv:2407.13647*, 2024c.

810 Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T Kwok, Zhen-  
811 guo Li, Adrian Weller, and Weiyang Liu. Metamath: Bootstrap your own mathematical questions  
812 for large language models. *arXiv preprint arXiv:2309.12284*, 2023.

813 Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, Sainbayar Sukhbaatar, Jing Xu, and Jason  
814 Weston. Self-rewarding language models. *arXiv preprint arXiv:2401.10020*, 2024.

815 Xiang Yue, Xingwei Qu, Ge Zhang, Yao Fu, Wenhao Huang, Huan Sun, Yu Su, and Wenhao Chen.  
816 Mammoth: Building math generalist models through hybrid instruction tuning. *arXiv preprint*  
817 *arXiv:2309.05653*, 2023.

818 Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah Goodman. Star: Bootstrapping reasoning with  
819 reasoning. *Advances in Neural Information Processing Systems*, 35:15476–15488, 2022.

## 822 A MORE RESULTS AND ANALYSIS

823 We list detailed results, the source of training prompts, as well as some variants not presented in the  
824 main content, in Table 9, 11, and 12.

### 825 A.1 COMPARISON BETWEEN DPO AND PDPO

826 From the results across different base models and tasks, we find that pDPO usually have signif-  
827 icantly better results than DPO. The only exception appears on HumanEval and MBPP, which is  
828 possibly due to the two benchmark does not require detailed reasoning process, and the mismatch in  
829 distribution causes performance degradation.

### 830 A.2 VARIANTS OF ITERATIVE DPO TRAINING

831 In our experiments, we have employed several variants of Iterative DPO training for different con-  
832 siderations. For NuminaMath-790K, collecting all prompts for single iteration of DPO training can  
833 make it more challenging to avoid policy shifting, as pointed out by Xiong et al. (2024). To this  
834 end, we split the whole dataset into several parts for iterative training. During each iteration, we  
835 use around 160K prompts to collect solutions, construct pseudo feedback, and optimize the policy  
836 model.

837 For MathScale-300K, since the dataset is much smaller, we use all prompts across different itera-  
838 tions, and the process ends when we cannot observe significant improvements.

839 For code generation, we use a hybrid strategy. During each iteration, we introduce new dataset of  
840 prompts to be mixed with the original prompts, and collect new solutions. The pseudo test cases are  
841 also constructed based on the new solution programs. The reason is that we do not have too much  
842 data but we also hope to approach on-policy training.

843 During each iteration, the reference model is initialized from the policy model in last iteration,  
844 instead of the original base model.

### 845 A.3 EFFECT FROM THE AMOUNT OF SYNTHETIC TEST CASES

846 As shown in Table 2 and 3, PFPO-Self achieves even higher results than the model optimized on  
847 ground-truth test cases. One possible reason behind this is we have synthesized more test cases for  
848 self-consistency to reduce the false positive rate. To further verify this point, we down-sampled the  
849 synthetic test cases to align with the number of ground truth test cases. Specifically, for each coding  
850 problem, if the synthesized test cases exceeded the number of ground truth test cases, we randomly  
851 selected a subset of synthesized cases to match the number of ground truth test cases. Otherwise, all  
852 synthetic test cases were retained.

853 As shown in Table 7, the results demonstrate that reducing the number of synthetic test cases sig-  
854 nificantly decreases performance. This highlights the importance of scaling test cases for verifying  
855 program correctness effectively.



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

#### A.4 QUALITY ANALYSIS OF SYNTHETIC TEST CASES

In this section, we will analyze the quality of the synthetic test cases. Specifically, for each question in the APPs training set, we selected a ground-truth program to evaluate the synthetic test cases. A synthetic test case is valid, if and only if given the test case input and its output matches the output of the ground-truth program. We calculate the pass rate as metric, which represents the percentage of valid test cases among the total generated. The results are shown in Table 8, where the *Model* column denotes the model optimized by the pseudo feedback constructed from the corresponding test cases.

From the table, we can conclude that: (1) The synthetic test cases constructed by the weaker policy models has already demonstrate good quality. For example, 62.68% of the test cases for first round training, i.e., PFPO-Self-Self Iter. 0 are correct. (2) Iterative training on fixed prompt set will lead to overfitting. As the saturation of the quality of synthetic test cases, the policy model will fail to collect more correct instances from the training set, and the model will thus stop learning. (3) Pseudo feedback from frontier LLMs usually demonstrates better quality. We find that the test cases synthesized from frontier LLM’s solutions achieve 82.41% pass rate. Yet, we find that PFPO-LLM cannot outperform PFPO-Self Iter. 2, which is possibly due to less training data and the off-policy optimization of DPO, and can be alleviated by more rounds of training.

#### A.5 ANALYSIS ABOUT PLATEAUED PERFORMANCE OVER ITERATIONS

As also observed by Xiong et al. (2024) and Chen et al. (2024a), iterative style DPO tends to saturate after several iterations, regardless of whether the preference pairs are human labeled or synthetic. We believe there are at least two reasons for the saturation. The first one (as mentioned in Appendix A.4) is the quality of synthetic test cases may plateau at some point. The second reason could be that as the model improves with each iteration, more problems in the training set become “too easy” for the model at its current stage. For extremely challenging problems where no correct solution can be generated despite many attempts (sampling), we are unable to create valid preference pairs and are thus unlikely to solve them in the next iteration. And the number of moderately challenging problems, which are most effective for improving the model, gradually decreases after each iteration.

One possible solution to this is introducing unseen problems in each iteration to increase the likelihood of hitting new moderately challenging problems. In the third block of Table 1, we simulated the above setting by introducing new problems across iterations. As detailed in Appendix A, we split the NuminaMath dataset into five subsets of equal size and use a different subset for each iteration. We observed consistent improvements across iterations (see w/ PFPO-Self Iter. 1 to 5 rows in Table 1). However, please note that applying the above method requires a significant amount of unseen high-quality questions in that domain. Therefore, we cannot apply this strategy with “Mathstral-7B-v0.1” based models (Mathstral-7B-v0.1 has been trained on Numina already) or in coding domain (high quality prompts in coding domain is very limited). Finally, it is worth noting that our method’s advantage lies in its scalability, as no additional human annotation for answers is required when incorporating new problems.

#### A.6 PERFORMANCE ON MORE CHALLENGING DATASET

We also evaluated our performance on a subset of challenging math problems to explore whether the iterative training procedure also makes continuous improvements.

For the dataset construction, we began by sampling 16 solutions generated by Mathstral-7B w/ SFT on the MATH test set. We then identified questions for which the predicted answers, even after applying majority voting, remained incorrect. This yielded a MATH-Challenging test set of 1,580 questions. As shown in Table 10, the results indicate that our method can achieve improvements even on challenging questions over iterations.

## B HYPER-PARAMETERS

All hyper-parameters are listed in Table 5.

## C PROMPTS

### C.1 TEST CASE INPUTS GENERATION

The 2-shot prompt for test case inputs generation is shown in Figure 6, 7, and 8.

### C.2 PROMPT FOR COMPETITION-LEVEL CODE GENERATION WITH RATIONALE

The 1-shot prompt for competition-level code generation is shown in Figure 9, 10, and 11. During SFT, we remove the 1-shot example.

## D MORE IMPLEMENTATION DETAILS ABOUT PDPO

Our own implementation of pDPO include the following steps: (1) Sample multiple solutions for each problem from the policy model. (2) Supposing a non-empty line is a reasoning step, we sample multiple prefixes for each solution following a fixed ratio, where each prefix is composed of several continuous reasoning steps from the beginning. If the prefix dataset is too large, we will control that each problem will have fixed amount of prefixes. In most experiments, we will set the ratio as 30% and choose the fixed amount in  $\{8, 10, 20\}$ , as shown in Table 5. (3) Take the prompt, problem, as well as each solution prefix to compose new input, and sample 3 completions for it. (4) Estimate the expected returns of each prefix by checking the results approached by the completions. (5) Here is the main difference with the original implementation of pDPO (Jiao et al., 2024): The authors choose to first sample small amount of prefixes as well as the estimated expected returns for training another reward model, and use the reward model to annotate the complete solutions. This is to reduce the computation resource usage. Instead, in this paper, we simply estimate expected returns for all sampled solutions within given budget, and train the policy model on the prefixes (incomplete solutions) via DPO directly.

## E RESULTS OF DIFFERENT TRAINING ORDER

Table 6 compares the results of different training order. Assuming MathScale-300K and MathScale-500K share similar difficulty, the results indicate that the better quality the training data has, the earlier iteration it should be employed in.

## F RELATED WORK

### F.1 MATHEMATICAL REASONING

To enhance the model’s critical capability in mathematical reasoning, researchers have explored various techniques to guide the reasoning process. These efforts include prompts engineering (Wei et al., 2022), tool usage (He-Yueya et al., 2023; Chen et al., 2021b; 2022), **process rewarding** (Lightman et al., 2024; Wang et al., 2024a; Jiao et al., 2024; Lai et al., 2024), and verifiers (Cobbe et al., 2021b; Li et al., 2022a; Weng et al., 2022). In addition to improve frozen LLMs, researchers also work on synthesizing math related data to fine-tune them (Yu et al., 2023; Luo et al., 2023; Mitra et al., 2024; Yue et al., 2023).

### F.2 CODE GENERATION

Code generation has long been recognized as challenging due to data sparsity. The emergence of LLMs (Guo et al., 2024; DeepSeek-AI et al., 2024; Nijkamp et al., 2023b;a; Zelikman et al., 2022; Li et al., 2023a) pretrained on corpus including rich code has partially reduced this problem. Yet, it is still far from enough due to the significant gap between code completion and programming to solve

972 complex problems. Li et al. (2022b) make earlier efforts towards competition-level code generation,  
 973 which is achieved by filtering solutions from large sampling space via example test cases. Wei et al.  
 974 (2024) propose to synthesize instruction tuning dataset for code generation via self-instruct (Wang  
 975 et al., 2023a). Besides, Le et al. (2022); Liu et al. (2023) and Dou et al. (2024) propose to use the  
 976 feedback from compiler and unit test to improve code generation through reinforcement learning.  
 977 However, these approaches still rely on human annotation to obtain test cases, and thus cannot  
 978 be scaled to large scale training. Dou et al. (2024) further employ the unit test based feedback for  
 979 process supervision. Weyssow et al. (2024) have constructed a preference dataset for code generation  
 980 annotated by GPT-3.5. Yet, the absence of test cases hinder it to be generally employed in on-policy  
 981 learning.

982 Table 5: The hyper-parameters used in our experiments. **K - S.C.** refers to the amount of sampled  
 983 solutions for each problem that are used to determine the pseudo label via self-consistency. Instead,  
 984 **K - DPO** indicates the amount of sampled solutions per question used for constructing preference  
 985 pairs. **No. Prefix** is the amount of sampled prefixes per question used for process DPO training.  
 986 There are two kinds of values for *No. Prefix*. The integers are the exact amount, while the percentage  
 987 indicates that we sampled the specific ratio of prefixes among all solutions for each problem. This is  
 988 used to address the problem of limited training data.  $\beta$  is the coefficient used in DPO training.  $\alpha$  is  
 989 the weight for NLL loss.  $\epsilon$  is the reward (feedback) lower bound of positive samples to control the  
 990 quality. For pDPO training, sometimes we use fraction to simply indicate the number of sampled  
 991 completions as well as the successful attempts. For example,  $\frac{1}{3}$  means that we sample 3 completions  
 992 for each prefix, and if there is at least 1 completion is successful, it is kept.  $\sigma$  is the margin to control  
 993 the gap between positive and negative samples.

Model	K - S.C.	K - DPO	No. Prefix	$\beta$	$\alpha$	$\epsilon$	$\sigma$
Mathstral-7B w/ SFT							
w/ DPO (M.S.-500k, Iter. 0)	10	10	—	0.5	0.2	1.0	0.0
w/ pDPO (M.S.-500k, Iter. 0)	10	—	10	0.1	0.2	$\frac{2}{3}$	0.0
w/ pDPO (M.S.-300k-S.C., Iter. 1)	10	10	10	0.5	1.0	$\frac{1}{3}$	0.0
w/ pDPO (M.S.-300k-S.C., Iter. 2)	10	10	10	0.5	1.0	$\frac{1}{3}$	0.0
Llama-3.1-8b w/ SFT							
w/ DPO (M.S.-500k, Iter. 0)	10	10	—	0.5	0.2	1.0	0.0
w/ pDPO (M.S.-500k, Iter. 0)	10	10	10	0.5	0.2	$\frac{1}{3}$	0.0
w/ pDPO (Numina-S.C. 160k, Iter. 1)	16	—	8	0.5	0.2	$\frac{1}{3}$	0.0
w/ pDPO (Numina-S.C. 320k, Iter. 2)	16	16	8	0.5	1.0	$\frac{1}{3}$	0.0
w/ pDPO (Numina-S.C. 480k, Iter. 3)	16	—	8	0.5	1.0	$\frac{1}{3}$	0.0
w/ pDPO (Numina-S.C. 640k, Iter. 3)	16	—	8	0.6	0.2	$\frac{1}{3}$	0.0
w/ pDPO (Numina-S.C. 790k, Iter. 3)	16	—	32	0.6	0.2	$\frac{2}{3}$	0.0
Deepseek-coder-v1.5-chat w/ SFT							
w/ DPO (APPs)	10	10	—	0.1	0.0	1.0	0.0
w/ pDPO (APPs)	10	—	30%	0.1	0.2	$\frac{1}{5}$	0.0
w/ DPO (APPs - S.C.)	10	10	—	0.1	0.0	0.5	0.6
w/ pDPO (APPs - S.C.)	10	—	30%	0.1	0.0	0.5	0.4
w/ DPO (APPs & M.C. - S.C.)	10	10	—	0.4	0.2	0.5	0.6
w/ DPO (APPs & M.C. & xCode. - S.C.)	10	10	—	0.4	0.2	0.5	0.6
w/ pDPO (APPs & M.C. & xCode. - S.C.)	10	10	10	0.5	0.2	0.5	0.4
w/ pDPO (APPs & M.C. - S.C.)	10	10	20	0.4	0.2	0.5	0.4

1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079

Table 6: The experimental results of models trained via different orders. The prompt set for DPO training keeps the same with the original setting, where PFPO-Self employs MathScale-300K and PFPO-LLM takes MathScale-500K.

	MATH	GSM8K	Colledge Math
Mathstral-7B-v0.1	58.3	85.6	34.3
w/ SFT-4o	61.4	87.3	38.4
w/ PFPO-LLM Iter. 0	66.7	90.0	41.3
w/ PFPO-Self Iter. 1	67.8	<b>90.8</b>	42.0
w/ PFPO-Self Iter. 0	64.6	<b>89.8</b>	<b>39.4</b>
w/ PFPO-LLM Iter. 1	<b>65.2</b>	89.1	39.3

Table 7: The experimental results exploring the influence by the amount of synthetic test cases. *w/ down-sampled synthetic test cases* refers to the setting that we reduce the amount of synthetic test cases to which is equal and less than the that of the ground-truth test cases annotated for each problem in the training set of APPs.

	APPs	LiveCodeBench
DeepSeek-coder-v1.5-chat	14.3	21.1
w/ SFT	15.4	22.9
w/ golden test cases	16.9	22.9
w/ synthetic test cases	<b>17.4</b>	<b>23.4</b>
w/ down-sampled synthetic test cases	14.8	21.1

Table 8: The pass rate evaluated by executing the annotated ground-truth solution program on our synthetic test cases. The *Model* column denotes the one using the corresponding test cases for preference optimization.

Model	Pass Rate
PFPO-LLM Iter. 0 (APPs)	<b>82.41</b>
PFPO-Self Iter. 0 (APPs)	62.68
PFPO-Self Iter. 1 (APPs & M.C.)	66.80
PFPO-Self Iter. 2 (APPs & M.C. & xCode.)	66.75

Table 9: Overall results on mathematical reasoning benchmarks. The indentation across different levels represents the corresponding model is initialized from the parent-level. *M.S.* refers to the short of *MathScale*, and *S.C.* means *Self-Consistency*.

	MATH	GSM8K	College Math
Llama-3.1-8B-base			
w/ SFT (M.S.-500k)	53.8	85.1	34.6
w/ DPO (M.S.-500k)	53.8	86.0	35.1
w/ pDPO (M.S.-500k)	55.0	86.6	35.8
w/ pDPO (Numina-S.C. 160k)	55.9	87.6	36.6
w/ pDPO (Numina-S.C. 320k)	56.6	88.9	37.0
w/ pDPO (Numina-S.C. 480k)	57.0	88.8	36.7
w/ pDPO (Numina-S.C. 640k)	57.4	89.1	37.6
w/ pDPO (Numina-S.C. 790K)	<b>57.8</b>	<b>89.6</b>	<b>38.0</b>
Mathstral-7B-v0.1	58.3	85.6	34.3
w/ SFT (M.S.-500k)	61.4	87.3	38.4
w/ DPO (M.S.-500k)	63.0	88.6	39.1
w/ pDPO (M.S.-500k)	66.7	90.0	41.3
w/ pDPO (M.S.-300k-S.C., Iter. 0)	67.8	<b>90.8</b>	42.0
w/ pDPO (M.S.-300k-S.C., Iter. 1)	<b>68.6</b>	90.3	42.2
w/ pDPO (M.S.-300k-S.C., Iter. 2)	68.2	90.4	<b>42.3</b>

Table 10: The performance on the more challenging sub-test set of MATH, which is constructed from the failed questions Mathstral-7B w/ SFT via Maj@16.

Mathstral-7B	Greedy Decoding
w/ SFT	19.6
PFPO-LLM Iter. 0	23.6
PFPO-Self Iter. 1	24.6
PFPO-Self Iter. 2	26.7

Table 11: Overall results (Pass@1) on program generation benchmarks. *M.C.* refers to the prompt set of Magicoder (Wei et al., 2024), and *xCode*. is the short for xCodeEval (Khan et al., 2024b). *Introductory*, *Interview*, and *Competition* indicate the three difficulty levels of APPs. *S.C.* indicates that the test cases used for constructing preference dataset are synthesized through self-consistency. On the contrary, the row without *S.C.* refers that the test cases the golden ones from the original dataset.

	APPs				HumanEval	MBPP
	Overall	Introductory	Interview	Competition		
Deepseek-coder-v1.5-chat	14.3	35.7	10.8	3.2	75.6	73.9
w/ SFT (APPs - GPT-4o)	15.4	37.8	11.6	4.1	72.0	72.8
w/ DPO (APPs)	16.3	36.2	13.3	5.3	74.4	74.3
w/ pDPO (APPs)	16.9	37.3	13.8	6.1	73.8	73.2
w/ DPO (APPs - S.C.)	16.8	39.2	13.2	5.3	76.2	75.9
w/ pDPO (APPs - S.C.)	17.4	37.5	14.8	5.4	73.2	75.1
w/ DPO (APPs & M.C. - S.C.)	18.0	39.2	14.9	6.2	<b>79.3</b>	75.5
w/ DPO (APPs & M.C. & xCode. - S.C.)	<b>19.2</b>	<b>42.2</b>	15.8	6.5	75.0	74.2
w/ pDPO (APPs & M.C. & xCode. - S.C.)	19.1	40.9	<b>15.9</b>	<b>6.9</b>	73.8	75.1
w/ pDPO (APPs & M.C. - S.C.)	18.5	40.0	15.3	6.5	75.6	<b>76.3</b>

1134  
 1135  
 1136  
 1137  
 1138  
 1139  
 1140  
 1141  
 1142  
 1143  
 1144  
 1145  
 1146  
 1147  
 1148  
 1149  
 1150  
 1151  
 1152  
 1153  
 1154  
 1155  
 1156  
 1157  
 1158  
 1159  
 1160  
 1161  
 1162  
 1163  
 1164  
 1165  
 1166  
 1167  
 1168  
 1169  
 1170  
 1171  
 1172  
 1173  
 1174  
 1175  
 1176  
 1177  
 1178  
 1179  
 1180  
 1181  
 1182  
 1183  
 1184  
 1185  
 1186  
 1187

Table 12: Overall results on LiveCodeBench. We follow the recommended setting by sampling 10 solutions for each problem with temperature as 0.2, and estimating the Pass@1 results. The cutoff date of the test questions is from **2023-09-01** to **2024-09-01**. All results except those of our models are referenced from the official leaderboard.

	Overall	Easy	Medium	Hard
Claude-3.5-Sonnet	51.3	87.2	45.3	11.0
DeepSeekCoder-V2	41.9	79.9	32.0	4.9
Claude-3-Sonnet	26.9	67.2	7.3	1.4
Claude-3-Haiku	24.0	61.3	5.5	0.9
GPT-3.5-Turbo-0125	24.0	55.0	11.6	0.3
LLama-3-70b-Instruct	27.4	59.4	15.6	1.3
Deepseek-Coder-33b-Chat	23.4	56.1	8.6	0.9
Deepseek-coder-v1.5-chat	21.1	51.3	7.4	0.2
w/ SFT (APPs - GPT-4o)	22.9	53.0	10.6	0.2
w/ DPO (APPs)	22.9	53.7	9.4	1.0
w/ pDPO (APPs)	22.9	55.0	8.1	1.3
w/ DPO (APPs - S.C.)	24.2	55.7	11.0	0.8
w/ pDPO (APPs - S.C.)	23.4	54.2	10.3	0.7
w/ DPO (APPs & M.C. - S.C.)	23.7	55.8	9.5	1.1
w/ DPO (APPs & M.C. & xCode. - S.C.)	23.7	55.8	9.4	1.3
w/ pDPO (APPs & M.C. & xCode. - S.C.)	24.3	56.8	9.8	<b>1.6</b>
w/ pDPO (APPs & M.C. - S.C.)	<b>24.6</b>	<b>56.9</b>	<b>11.4</b>	0.2

1188  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241

You are an expert programmer. Your task is to write some test cases to the programming problems to help verify the expected program solutions. You only need to give me the inputs in the required format. Now, let me introduce the details to you:

#### ## Program Format

There will be two kinds of programming problems. One type of problem accepts standard input-output stream. As a result, the test case inputs should contain only the inputs text stream.

Another kind of problem is based on function calling, which shows a segment of starter code to illustrate the function head, defining the name of the arguments to be accepted. In this case, you should return me the inputs in the format of function calling, like `function\_name(\*arguments)`.

#### ## Response Format

You should return me the test case inputs in `json\_object` format. You need to generate **10** groups of test case inputs, and each key field is named as `test\_case\_i`, where `i` is the index of the test case. The value of each key is the test case inputs in the required format, which should be a string.

#### ## Examples for Standard Input-Output and Function Calling.

##### ### Standard Input-Output Stream

##### #### Programming Problem

Polycarp has  $n$  different binary words. A word called binary if it contains only characters '0' and '1'. For example, these words are binary: "0001", "11", "0" and "0011100".

Polycarp wants to offer his set of  $n$  binary words to play a game "words". In this game, players name words and each next word (starting from the second) must start with the last character of the previous word. The first word can be any. For example, these sequence of words can be named during the game: "0101", "1", "10", "00", "00001".

Word reversal is the operation of reversing the order of the characters. For example, the word "0111" after the reversal becomes "1110", the word "11010" after the reversal becomes "01011".

Probably, Polycarp has such a set of words that there is no way to put them in the order correspondent to the game rules. In this situation, he wants to reverse some words from his set so that: the final set of  $n$  words still contains different words (i.e. all words are unique); there is a way to put all words of the final set of words in the order so that the final sequence of  $n$  words is consistent with the game rules.

Polycarp wants to reverse minimal number of words. Please, help him.

-----Input-----

The first line of the input contains one integer  $t$  ( $1 \leq t \leq 10^4$ ) — the number of test cases in the input. Then  $t$  test cases follow.

The first line of a test case contains one integer  $n$  ( $1 \leq n \leq 2 \cdot 10^5$ ) — the number of words in the Polycarp's set. Next  $n$  lines contain these words. All of  $n$  words aren't empty and contains only characters '0' and '1'. The sum of word lengths doesn't exceed  $4 \cdot 10^6$ . All words are different.

Guaranteed, that the sum of  $n$  for all test cases in the input doesn't exceed  $2 \cdot 10^5$ . Also, guaranteed that the sum of word lengths for all test cases in the input doesn't exceed  $4 \cdot 10^6$ .

-----Output-----

Print answer for all of  $t$  test cases in the order they appear.

Figure 6: 2-shot prompt for test case inputs generation. Page 1.

1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295

If there is no answer for the test case, print -1. Otherwise, the first line of the output should contain  $k$  ( $0 \leq k \leq n$ ) — the minimal number of words in the set which should be reversed. The second line of the output should contain  $k$  distinct integers — the indexes of the words in the set which should be reversed. Words are numerated from 1 to  $n$  in the order they appear. If  $k=0$  you can skip this line (or you can print an empty line). If there are many answers you can print any of them.

----Example----

Input

```
4
4
0001
1000
0011
0111
3
010
101
0
2
00000
00001
4
01
001
0001
00001
```

Output

```
1
3
-1
0

2
1 2
```

#### Response

```
{
  "test_case_0": "3\n3\n101\n110\n011\n2\n01\n10\n4\n0001\n1000\n0011\n0111",
  "test_case_1": "2\n2\n01\n10\n3\n000\n111\n110",
  ...
}
```

### Function Calling

#### Programming Problem

Given a single positive integer  $x$ , we will write an expression of the form  $x (op1) x (op2) x (op3) x \dots$  where each operator  $op1, op2, \dots$  is either addition, subtraction, multiplication, or division (+, -, \*, or /). For example, with  $x = 3$ , we might write  $3 * 3 / 3 + 3 - 3$  which is a value of 3.

When writing such an expression, we adhere to the following conventions:

The division operator (/) returns rational numbers.

There are no parentheses placed anywhere.

We use the usual order of operations: multiplication and division happens before addition and subtraction.

It's not allowed to use the unary negation operator (-). For example, " $x - x$ " is a valid expression as it only uses subtraction, but " $-x + x$ " is not because it uses negation.

Figure 7: 2-shot prompt for test case inputs generation. Page 2.



1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349

We would like to write an expression with the least number of operators such that the expression equals the given target. Return the least number of operators used.

Example 1:

Input:  $x = 3$ , target = 19

Output: 5

Explanation:  $3 * 3 + 3 * 3 + 3 / 3$ . The expression contains 5 operations.

Example 2:

Input:  $x = 5$ , target = 501

Output: 8

Explanation:  $5 * 5 * 5 * 5 - 5 * 5 * 5 + 5 / 5$ . The expression contains 8 operations.

Example 3:

Input:  $x = 100$ , target = 100000000

Output: 3

Explanation:  $100 * 100 * 100 * 100$ . The expression contains 3 operations.

Note:

$2 \leq x \leq 100$

$1 \leq \text{target} \leq 2 * 10^8$

class Solution:

def leastOpsExpressTarget(self, x: int, target: int) -> int:

#### Response

```
{
  "test_case_0": "leastOpsExpressTarget(3, 19)",
  "test_case_1": "leastOpsExpressTarget(3, 32)",
  "test_case_2": "leastOpsExpressTarget(6, 100)",
  ...
}
```

## Get Started

Note that in the above examples, I omit some test case inputs. You should return **10** groups of inputs to me in 'json\_object' format.

#### Programming Problem

[[Question]]

#### Response

Figure 8: 2-shot prompt for test case inputs generation. Page 3.

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

You are an expert programmer. I will show you a programming problem. Please carefully comprehend the requirements in the problem, and write down the solution program to pass it under the given time and memory constraints.

**\*\*REMEMBER\*\*** to strictly follow the steps below to help reduce the potential flaws:

- (1) According to the input scale and the time/memory constraints, think about the time complexity and space complexity of your solution.
- (2) Think **\*\*step-by-step\*\*** to design the algorithm.
- (3) Translate your thoughts into Python program to solve it.

Besides, your Python solution program should be located between <BEGIN> and <END> tags:

```
<BEGIN>
t = int(input())
...
print(ans)
<END>
```

Here is an example:

```
## Problem
```

You are given an array  $a$  of length  $n$  consisting of zeros. You perform  $n$  actions with this array: during the  $i$ -th action, the following sequence of operations appears: Choose the maximum by length subarray (continuous subsegment) consisting only of zeros, among all such segments choose the leftmost one; Let this segment be  $[l; r]$ . If  $r-l+1$  is odd (not divisible by 2) then assign (set)  $a[\frac{l+r}{2}] := i$  (where  $i$  is the number of the current action), otherwise (if  $r-l+1$  is even) assign (set)  $a[\frac{l+r-1}{2}] := i$ .

Consider the array  $a$  of length 5 (initially  $a=[0, 0, 0, 0, 0]$ ). Then it changes as follows: Firstly, we choose the segment  $[1; 5]$  and assign  $a[3] := 1$ , so  $a$  becomes  $[0, 0, 1, 0, 0]$ ; then we choose the segment  $[1; 2]$  and assign  $a[1] := 2$ , so  $a$  becomes  $[2, 0, 1, 0, 0]$ ; then we choose the segment  $[4; 5]$  and assign  $a[4] := 3$ , so  $a$  becomes  $[2, 0, 1, 3, 0]$ ; then we choose the segment  $[2; 2]$  and assign  $a[2] := 4$ , so  $a$  becomes  $[2, 4, 1, 3, 0]$ ; and at last we choose the segment  $[5; 5]$  and assign  $a[5] := 5$ , so  $a$  becomes  $[2, 4, 1, 3, 5]$ .

Your task is to find the array  $a$  of length  $n$  after performing all  $n$  actions. Note that the answer exists and unique.

You have to answer  $t$  independent test cases.

```
----Input----
```

The first line of the input contains one integer  $t$  ( $1 \leq t \leq 10^4$ ) — the number of test cases. Then  $t$  test cases follow.

The only line of the test case contains one integer  $n$  ( $1 \leq n \leq 2 \cdot 10^5$ ) — the length of  $a$ .

It is guaranteed that the sum of  $n$  over all test cases does not exceed  $2 \cdot 10^5$  ( $\sum n \leq 2 \cdot 10^5$ ).

```
----Output----
```

For each test case, print the answer — the array  $a$  of length  $n$  after performing  $n$  actions described in the problem statement. Note that the answer exists and unique.

Figure 9: 1-shot competition-level code generation prompt. When being applied to SFT, the 1-shot example is removed. Page 1.

1404  
1405  
1406  
1407 If there is no answer for the test case, print -1. Otherwise, the first line of the output should contain  $k$  ( $0 \leq k \leq n$ ) — the minimal number of words in the set which should be reversed. The second line of the output  
1408 should contain  $k$  distinct integers — the indexes of the words in the set which should be reversed. Words are  
1409 numerated from  $1$  to  $n$  in the order they appear. If  $k=0$  you can skip this line (or you can print an empty  
1410 line). If there are many answers you can print any of them.  
1411  
1412 -----Example-----  
1413 Input  
1414 4  
1415 4  
1416 0001  
1417 1000  
1418 0011  
1419 0111  
1420 3  
1421 010  
1422 101  
1423 0  
1424 2  
1425 00000  
1426 00001  
1427 4  
1428 01  
1429 001  
1430 0001  
1431 00001  
1432  
1433 Output  
1434 1  
1435 3  
1436 -1  
1437 0  
1438  
1439 2  
1440 1 2  
1441  
1442 ##### Response  
1443  
1444 {  
1445 "test\_case\_0": "3\n3\n101\n110\n011\n2\n01\n10\n4\n0001\n1000\n0011\n0111",  
1446 "test\_case\_1": "2\n2\n01\n10\n3\n000\n111\n110",  
1447 ...  
1448 }  
1449  
1450 ### Function Calling  
1451  
1452 ##### Programming Problem  
1453  
1454 Given a single positive integer  $x$ , we will write an expression of the form  $x (op1) x (op2) x (op3) x \dots$  where each  
1455 operator  $op1, op2, \dots$  is either addition, subtraction, multiplication, or division (+, -, \*, or /). For example, with  
1456  $x = 3$ , we might write  $3 * 3 / 3 + 3 - 3$  which is a value of 3.  
1457 When writing such an expression, we adhere to the following conventions:  
  
The division operator (/) returns rational numbers.  
There are no parentheses placed anywhere.  
We use the usual order of operations: multiplication and division happens before addition and subtraction.  
It's not allowed to use the unary negation operator (-). For example, " $x - x$ " is a valid expression as it only uses  
subtraction, but " $-x + x$ " is not because it uses negation.

Figure 10: 1-shot competition-level code generation prompt. When being applied to SFT, the 1-shot example is removed. Page 2.

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

We would like to write an expression with the least number of operators such that the expression equals the given target. Return the least number of operators used.

Example 1:

Input:  $x = 3$ , target = 19

Output: 5

Explanation:  $3 * 3 + 3 * 3 + 3 / 3$ . The expression contains 5 operations.

Example 2:

Input:  $x = 5$ , target = 501

Output: 8

Explanation:  $5 * 5 * 5 * 5 - 5 * 5 * 5 + 5 / 5$ . The expression contains 8 operations.

Example 3:

Input:  $x = 100$ , target = 100000000

Output: 3

Explanation:  $100 * 100 * 100 * 100$ . The expression contains 3 operations.

Note:

$2 \leq x \leq 100$

$1 \leq \text{target} \leq 2 * 10^8$

class Solution:

def leastOpsExpressTarget(self, x: int, target: int) -> int:

#### Response

```
{
  "test_case_0": "leastOpsExpressTarget(3, 19)",
  "test_case_1": "leastOpsExpressTarget(3, 32)",
  "test_case_2": "leastOpsExpressTarget(6, 100)",
  ...
}
```

## Get Started

Note that in the above examples, I omit some test case inputs. You should return **\*\*10\*\*** groups of inputs to me in 'json\_object' format.

#### Programming Problem

[[Question]]

#### Response

Figure 11: 1-shot competition-level code generation prompt. When being applied to SFT, the 1-shot example is removed. Page 3.