ACECODER: Acing Coder RL via Automated Test-Case Synthesis

Anonymous ACL submission

Abstract

Most progress in recent coder models has been driven by supervised fine-tuning (SFT), while the potential of reinforcement learning (RL) remains largely unexplored, primarily due to the lack of reliable reward data/model in the code domain. In this paper, we address this challenge by leveraging automated large-scale testcase synthesis to enhance code model training. Specifically, we design a pipeline that generates extensive (question, test-cases) pairs from existing code data. Using these test cases, we construct preference pairs based on pass rates over sampled programs to train reward models with Bradley-Terry loss. It shows an average of 10-point improvement for Llama-3.1-8B-Ins and 5-point improvement for Qwen2.5-Coder-7B-Ins through best-of-32 sampling, making the 7B model on par with 236B DeepSeek-V2.5. Furthermore, we conduct reinforcement learning with both reward models and testcase pass rewards, leading to consistent improvements across HumanEval, MBPP, Big-CodeBench, and LiveCodeBench (V4). Notably, we follow the R1-style training to start from Owen2.5-Coder-base directly and show that our RL training can improve model on HumanEval-plus by over 25% and MBPP-plus by 6% for merely 80 optimization steps. We believe our results highlight the huge potential of reinforcement learning in coder models.

1 Introduction

002

011

017

021

022

024

042

In recent years, code generation models have advanced significantly with compute scaling (Kaplan et al., 2020) and training data quality improvement (Huang et al., 2024; Lozhkov et al., 2024; Guo et al., 2024b). The state-of-the-art coder models, including Code-Llama (Rozière et al., 2023), Qwen2.5-Coder (Hui et al., 2024a), DeepSeek-Coder (Guo et al., 2024a) and so on, have shown unprecedented performance across a wide range of coding tasks like program synthesis (Chen et al., 2021), program repair (Zheng et al.,



Figure 1: Overall Workflow of our model: we start from the seed code dataset to create well-formatted questions and corresponding test cases. Then we adopt strong models to filter the noisy test cases. Finally, we adopt these test cases to harvest positive and negative program pairs for reward model training and RL.

2024a), optimization (Shypula et al., 2023), test generation (Steenhoek et al., 2023), SQL (Yu et al., 2018), issue fix (Jimenez et al., 2024). These models are all pre-trained and further supervised finetuned (SFT) on large-scale coding data from web resources like Common Crawl or Github.

Though strong performance has been achieved through SFT (Luo et al., 2023; Wei et al., 2024), very few models have explored the potential of reinforcement learning (RL) (Ouyang et al., 2022a), which has proven effective in other domains such as mathematical reasoning like DeepSeek-R1 (Shao et al., 2024). We argue that this absence of RLbased training in coder models is primarily due to two key challenges:

(1) Lack of reliable reward signals for code generation. In tasks such as mathematical problem-solving, rewards can be easily derived from rule-based string matches with reference answers (Guo et al., 2025) or large-scale human annotations (Ouyang et al., 2022b). In contrast, evaluating code quality typically requires executing test cases to measure the pass rate, making reward signal design more complex. This also explains why existing reward models like Skywork (Liu et al., 2024a) can hardly generalize to the coding domain (see subsection 4.4).

058

059

060

063

064

067

081

084

101

102

103

104 105

106

107

109

(2) Scarcity of large-scale coding datasets with reliable test cases. Most existing coding datasets like APPS (Hendrycks et al., 2021; Chen et al., 2021) heavily rely on costly human expert annotations for test cases, which limits their scalability for training purposes. The largest data is TACO (Li et al., 2023) with 25K examples, which are crawled from the popular coding competition websites, which were already heavily exploited during the pre-training phase.

Therefore, we curate ACECODE-87K, on which we trained our reward models: ACECODE-RM-7B and ACECODE-RM-32B. Comprehensive experiments of best-of-N sampling show that ACECODE-RM can significantly boost existing LLM's performance on coding benchmarks. For example, ACECODE-RM-7B can improve the performance of Llama-3.1-8B-Instruct by an average of 8.4 points across the 4 coding benchmarks, i.e. HumanEval (Liu et al., 2023), MBPP (Liu et al., 2023), BigCodeBench (Zhuo et al., 2024) and LiveCodeBench (Jain et al., 2024). Even for the stronger coder model Qwen2.5-Coder-7B-Instruct, our "7B+7B" combination still gets an average of 2.6 improvements. ACECODE-RM-32B is even more powerful, which pushes the former two numbers to 10.7 and 4.7 respectively, showcasing the effectiveness of ACECODE-RM.

Furthermore, we adopt ACECODE-RM-7B and test case pass rate separately to do reinforcement learning with reinforce++ (Hu, 2025) over coder models. Experiments show 2.1 and 0.7 points of average improvement when starting from Qwen2.5-7B-Ins and the Qwen2.5-Coder-7B-Ins respectively, making the latter even more powerful than GPT-4-Turbo on benchmarks like MBPP. Inspired by the recent DeepSeek-R1 (Guo et al., 2025), we also perform RL training directly from the Qwen2.5-Coder-7B-base model and saw a surprising 25% improvement on HumanEval-plus and 6% improvement on MBPP-plus (Liu et al., 2023) with merely 80 optimization steps (48 H100 GPU hours). These improvements are also generalizable to other more difficult benchmarks.

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

146

147

148

149

150

151

152

153

To our knowledge, this is the first work to perform reward model training and reinforcement learning for code generation using a fully automated pipeline that synthesizes large-scale reliable tests. We believe our ACECODE-87K will unlock the potential of RL training for code generation models and help the community to further push the boundaries of LLM's coding abilities.

2 Methodology

In this section, we will introduce the overall methodology of ACECODER. We begin with formulations of the problems we are investigating, including reward model training and reinforcement learning for LLMs. We then elaborate on how we synthesize the test cases and construct the ACECODE-87K. Finally, we explain how we perform the reinforcement learning using our ACECODE-RM trained on the ACECODE-87K.

2.1 Problem Formulation

Reward Model Training Let x denote the coding question and $\mathbf{y} = \{y_1, \dots, y_t\}$ denote the program solution, where y_i represents the *i*-th token of the program solution and $(\mathbf{x}, \mathbf{y}) \in D$. Assuming θ represents the parameters of the model, then n responses $(\mathbf{y}^1, \dots, \mathbf{y}^n)$ will be sampled from the model π_{θ} given the input \mathbf{x} . Let (s_1, \dots, s_n) be the target rewards, i.e. the test case pass rates in our scenario, then we define the Bradley-Terry loss (Bradley and Terry, 1952) for every pair of responses \mathbf{y}^i and \mathbf{y}^j with scores of s_i and s_j when we are training a reward model R_{ϕ} as follows:

$$\mathcal{L}_{\phi}(\mathbf{x}, s_i, s_j) = \mathbb{1}[s_i > s_j] \log \sigma(R_{\phi}(\mathbf{x}, \mathbf{y}^i) - R_{\phi}(\mathbf{x}, \mathbf{y}^j))$$
¹⁴⁸

where $\mathbb{1}[\cdot] = 1$ if the expression inside the brackets is true, otherwise, it's 0. The final loss function for the reward training is:

$$\mathcal{L}(\phi) = -\frac{1}{n(n-1)} \sum_{i=1}^{n} \sum_{j=1}^{n} \mathcal{L}_{\phi}(\mathbf{x}, s_i, s_j) \quad (1)$$

That means the reward model is trained to assign higher values to preferred responses and lower values to non-preferred ones, maximizing the difference between these ratings.

238

239

240

241

242

243

244

245

199

200

Best-of-N Sampling After we get the trained reward model R_{ϕ} , one way to quickly test the performance of the reward model is Best-of-N sampling, which is usually used as a test-time scaling approach. We will simply select the best response according to the predicted value of R_{ϕ} . That is $\mathbf{y}^* = \arg \max_{\mathbf{v}^i \in \mathbf{v}^1, \dots, \mathbf{v}^N} R_{\phi}(\mathbf{x}, \mathbf{y}^i)$.

161

162

163

165

168

169

170

171

175

176

177

178

179

180

181

182

183

184

186

187

190

192

194

196

198

Reinforcement Learning We can finally conduct reinforcement learning for the original policy model π_{θ} after we get a well-trained reward model R_{ϕ} . Proximal Policy Optimization (PPO) is an actor-critic RL algorithm that is widely used for LLM's RL process. Let $\pi_{\theta_{old}}$ be the reference model and π_{θ} be the current policy model that we are updating frequently during the RL training. We denote $r_t(\theta)$ as the probability ratio of the current policy model over the old policy model on the *t*-th generated token:

$$r_t(\theta) = \frac{\pi_{\theta}(y_t | \mathbf{x}, \mathbf{y}_{< t})}{\pi_{\theta_{old}}(y_t | \mathbf{x}, \mathbf{y}_{< t})}$$
(2)

Then the PPO algorithms optimize the LLM by thefollowing surrogate objective:

$$\mathcal{L}_{PPO}(\theta) = -\frac{1}{|\mathbf{y}|} \sum_{t=1}^{|\mathbf{y}|} \min \left[r_t(\theta) A_t, \operatorname{clip}\left(r_t(\theta), 1-\epsilon, 1+\epsilon\right) A_t \right]$$

where $\mathbf{y} \sim \pi_{\theta_{old}}(\cdot|\mathbf{x})$, and A_t is the advantage computed through the Generalized Advantage Estimation (GAE) (Schulman et al., 2015) via the rewards generated by R_{ϕ} and the learned value function V_{ψ} . The PPO training objective will force the policy model π to increase the probability of generating tokens with higher A_t and decrease the probability ratio of generating tokens with lower A_t until the clipped bounds $1 + \epsilon$ and $1 - \epsilon$ are reached respectively.

However, PPO usually requires training an additional value model V_{ψ} and thus makes the training inefficient. Recently, there are some other works like Reinforecement++ (Hu, 2025) that eliminate the need for value model but instead compute advantage only using the rewards generated by R_{ϕ} and the KL-divergence of the tokens after the *t*-th tokens. This makes the RL process more efficient and has also proved to be more stable.

3 ACECODE-87K

To be able to train a reward model specifically designed for code generation, the first thing is to synthesize reliable test cases for each coding problem and use them as training signals. In this section, we explain the whole procedure of constructing ACECODE-87K step by step. We show the overall statistics in Table 1.

Test Case Synthesis from Seed Dataset We start from existing coding datasets with provided question x and corresponding program y. Specifically, we combine Magicoder-Evol-Instruct¹, Magicoder-OSS-Instruct-75K², and StackPyFunction³ as our seed dataset. We only keep the questions written in Python that contain either a function or a class, resulting in a total of 124K entries. We find that these datasets contain highly noisy questions that could not be easily evaluated using test cases. Therefore, we feed every question-solution pair (x, y) into a GPT-40-mini (Hurst et al., 2024) to propose a refined LeetCode-style question $\mathbf{x}_{\mathbf{r}}$ with highly structured instructions. Meanwhile, we also prompt it to 'imagine' around 20 test cases $(t_1, ..., t_m)$ for each refined coding question $\mathbf{x}_{\mathbf{r}}$ based on its understanding of the expected behavior of the desired program. See prompt template used in subsection A.2. Please note that we do not use the program solution y from the existing datasets at all in our final curated ACECODE-87K. These datasets are purely used as seeds to help LLM formulate well-structured coding problems.

Test Case Filtering These 'imagined' test cases generated from the LLM contain severe hallucinations. To filter out those hallucinated test cases, we facilitated a stronger coder model Qwen2.5-Coder-32B-Instruct (Hui et al., 2024a) as a proxy to perform quality control. Specifically, we prompt it for each $\mathbf{x}_{\mathbf{r}}$ to generate a program \mathbf{y}' and then run these programs over the test cases to approximate their quality. We removed all test cases t_i where the generated solution program y' could not pass. Furthermore, we removed questions with fewer than 5 tests after filtering, as these questions might be overly ambiguous. With the above filtering, we constructed the ACECODE-87K with 87.1K distinct coding questions and 1.38M cleaned test cases, as represented by $(\mathbf{x}_{\mathbf{r}}, (t_1, ..., t_{m_c}))$, where m_c represents the number of test cases after filtering.

Preference Pairs Construction We propose to use the Bradley-Terry model to train the reward model as defined in Equation 1. Therefore, we

¹ise-uiuc/Magicoder-Evol-Instruct-110K

²ise-uiuc/Magicoder-OSS-Instruct-75K

³bigcode/stack-dedup-python-fns

Subset	Evol	OSS	Stack Python	Overall						
Before Filtering										
# Examples # Avg Test Cases	36,256 19.33	37,750 17.21	50,000 18.27	124,006 18.26						
After Filtering										
# Examples # Avg Test Cases # Pairs	26,920 15.14 89,089	25,862 16.33 91,636	34,367 16.08 126,784	87,149 15.87 307,509						

Table 1: Dataset statistics of ACECODE-87K before and after test-case filtering.

need to construct (question, [positive program, negative program]) data from ACECODE-87K. Specifically, we sample programs $(\mathbf{y}^1, ..., \mathbf{y}^n)$ from existing models (e.g. Llama-3.1 (Grattafiori et al., 2024)) w.r.t \mathbf{x}_r and utilize the test-case pass rate to distinguish positive and negative programs. Since the pass rate s_i for the sampled program \mathbf{y}^i can be any number between [0, 1], a minor difference in pass rate may not represent that one program is more accurate than another. Therefore, instead of using $\mathbb{1}[s_i > s_j]$ to select the preference pairs, we have thus modified the selection rules to be:

246

247

253

254

257

261

263

265

269

271

272

273

276

278

279

281

$$\mathbb{1}[s_i > s_j + 0.4, s_i > 0.8, s_j > 0]$$
(3)

This is to ensure the preferred program has at least a 0.8 pass rate to make sure it represents a more correct program. Also, we find many sampled programs with 0 pass rates can be caused by some small syntax errors or some Python packaging missing errors during evaluation, we chose to not include them as the preference pair to make sure our constructed datasets represent only the preferencebased on the valid pass rate. We also ensure the sampled programs all come from the backbone of R_{ϕ} so the reward model is trained in an on-policy way. After that, we train our reward model R_{ϕ} by fully fine-tuning an instruct coding model. Specifically, We extract the last token's final hidden representations and pass it through a linear model head that generates a single scalar output, which is optimized via the loss function defined in Equation 1.

4 Experiments

4.1 Reward Model Training Setup

We mainly use Qwen2.5-Coder-7B-Instruct ⁴ as the backbone of the reward model and sample 16 responses from it for each question in ACECODE-87K. Finally, following the rule defined in Equation 3, around 300K preference pairs were created out of 46,618 distinct questions (37.34% of the total questions) that have at least one pair satisfying the condition, and other questions are not used. 283

284

285

289

290

291

292

293

294

296

297

298

300

301

302

303

304

305

306

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

329

330

331

Our reward model is trained using LlamaFactory (Zheng et al., 2024b). We apply full finetuning with DeepSpeed stage 3. We train for 1 epoch using a cosine learning rate schedule, starting at 1e-5 with a warmup ratio of 0.1 to gradually increase the learning rate in the initial training phase. Training batch size is set to 128. We enable bf16 precision to reduce memory overhead without compromising model fidelity. The training takes 24 hours on 8 x A100 GPUs.

4.2 Reinforcement Learning Setup

We perform RL training from three policy models: Owen2.5-7B-Instruct⁵ and Owen2.5-Coder-7B-Base ⁶ and Qwen2.5-Coder-7B-Instruct. Two types of reward can be used, i.e. the trained reward model ACECODE-RM-7B and the rule-based reward, i.e. pass rate over the test cases in ACECODE-87K. During training, we set the pass rate to be a binary reward, which is 1.0 when all test cases passed, otherwise 0. This is similar to the verfiable reward used in Tulu3 (Lambert et al., 2024a) and DeepSeek-R1 (Guo et al., 2025). Similar to DeepSeek-R1 (Guo et al., 2025), we also experiment with RL from the base model because SFT may cause the search space of the model to be stuck in the local minimum. Since coding is also a highly verifiable task like math, we include the Qwen2.5-Coder-7B-Base in our experiments.

We have trained different policy model backbones with different rewards, resulting in 6 RL models in total. All the RL-tuning are based on OpenRLHF (Hu et al., 2024). We adopt the Reinforcement++ (Hu, 2025) algorithm instead of PPO to improve the training efficiency without training the value model. It's also proved to be more stable than PPO and GRPO. We train our model on a subsampled hard version of ACECODE-87K, where we keep the 25% of the questions with lower average pass rates and higher variance. This is to ensure the question is hard and that the sampled programs are diverse enough. For the training hyperparameters, we set the rollout batch size to 256, and 8 programs are sampled from per question. The training batch size is 128 with a learning rate of 5e-7. All the models are trained for 1 episode and finished in 6 hours on 8 x H100 GPUs.

⁴Qwen/Qwen2.5-Coder-7B-Instruct

⁵Qwen/Qwen2.5-7B-Instruct

⁶Qwen/Qwen2.5-Coder-7B

Mehod	# N	Huma -	anEval Plus	_ MI	3PP Plus	BigCod Full	eBench-C Hard	BigCod Full	eBench-I Hard	LiveCodeBench V4	Average
GPT-40 (0806)	1	92.7	87.2	87.6	72.2	58.9	36.5	48.0	25.0	43.6	61.3
DeepSeek-V2.5	1	90.2	83.5	87.6	74.1	53.2	29.1	48.9	27.0	41.8	59.5
DeepSeek-V3	1	91.5	86.6	87.6	73.0	62.2	39.9	50.0	27.7	63.5	64.6
Qwen2.5-Coder-32B	1	92.1	87.2	90.5	77.0	58.0	33.8	49.0	27.7	48.3	62.6
Inference Model = Mistral-7B-Instruct-V0.3											
Greedy	1	36.6	31.1	49.5	41.3	25.9	6.1	20.1	5.4	7.3	24.8
Average	64	37.1	30.8	45.1	38.0	21.7	4.2	17.6	3.0	4.0	22.4
Oracle	64	87.2	78.0	83.9	73.5	68.4	37.8	58.5	31.1	24.3	60.3
	16	65.9	56.7	59.3	52.4	35.1	10.1	29.3	8.8	11.9	36.6
AceCodeRM-7B	32	68.3	58.5	59.8	51.6	37.4	8.8	30.7	10.8	14.6	37.8
	64	71.3	61.6	59.8	51.6	39.4	6.8	31.8	9.5	15.4	38.6
Δ (RM-greedy)	-	+34.8	+30.5	+10.3	+11.1	+13.5	+4.1	+11.7	+5.4	+8.1	+13.8
	16	68.3	61.0	58.7	49.5	37.7	11.5	30.9	10.1	12.9	37.8
AceCodeRM-32B	32	72.6	65.9	61.6	51.6	40.5	9.5	33.9	13.5	16.1	40.6
	64	75.0	64.6	60.6	50.0	42.7	15.5	35.6	13.5	17.4	41.7
Δ (RM-greedy)	-	+38.4	+34.8	+12.2	+11.1	+16.8	+9.5	+15.5	+8.1	+10.1	+16.9
				Inference	e Model =	= Llama-3	.1-8B-Instru	ıct			
Greedy	1	68.9	62.2	67.2	54.8	38.5	12.8	31.8	13.5	18.0	40.9
Average	64	61.7	54.9	64.5	54.5	32.8	10.1	26.6	9.0	13.8	36.4
Oracle	64	93.9	90.2	92.1	82.3	80.0	54.7	67.9	48.6	40.8	72.3
	16	77.4	70.7	76.5	64.3	45.8	20.3	36.4	12.2	26.1	47.7
AceCodeRM-7B	32	79.9	72.6	76.2	62.4	47.6	23.0	37.3	13.5	27.3	48.9
	64	81.7	74.4	74.6	61.9	47.8	23.6	38.1	13.5	27.6	49.3
Δ (RM-greedy)	-	+12.8	+12.2	+9.3	+9.5	+9.3	+10.8	+6.2	0.0	+9.6	+8.4
	16	82.3	74.4	72.8	60.6	49.8	20.3	38.4	13.5	27.5	48.8
AceCodeRM-32B	32	81.7	76.2	72.8	60.6	50.4	22.3	39.1	13.5	30.3	49.6
	64	85.4	79.3	72.0	59.0	48.5	19.6	40.0	13.5	31.0	49.8
Δ (RM-greedy)	-	+16.5	+17.1	+9.3	+9.5	+11.8	+10.8	+8.2	+0.0	+13.0	+9.0
			Inf	erence M	Iodel = Q	wen2.5-C	oder-7B-In	struct			
Greedy	1	91.5	86.0	82.8	71.4	49.5	19.6	41.8	20.3	34.2	55.2
Average	64	86.0	80.1	77.9	65.6	45.3	18.6	37.3	16.2	31.8	51.0
Oracle	64	98.2	95.7	97.4	90.7	80.9	62.8	73.5	53.4	57.4	78.9
	16	90.2	82.9	88.6	74.9	53.8	20.9	45.0	21.6	40.1	57.6
AceCodeRM-7B	32	90.9	86.0	87.8	74.1	53.4	25.0	43.9	19.6	39.8	57.8
	64	90.9	85.4	87.6	73.8	52.9	24.3	43.5	21.6	40.1	57.8
Δ (RM-greedy)	-	-0.6	0.0	+5.8	+3.4	+4.3	+5.4	+3.2	+1.4	+5.9	+2.6
	16	90.2	86.6	88.4	74.9	53.9	25.0	45.4	19.6	44.0	58.7
AceCodeRM-32B	32	90.2	86.6	88.4	75.4	55.4	29.7	45.6	21.6	43.5	59.6
	64	89.6	86.0	87.8	75.1	55.0	26.4	46.1	22.3	44.5	59.2
Δ (RM-greedy)	-	-0.6	+0.6	+5.8	+4.0	+6.0	+10.1	+4.3	+2.0	+10.3	+4.4

Table 2: ACECODE-RM's best-of-n results on several benchmarks. Specifically, -C means completion split and -I means instruct split of BigCodeBench. The Δ might be off by 0.1 due to rounding.

4.3 Evaluation Setup

332

333

334

335

336

338

341 342

343

345

346

We evaluate our method on four established code-focused benchmarks: HumanEval(+) (Chen et al., 2021), MBPP(+) (Austin et al., 2021), BigCodeBench (Zhuo et al., 2024) and Live-CodeBench (V4) (Jain et al., 2024). These benchmarks collectively cover a diverse array of coding tasks, enabling us to assess both the correctness and quality of generated code. For Best-of-N sampling, we adopt top-p sampling with a temperature of 1.0 to generate multiple (16/32/64) candidate solutions per question and then select the response with the highest reward for evaluation. For RL experiments, we use each benchmark's default setting, which is greedy sampling most of the time.

4.4 Main Results

Here we show the experimental results of reward model and RL-trained model.

RM Results We conduct Best-of-N experiments on 3 inference models, specifically Mistral-Instruct-V0.3-7B(AI, 2023), Llama-3.1-Instruct-8B (Grattafiori et al., 2024), and Qwen2.5-Coder-7B-Insutret (Hui et al., 2024b; Yang et al., 2024a). We additionally report the average score across all generated samples and also the oracle score (pass@N) for better comparison.

According to Table 2, ACECODE-RM can consistently boost the performance of inference models by a large margin compared to the greedy decoding results. On weaker models like Mistral (AI, 2023)

357

358

360

361

347

348

Model	Huma	nEval Plus	MI	BPP Plus	BigCo Full	deBench (C) Hard	BigCoo Full	leBench (I) Hard	LiveCodeBench V4	Average	
	-	1 103	_	Tius	1 uli	Haru	1 ull	Thatu	* 7		
RLEF-8B	-	67.5	-	57.0	-	-	-	-	-	-	
RLEF-70B	-	78.5	-	67.6	-	-	-	-	-	-	
PPOCoder-7B	78.7	-	67.0	-	-	-	-	-	-	-	
StepCoder-7B	76.8	-	63.8	-	-	-	-	-	-	-	
CodeGemma-7B	60.5	-	55.2	-	-	-	-	-	-	-	
DSTC-33B	79.9	72.0	82.5	70.4	51.6	22.3	41.0	18.2	-	-	
				Basel	ine = Qv	en2.5-7B-Ins	truct				
Baseline	81.7	73.2	79.4	67.7	45.6	16.9	38.4	14.2	29.0	49.6	
AceCoder _{<i>BM</i>}	83.5	77.4	83.1	71.2	46.8	16.9	39.0	14.9	30.3	51.5	
AceCoder _{Rule}	84.1	77.4	80.2	68.3	46.8	15.5	40.2	15.5	30.1	50.9	
Δ (RL-baseline)	+2.4	+4.3	+3.7	+3.4	+1.2	0.0	+1.8	+1.4	+1.3	+2.0	
				Baselin	e = Qwe	n2.5-Coder-7I	B-Base				
Baseline	61.6	53.0	76.9	62.9	45.8	16.2	40.2	14.2	28.7	44.4	
AceCoder _{RM}	83.5	75.6	80.2	67.2	41.9	14.9	36.8	16.2	25.7	49.1	
AceCoder _{Rule}	84.1	78.0	82.3	69.3	48.6	18.2	43.2	18.2	28.5	52.3	
Δ (RL-baseline)	+22.5	+25.0	+5.4	+6.4	+2.8	+2.0	+3.1	+4.1	-0.2	+7.9	
Baseline = Qwen2.5-Coder-7B-Instruct											
Baseline	91.5	86.0	82.8	71.4	49.5	19.6	41.8	20.3	34.2	55.2	
$AceCoder_{RM}$	89.0	84.1	86.0	72.8	50.4	18.9	42.0	19.6	35.0	55.3	
AceCoder _{Rule}	90.9	84.8	84.1	71.7	50.9	23.0	43.3	19.6	34.9	55.9	
Δ (RL-baseline)	-0.6	-1.2	+3.2	+1.3	+1.4	+3.4	+1.5	-0.7	+0.8	+0.7	

Table 3: ACECODER's Performance after RL tuning using Reinforcement++ algorithm. We start with 3 different initial policy models and 2 kind of reward types, where *RM* means using our trained ACECODE-RM and *Rule* means using the binary pass rate. Results show consistent improvement across various benchmarks.

and Llama-3.1 (Zheng et al., 2024b), the overall improvements are greater than 10 points. These improvements can be attributed to our reward model's ability to identify high-quality completions among multiple candidates, thereby reducing the impact of suboptimal sampling on the final output. Notably, these gains become more pronounced on benchmarks where the gap between greedy decoding and oracle performance (i.e., the best possible completion among all samples) is larger. In such cases, the variance among sampled completions is relatively high, providing greater opportunities for the reward model to pinpoint and elevate top-tier responses.

363

367

370

375

377

381

Greedy decoding systematically outperforms the average sampled performance, reflecting the strong code generation capability of these inference models. Consequently, while most reward models achieve best-of-N results above the average, we consider a reward model effective only if it surpasses the performance of greedy decoding.

382**RL Results**We perform RL training over 3 dif-383ferent initial policy models in Table 3 with model-384based and rule-based rewards. When starting from385Qwen2.5-Instruct-7B, we can see the RL tuning can386consistently improve the performance, especially387on HumanEval and MBPP. Even for the Plus ver-388sion with more and harder test cases, the RL-tuned

model also has more than 3 points of improvement.

390

391

392

393

394

395

396

397

398

400

401

402

403

404

405

406

407

408

409

410

411

When starting from the Qwen2.5-Coder-Instruct-7B itself, we can still observe improvements, especially when using the rule-based reward. For example, we get more than 3.4 improvement on BigCodeBench-Full-Hard. Using the reward model for RL can also bring 3.2 improvement on MBPP. This highlights the charm of self-improvement given the reward model backbone is the same with the initial policy model. We compare our method with other RL-based models like RLEF (Chen et al., 2024), PPOCoder (Shojaee et al., 2023a), Step-Coder (Dou et al., 2024b), DSTC (Liu et al., 2024c), etc. We show that our 7B model is able to beat these competitors across the evaluation benchmarks.

Another experiment we conduct is to perform RL training directly from base model Qwen2.5-Coder-7B-base. We show significant improvement, especially through test-case pass rewards on HumanEval, MBPP, and BigCodeBench-I. These results are achieved by only training for 80 steps. We believe further scaling up the training will lead to much larger gains.

Comparison with Other RMsWe compare412our ACECODE-RM with 3 top-ranked RM on413the RewardBench, including InternLM2-RM-4148B (Cai et al., 2024), Skywork-Llama-3.1-8B, and415

Method & RM	Huma -	ınEval Plus	_ME _	3PP Plus	BigCoo Full	leBench-C Hard	BigCoo Full	deBench-I Hard	LiveCodeBench V4	Average
Greedy	68.9	62.2	67.2	54.8	38.5	12.8	31.8	13.5	18.0	40.9
Average	50.1	42.2	57.9	47.2	22.0	10.6	18.2	12.0	14.9	30.6
InternLM2-RM-8B	57.9	55.5	66.7	54.0	38.7	8.8	29.8	8.8	15.1	37.3
Skywork-Gemma-27B	73.8	67.1	64.3	53.4	40.1	14.9	32.5	12.8	23.6	42.5
Skywork-Llama-3.1-8B	67.7	61.6	69.6	56.9	40.6	10.8	31.8	12.2	18.8	41.1
∆ (max(other RM)-greedy)	+4.9	+4.9	+2.4	+2.1	+2.1	+2.0	+0.6	-0.7	+5.6	+2.6
ACECODE-RM-7B Δ (RM-greedy)	77.4 +8.5	70.7 +8.5	76.5 +9.3	64.3 +9.5	45.8 +7.3	20.3 +7.4	36.4 +4.6	12.2 -1.4	26.1 +8.1	47.7 +6.8

Table 4: ACECODE-RM's performance against other open-sourced reward models in terms of Best-of-16 sampling for Llama-3.1-8B-Inst. We can see the top-ranked RM on Reward Bench get little improvements compared to ours.

Skywork-Gemma-27B (Liu et al., 2024a), where re-416 sults are reported in Table 4. We can see that these 417 general-purpose RM can hardly improve and some-418 times decrease the performance through Best-of-N 419 sampling compared to greedy sampling, showcas-420 421 ing the incapability in identifying the correct generated programs. On the other hand, our ACECODE-422 RM surpasses all other publicly released reward 423 models in our evaluation and consistently gets pos-424 itive gains. These findings further underscore our 425 assumption that previous RM training lacks of reli-426 able signals for codes and prove that our RMs can 427 generate reliable and state-of-the-art reward signals 428 429 in code generation tasks.

4.5 Ablation Studies

430

431

432 433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

Test Case Quality Matters We also conduct experiments to investigate how filtering the test cases with a proxy model can affect the results. As shown in Table 5, training RM on data after the filtering improve the performance significantly, especially for those hard code questions like MBPP-Plus and BigCodeBench-Hard (C/I). We believe this is because the test case filtering can ensure the remaining ones are consistent with each other and thus point to the same implicit program, which improves the quality of the rewards.

RM Backbone Matters Our results in Table 6 clearly show that changing the backbone of the reward model from Llama-3.1 to Qwen2.5 can significantly improve the Best-of-16 performance. This is because the Qwen2.5-Coder models have been pretrained on way more code-related data compared to the Llama-3.1 models, and thus more knowledgeable when tuning it into a reward model.

450 Does R1-style Tuning Work? Inspired by the
451 recent DeepSeek-R1 (Guo et al., 2025), we also
452 conduct the RL directly from the base model with453 out any SFT. It turns out we get huge improve-

ments when using rule-based rewards. For example, we get 25.0 points of improvements on HumanEval-Plus after training only 6 hours from the Base Model, which is way more efficient that the large-scale SFT. What's more, the ACE-CODER *Rule* improve the BigCodeBench-Instruct-Full's performance from 40.2 to 43.2, nearly the same performance with DeepSeek-R1-Distill-Qwen-32B (43.9) which was directly distilled from the DeepSeek-R1 Model. This further consolidates the finding of DeepSeek-Zero. However, we do find that using reward models for RL tuning can lead to worse results. We attribute this to the potential reward hacking during the tuning.

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

486

487

488

489

490

491

5 Related Works

5.1 Synthesizing Test Cases

Automatic test generation is a widely used approach for verifying the correctness of LLMgenerated programs. Prior work has commonly employed the same LLM that generates the programs to also generate test cases, selecting the most consistent program from multiple sampled outputs in a self-consistency manner (Chen et al., 2022; Huang et al., 2023; Jiao et al., 2024). However, these generated test cases often suffer from significant hallucinations. To address this issue, Algo (Zhang et al., 2023) introduced the use of an oracle program solution to improve test case quality. While similar in spirit to our test case filtering approach, Algo constructs its oracle solution by exhaustively enumerating all possible combinations of relevant variables, whereas we leverage a stronger coder LLM to generate the oracle solution. Beyond using test cases as verification signals, Clover (Sun et al., 2023) enhances program verification by performing consistency checks between code, docstrings, and formal annotations, incorporating formal verification tools alongside LLMs.

Method	Huma -	anEval Plus	_MF	3PP Plus	BigCoo Full	leBench-C Hard	BigCoo Full	deBench-I Hard	LiveCodeBench V4	Average		
Inference Model = Llama-3.1-8B-Instruct												
RM w/o Test Case Filter RM w/ Test Filter Δ (w/ Filter - w/o Filter)	73.8 77.4 +3.7	65.9 70.7 +4.9	73.3 76.5 +3.2	61.4 64.3 +2.9	44.6 45.8 +1.2	17.6 20.3 +2.7	35.5 36.4 +0.9	9.5 12.2 +2.7	25.1 26.1 +1.0	45.2 47.7 +2.5		
Inference Model = Qwen2.5-Coder-7B-Instruct												
RM w/o Test Case Filter RM w/ Test Filter Δ (w/ Filter - w/o Filter)	91.5 90.2 -1.2	86.0 82.9 -3.0	86.0 88.6 +2.6	72.2 74.9 +2.6	52.5 53.8 +1.3	21.6 20.9 -0.7	43.4 45.0 +1.6	19.6 21.6 +2.0	36.9 40.1 +3.2	56.6 57.6 +1.0		

Table 5: Ablation study on test-case filtering. Results are Best-of-16 sampling performance.

Method	Huma -	nEval Plus	_ME _	BPP Plus	BigCoo Full	deBench-C Hard	BigCoo Full	leBench-I Hard	LiveCodeBench V4	Average		
Inference Model = Llama-3.1-8B-Instruct												
ACECODE-RM (LLama) ACECODE-RM (Qwen) Δ (Qwen-Llama)	65.9 77.4 +11.6	59.1 70.7 +11.6	69.6 76.5 +6.9	57.9 64.3 +6.3	42.7 45.8 +3.1	12.8 20.3 +7.4	32.9 36.4 +3.5	13.5 12.2 -1.4	19.9 26.1 +6.2	41.6 47.7 +6.1		
Inference Model = Qwen2.5-Coder-7B-Instruct												
ACECODE-RM (LLama) ACECODE-RM (Qwen) Δ (Qwen-Llama)	87.8 90.2 +2.4	81.7 82.9 +1.2	82.0 88.6 +6.6	67.7 74.9 +7.1	50.5 53.8 +3.2	25.0 20.9 -4.1	39.0 45.0 +6.0	19.6 21.6 +2.0	32.4 40.1 +7.7	54.0 57.6 +2.4		

Table 6: Comparison of ACECODE-RM's performance trained on different base model, where ACECODE-RM (Llama) is based on Llama-3.1-Inst-8B and ACECODE-RM (Qwen) is based on Qwen-Coder-2.5-7B-Inst. Results are Best-of-16 sampling performance.

493 494 495 496 497 498 499

492

5.2

504 505 506

507 508

510

511

512

513 514

515

516

517

518

Reinforcement Learning from Human Feedback (RLHF)(Ouyang et al., 2022b) has been widely

Reinforcement Learning for LLM

adopted to enhance the capabilities of large language models (LLMs) in various tasks, including conversational interactions and mathematical reasoning(Yang et al., 2024b). Reinforcement learning (RL) algorithms such as PPO(Schulman et al., 2017), GRPO(Shao et al., 2024), and Reinforcement++(Hu, 2025) have been employed to finetune models using reward signals derived from either learned reward models(Shao et al., 2024) or predefined rule-based heuristics (Guo et al., 2025; Wang et al., 2025).

Given that coding is an inherently verifiable task, recent studies have explored RL techniques that leverage direct execution accuracy as a reward signal. PPOCoder (Shojaee et al., 2023b) and CodeRL (Le et al., 2022) demonstrated the effectiveness of PPO-based RL for coding tasks, while RLEF (Gehring et al., 2024) extended this approach to multi-turn settings by incorporating execution feedback at each step. StepCoder (Dou et al., 2024a) refined the reward mechanism by assigning rewards at a more granular level, considering only successfully executed lines of code. Additionally, DSTC (Liu et al., 2024d) explored

the application of Direct Preference Optimization (DPO) to code generation by using self-generated test cases and programs.

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

Despite these advancements, most prior RLbased approaches for coding have been constrained by the use of pre-annotated datasets such as APPS (Hendrycks et al., 2021), which consists of only 5,000 examples, with most problems having a single test case. This limited data availability poses challenges to scalable RL training. Furthermore, the potential of reward models for coding remains largely unexplored. In this work, we address these limitations by automatically synthesizing test cases and leveraging trained reward models for reinforcement learning, demonstrating the scalability and effectiveness of our approach.

6 Conclusion

We introduced ACECODER as the first approach to reward model training and RL tuning for code generation using large-scale, reliable test case synthesis. Our data pipeline produces high-quality verifiable code without relying on the most advanced models, enabling effective reward model training and reinforcement learning. Our method significantly improves Best-of-N performance. However, RL training gains are less pronounced, leaving it as a future work to enhance.

Limitations

Test Case Synthesis Despite our efforts to enhance the reliability of synthesized test cases 548 through prompt engineering and filtering with a reference solution, inaccuracies can still arise. These errors may stem from an incorrect reference solu-551 552 tion or test cases that are too simple, failing to capture challenging edge cases. Consequently, passing all test cases does not necessarily guarantee a pro-554 gram's correctness, leading to noise in the reward model training and reinforcement learning (RL) 556 tuning signals. To address this, future work can 557 leverage stronger large language models (LLMs) to 558 synthesize more rigorous test cases, ensuring the in-560 clusion of harder corner cases. Additionally, using more advanced coding LLMs to generate reference solutions could further improve test case filtering, preserving only high-quality examples.

Reinforcement Learning for Coding In this paper, we explored RL tuning using three models and 565 two types of rewards: RM-based and rule-based. While significant improvements are observed when 567 tuning Qwen2.5-7B-Instruct and Qwen2.5-Coder-7B-Base, tuning on Qwen2.5-Coder-7B-Instruct 569 exhibited less pronounced gains due to it's strong 570 ability originally. This suggests that the current 571 reward signals may still contain noise. Furthermore, there remains considerable room for improvement, particularly in tuning the Qwen2.5-574 Coder-7B-Base. Given recent advancements in models such as DeepSeek-R1, future work could 576 further refine RL tuning strategies to achieve better performance with more fine-grained reward design.

Ethical Statements

This work fully complies with the ACL Ethics Policy. We declare that there are no ethical issues in this paper, to the best of our knowledge.

References

582

583

585

586

587

589

590

592

- Mistral AI. 2023. Mistral-7b-instruct-v0.3. https://huggingface.co/mistralai/ Mistral-7B-Instruct-v0.3.
- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. 2021.
 Program synthesis with large language models. arXiv preprint arXiv:2108.07732.
- Ralph Allan Bradley and Milton E. Terry. 1952. Rank

analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39:324.

- Zheng Cai, Maosong Cao, Haojiong Chen, Kai Chen, Keyu Chen, Xin Chen, Xun Chen, Zehui Chen, Zhi Chen, Pei Chu, Xiaoyi Dong, Haodong Duan, Qi Fan, Zhaoye Fei, Yang Gao, Jiaye Ge, Chenya Gu, Yuzhe Gu, Tao Gui, Aijia Guo, Qipeng Guo, Conghui He, Yingfan Hu, Ting Huang, Tao Jiang, Penglong Jiao, Zhenjiang Jin, Zhikai Lei, Jiaxing Li, Jingwen Li, Linyang Li, Shuaibin Li, Wei Li, Yining Li, Hongwei Liu, Jiangning Liu, Jiawei Hong, Kaiwen Liu, Kuikun Liu, Xiaoran Liu, Chengqi Lv, Haijun Lv, Kai Lv, Li Ma, Runyuan Ma, Zerun Ma, Wenchang Ning, Linke Ouyang, Jiantao Qiu, Yuan Qu, Fukai Shang, Yunfan Shao, Demin Song, Zifan Song, Zhihao Sui, Peng Sun, Yu Sun, Huanze Tang, Bin Wang, Guoteng Wang, Jiaqi Wang, Jiayu Wang, Rui Wang, Yudong Wang, Ziyi Wang, Xingjian Wei, Qizhen Weng, Fan Wu, Yingtong Xiong, Chao Xu, Ruiliang Xu, Hang Yan, Yirong Yan, Xiaogui Yang, Haochen Ye, Huaiyuan Ying, Jia Yu, Jing Yu, Yuhang Zang, Chuyu Zhang, Li Zhang, Pan Zhang, Peng Zhang, Ruijie Zhang, Shuo Zhang, Songyang Zhang, Wenjian Zhang, Wenwei Zhang, Xingcheng Zhang, Xinyue Zhang, Hui Zhao, Qian Zhao, Xiaomeng Zhao, Fengzhe Zhou, Zaida Zhou, Jingming Zhuo, Yicheng Zou, Xipeng Qiu, Yu Qiao, and Dahua Lin. 2024. Internlm2 technical report. Preprint, arXiv:2403.17297.
- Angelica Chen, Jérémy Scheurer, Jon Ander Campos, Tomasz Korbak, Jun Shern Chan, Samuel R Bowman, Kyunghyun Cho, and Ethan Perez. 2024. Learning from natural language feedback. *Transactions on Machine Learning Research*.
- Bei Chen, Fengji Zhang, A. Nguyen, Daoguang Zan, Zeqi Lin, Jian-Guang Lou, and Weizhu Chen. 2022. Codet: Code generation with generated tests. *ArXiv*, abs/2207.10397.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. 2021. Evaluating large language models trained on code. Preprint, arXiv:2107.03374.

595

596

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

763

764

Shihan Dou, Yan Liu, Haoxiang Jia, Limao Xiong, Enyu Zhou, Junjie Shan, Caishuang Huang, Wei Shen, Xiaoran Fan, Zhiheng Xi, Yuhao Zhou, Tao Ji, Rui Zheng, Qi Zhang, Xuanjing Huang, and Tao Gui. 2024a. Stepcoder: Improve code generation with reinforcement learning from compiler feedback. *ArXiv*, abs/2402.01391.

652

653

663

664

665

666

670

673

674

675

677

679

681

684

694

701

707

- Shihan Dou, Yan Liu, Haoxiang Jia, Limao Xiong, Enyu Zhou, Wei Shen, Junjie Shan, Caishuang Huang, Xiao Wang, Xiaoran Fan, et al. 2024b. Stepcoder: Improve code generation with reinforcement learning from compiler feedback. *arXiv preprint arXiv:2402.01391*.
- Jonas Gehring, Kunhao Zheng, Jade Copet, Vegard Mella, Taco Cohen, and Gabriele Synnaeve. 2024. Rlef: Grounding code llms in execution feedback with reinforcement learning. *ArXiv*, abs/2410.02089.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, and etc. 2024. The llama 3 herd of models. *Preprint*, arXiv:2407.21783.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, Yu Wu, Y. K. Li, Fuli Luo, Yingfei Xiong, and Wenfeng Liang. 2024a. Deepseek-coder: When the large language model meets programming - the rise of code intelligence. *ArXiv*, abs/2401.14196.
- Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, Yu Wu, YK Li, et al. 2024b. Deepseekcoder: When the large language model meets programming-the rise of code intelligence. *arXiv preprint arXiv:2401.14196*.
- Dan Hendrycks, Steven Basart, Saurav Kadavath, Mantas Mazeika, Akul Arora, Ethan Guo, Collin Burns, Samir Puranik, Horace He, Dawn Song, et al. 2021. Measuring coding challenge competence with apps. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track* (*Round 2*).
- Jian Hu. 2025. Reinforce++: A simple and efficient approach for aligning large language models. *arXiv* preprint arXiv:2501.03262.
- Jian Hu, Xibin Wu, Weixun Wang, Dehao Zhang, Yu Cao, OpenLLMAI Team, Netease Fuxi, AI Lab, and Alibaba Group. 2024. Openrlhf: An easy-touse, scalable and high-performance rlhf framework. *ArXiv*, abs/2405.11143.
 - Baizhou Huang, Shuai Lu, Weizhu Chen, Xiaojun Wan, and Nan Duan. 2023. Enhancing large language

models in coding through multi-perspective selfconsistency. In Annual Meeting of the Association for Computational Linguistics.

- Siming Huang, Tianhao Cheng, Jason Klein Liu, Jiaran Hao, Liuyihan Song, Yang Xu, J Yang, JH Liu, Chenchen Zhang, Linzheng Chai, et al. 2024. Open-coder: The open cookbook for top-tier code large language models. *arXiv preprint arXiv:2411.04905*.
- Binyuan Hui, Jian Yang, Zeyu Cui, Jiaxi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Kai Dang, An Yang, Rui Men, Fei Huang, Shanghaoran Quan, Xingzhang Ren, Xuancheng Ren, Jingren Zhou, and Junyang Lin. 2024a. Qwen2.5coder technical report. *ArXiv*, abs/2409.12186.
- Binyuan Hui, Jian Yang, Zeyu Cui, Jiaxi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Kai Dang, et al. 2024b. Qwen2. 5-coder technical report. *arXiv preprint arXiv:2409.12186*.
- Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. 2024. Gpt-40 system card. *arXiv preprint arXiv:2410.21276*.
- Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. 2024. Livecodebench: Holistic and contamination free evaluation of large language models for code. *arXiv preprint arXiv:2403.07974*.
- Fangkai Jiao, Geyang Guo, Xingxing Zhang, Nancy F. Chen, Shafiq Joty, and Furu Wei. 2024. Preference optimization for reasoning with pseudo feedback. *ArXiv*, abs/2411.16345.
- Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R Narasimhan. 2024. Swe-bench: Can language models resolve real-world github issues? In *The Twelfth International Conference on Learning Representations*.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*.
- Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, et al. 2024a. T\" ulu 3: Pushing frontiers in open language model post-training. *arXiv preprint arXiv:2411.15124*.
- Nathan Lambert, Valentina Pyatkin, Jacob Daniel Morrison, Lester James Validad Miranda, Bill Yuchen Lin, Khyathi Raghavi Chandu, Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, Noah A. Smith, and Hanna Hajishirzi. 2024b. Rewardbench: Evaluating reward models for language modeling. *ArXiv*, abs/2403.13787.

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

Hung Le, Yue Wang, Akhilesh Deepak Gotmare, Silvio Savarese, and Steven C. H. Hoi. 2022. Coderl: Mastering code generation through pretrained models and deep reinforcement learning. *ArXiv*, abs/2207.01780.

765

766

770

771

772

777

790

792

793

802

803

810

811

812

813

814

815

816

817

818 819

- Rongao Li, Jie Fu, Bo-Wen Zhang, Tao Huang, Zhihong Sun, Chen Lyu, Guang Liu, Zhi Jin, and Ge Li. 2023.
 Taco: Topics in algorithmic code generation dataset. *arXiv preprint arXiv:2312.14852.*
- Chris Yuhao Liu, Liang Zeng, Jiacai Liu, Rui Yan, Jujie He, Chaojie Wang, Shuicheng Yan, Yang Liu, and Yahui Zhou. 2024a. Skywork-reward: Bag of tricks for reward modeling in llms. *arXiv preprint arXiv:2410.18451*.
- Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. 2023. Is your code generated by chat-GPT really correct? rigorous evaluation of large language models for code generation. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Jiawei Liu, Songrun Xie, Junhao Wang, Yuxiang Wei, Yifeng Ding, and Lingming Zhang. 2024b. Evaluating language models for efficient code generation. In *First Conference on Language Modeling*.
- Zhihan Liu, Shenao Zhang, Yongfei Liu, Boyi Liu, Yingxiang Yang, and Zhaoran Wang. 2024c. Dstc: Direct preference learning with only self-generated tests and code to improve code lms. *arXiv preprint arXiv:2411.13611*.
- Zhihan Liu, Shenao Zhang, Yongfei Liu, Boyi Liu, Yingxiang Yang, and Zhaoran Wang. 2024d. Dstc: Direct preference learning with only self-generated tests and code to improve code lms. *ArXiv*, abs/2411.13611.
- Anton Lozhkov, Raymond Li, Loubna Ben Allal, Federico Cassano, Joel Lamy-Poirier, Nouamane Tazi, Ao Tang, Dmytro Pykhtar, Jiawei Liu, Yuxiang Wei, et al. 2024. Starcoder 2 and the stack v2: The next generation. *arXiv preprint arXiv:2402.19173*.
- Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xiubo Geng, Wenxiang Hu, Chongyang Tao, Jing Ma, Qingwei Lin, and Daxin Jiang. 2023. Wizardcoder: Empowering code large language models with evolinstruct.
- Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke E. Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Francis Christiano, Jan Leike, and Ryan J. Lowe. 2022a. Training language models to follow instructions with human feedback. *ArXiv*, abs/2203.02155.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al.

2022b. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744.

- Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, I. Evtimov, Joanna Bitton, Manish P Bhatt, Cristian Cantón Ferrer, Aaron Grattafiori, Wenhan Xiong, Alexandre D'efossez, Jade Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier, Thomas Scialom, and Gabriel Synnaeve. 2023. Code Ilama: Open foundation models for code. *ArXiv*, abs/2308.12950.
- John Schulman, Philipp Moritz, Sergey Levine, Michael I. Jordan, and P. Abbeel. 2015. Highdimensional continuous control using generalized advantage estimation. *CoRR*, abs/1506.02438.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *ArXiv*, abs/1707.06347.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Jun-Mei Song, Mingchuan Zhang, Y. K. Li, Yu Wu, and Daya Guo. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *ArXiv*, abs/2402.03300.
- Parshin Shojaee, Aneesh Jain, Sindhu Tipirneni, and Chandan K Reddy. 2023a. Execution-based code generation using deep reinforcement learning. *Transactions on Machine Learning Research*.
- Parshin Shojaee, Aneesh Jain, Sindhu Tipirneni, and Chandan K. Reddy. 2023b. Execution-based code generation using deep reinforcement learning. *ArXiv*, abs/2301.13816.
- Alexander Shypula, Aman Madaan, Yimeng Zeng, Uri Alon, Jacob Gardner, Milad Hashemi, Graham Neubig, Parthasarathy Ranganathan, Osbert Bastani, and Amir Yazdanbakhsh. 2023. Learning performance-improving code edits. *arXiv preprint arXiv:2302.07867*.
- Benjamin Steenhoek, Michele Tufano, Neel Sundaresan, and Alexey Svyatkovskiy. 2023. Reinforcement learning from automatic feedback for high-quality unit test generation. *arXiv preprint arXiv:2310.02368*.
- Chuyue Sun, Ying Sheng, Oded Padon, and Clark W. Barrett. 2023. Clover: Closed-loop verifiable code generation. In *SAIV*.
- Haozhe Wang, Long Li, Chao Qu, Fengming Zhu, Weidi Xu, Wei Chu, and Fangzhen Lin. 2025. Learning autonomous code integration for math lanuguage models. *ArXiv*.
- Yuxiang Wei, Zhe Wang, Jiawei Liu, Yifeng Ding, and Lingming Zhang. 2024. Magicoder: Empowering code generation with oss-instruct. In *Forty-first International Conference on Machine Learning*.

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

875

876

891

895

900

901 902

903

904

905

907

908

910

911

912 913

914

915

916

917

918

919

921

927

928

- An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu, Jianhong Tu, Jingren Zhou, Junyang Lin, Keming Lu, Mingfeng Xue, Runji Lin, Tianyu Liu, Xingzhang Ren, and Zhenru Zhang. 2024b. Qwen2.5-math technical report: Toward mathematical expert model via self-improvement. arXiv preprint arXiv:2409.12122.
- Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene Li, Qingning Yao, Shanelle Roman, et al. 2018. Spider: A large-scale human-labeled dataset for complex and cross-domain semantic parsing and text-to-sql task. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3911–3921.
- Kexun Zhang, Danqing Wang, Jingtao Xia, William Yang Wang, and Lei Li. 2023. Algo: Synthesizing algorithmic programs with generated oracle verifiers. *ArXiv*, abs/2305.14591.
- Zhenru Zhang, Chujie Zheng, Yangzhen Wu, Beichen Zhang, Runji Lin, Bowen Yu, Dayiheng Liu, Jingren Zhou, and Junyang Lin. 2025. The lessons of developing process reward models in mathematical reasoning. *arXiv preprint arXiv:2501.07301*.
- Tianyu Zheng, Ge Zhang, Tianhao Shen, Xueling Liu, Bill Yuchen Lin, Jie Fu, Wenhu Chen, and Xiang Yue. 2024a. Opencodeinterpreter: Integrating code generation with execution and refinement. *arXiv preprint arXiv:2402.14658*.
- Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyan Luo, Zhangchi Feng, and Yongqiang Ma. 2024b. Llamafactory: Unified efficient fine-tuning of 100+ language models. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations), Bangkok, Thailand. Association for Computational Linguistics.
- Terry Yue Zhuo, Minh Chien Vu, Jenny Chim, Han Hu, Wenhao Yu, Ratnadira Widyasari, Imam Nur Bani Yusuf, Haolan Zhan, Junda He, Indraneil Paul, et al. 2024. Bigcodebench: Benchmarking code generation with diverse function calls and complex instructions. arXiv preprint arXiv:2406.15877.

A Appendix

934

935

A.1 More related works

LLM for Code Generation Large language 936 models (LLMs) have demonstrated significant po-937 tential in code generation. Due to the unique na-938 ture of coding tasks, specialized coding models 939 such as Code Llama (Rozière et al., 2023) and Qwen Coder (Hui et al., 2024b; Yang et al., 2024a) were developed shortly after the emergence of 942 general-purpose LLMs. These models typically 943 undergo a two-phase training process: pre-training and fine-tuning. During pre-training, they are ex-945 posed to extensive coding corpora sourced from various internet platforms, including raw text, GitHub 947 repositories, and pull requests. This is followed 948 by supervised fine-tuning, which enhances their instruction-following capabilities. To assess the performance of these models in code generation, 951 several benchmarks have been established, includ-952 ing MBPP (Austin et al., 2021), HumanEval (Chen et al., 2021), EvalPlus (Liu et al., 2023, 2024b), Big Code Bench (Zhuo et al., 2024), and Live Code 955 Bench (Jain et al., 2024). These benchmarks usually include a series of prompts or problems for the 957 LLMs to solve, and they also contain test cases to assess the correctness of the generated code.

Reward Models Reward models play a crucial 960 961 role in aligning LLMs by assigning scalar values to response pairs based on specific evaluation cri-962 teria, such as human preference (Ouyang et al., 963 2022b) and accuracy (Zhang et al., 2025). They are widely used in reinforcement learning with human 965 feedback (RLHF) to refine model behavior and in 966 Best-of-N sampling to enhance test-time perfor-967 mance. However, while general-purpose reward models are effective for assessing human prefer-969 ence, they often struggle with specialized domains 970 like mathematics and coding due to the complex-971 ity of these tasks. For instance, even top-ranked reward models from Reward Bench (Lambert et al., 973 2024b), such as Skywork-RM (Liu et al., 2024a), 974 have difficulty providing reliable rewards for these 975 domains. To address this issue, task-specific reward 976 models have been developed, such as Qwen-2.5-Math-PRM (Zhang et al., 2025) for mathematical 978 reasoning. However, coding reward models have 979 remained largely absent due to the lack of reli-980 able training signals—an issue that our proposed ACECODE-RM aims to address. 982

A.2 Prompt Template

system:

You are an AI assistant that helps people with python coding tasks.

user:

You are the latest and best bot aimed at transforming some code snippet into a leetcode style question. You will be provided with a prompt for writing code, along with a reference program that answers the question. Please complete the following for me:

1. Come up with a leetcode style question which consists of a well-defined problem. The generated question should meet the following criteria:

a. The question is clear and understandable, with enough details to describe what the input and output are.

b. The question should be solvable by only implementing 1 function instead of multiple functions or a class. Therefore, please avoid questions which require complicated pipelines.

c. The question itself should not require any access to external resource or database.

d. Feel free to use part of the original question if necessary. Moreover, please do not ask for runtime and space complexity analysis or any test cases in your response.

2. Based on the modified question that you generated in part 1, you need to create around 20 test cases for this modified question. Each test case should be independent assert clauses. The parameters and expected output of each test case should all be constants, **without accessing any external resources**.

Here is the original question: {instruction}

Here is the reference program that answers the question:

```python
{program}

Now give your modified question and generated test cases in the following json format: {"question": ..., "tests":["assert ...", "assert ..."]}.

Table 7: Prompt Used for Converting Seed Code Dataset into LeetCode-style Questions and Test Cases

### system:

You are an AI assistant that helps people with python coding tasks.

# user:

You are the latest and best bot aimed at transforming some code snippet into a leetcode style question. You will be provided with a reference program. Please complete the following for me:

1. Come up with a leetcode style question which consists of a well-defined problem. The generated question should meet the following criteria:

a. The question is clear and understandable, with enough details to describe what the input and output are.

b. The question should be solvable by only implementing 1 function instead of multiple functions or a class. Therefore, please avoid questions which require complicated pipelines.

c. The question itself should not require any access to external resource or database.

d. Feel free to use part of the original question if necessary. Moreover, please do not ask for runtime and space complexity analysis or any test cases in your response.

2. Based on the modified question that you generated in part 1, you need to create around 20 test cases for this modified question. Each test case should be independent assert clauses. The parameters and expected output of each test case should all be constants, \*\*without accessing any external resources\*\*.

Here is the reference program:

``` python
{program}

•••

Table 8: Prompt Used for Converting Seed Code Dataset using only the reference program without instruction into LeetCode-style Questions and Test Cases