

CodeRise: Bootstrapping LLMs for Ultra Low-Resource Programming Languages via Progressive Self-Refinement Curriculum

Anonymous ACL submission

Abstract

Large Language Models (LLMs) struggle with code generation for Ultra Low-Resource Programming Languages (ULRPLs) due to the scarcity of training data. Existing synthetic data generation methods fail in this context, suffering from a severe cold-start problem and resulting in samples that lack diversity. To overcome these challenges, we propose CodeRise, a novel two-stage framework that autonomously generates a high-quality, diverse, and progressively complex curriculum for ULRPLs. The framework first tackles the cold-start and distribution issues by leveraging the full formal syntax of the target language as structural guidance and applying a biased sampling strategy over library modules. Building on this foundation, we fine-tune the model to generate increasingly complex code without explicit syntax input, using an adaptive curriculum and multi-turn self-debugging to progressively improve code quality. We evaluate on two ULRPLs, Tingo and Janet, using migrated HumanEval-Tingo and MBPP-Tingo, as well as our new benchmarks, TingoEval and JanetEval. Experiments show that CodeRise significantly outperforms both training-free and training-based baselines in ultra low-resource environments.

1 Introduction

Large Language Models (LLMs) (OpenAI et al., 2024b; Hui et al., 2024; Guo et al., 2024) have demonstrated outstanding performance across a wide range of code-related tasks, including code generation (Chen et al., 2021; Lai et al., 2023), code repair (Chen et al., 2024; Zhong et al., 2024; Tian et al., 2024), code completion (Bavarian et al., 2022; Ding et al., 2024), and program analysis (First et al., 2023). This remarkable success has been predominantly observed in the context of high-resource languages such as Python, where vast and high-quality codebases are readily available for training. The success is even beginning to extend

to many Low-Resource Programming Languages (LRPLs) such as R and Perl (Cassano et al., 2022; Giagnorio et al., 2025; Joel et al., 2024; Cassano et al., 2024). However, the powerful capabilities of LLMs do not universally extend to the Ultra Low-Resource Programming Languages (Mora et al., 2024), despite their importance in various specialized domains. These languages suffer from multi-dimensional data scarcity: the volume of publicly available code is inherently low, and systematic collection is also challenging. For instance, many ULRPLs lack unique file extensions, making them invisible to automatic identification tools such as GitHub Linguist¹. This severe lack of data directly limits the ability of models to learn the syntax, libraries, and coding patterns specific to ULRPLs.

A promising direction to overcome this data scarcity is to leverage the LLM itself to generate synthetic training data, a principle exemplified by methods like *Self-Instruct* (Wang et al., 2023) and *SelfCodeAlign* (Wei et al., 2024a). The Self-Instruct paradigm typically operates by iteratively bootstrapping instructions from a seed corpus, prompting a powerful teacher model to generate corresponding responses, and then validating these generations. However, existing relevant techniques, which are primarily designed for high-resource languages, exhibit limitations when directly applied to ULRPLs. A primary obstacle for ULRPLs is the cold-start problem. Even with state-of-the-art LLMs, their initial proficiency in these languages is often so low that they fail to produce even a single syntactically valid code sample, which prevents synthetic data generation from bootstrapping.

A straightforward workaround to launch this process is to provide the model with the complete language syntax in context. However, this approach is fraught with issues: the extremely long context

¹<https://github.com/github-linguist/linguist>

Module	Frequency	Proportion(%)
text	6202	47.8
math	1438	11.1
rand	1291	10.0
fmt	1354	10.4
json	935	7.2
enum	791	6.1
times	650	5.0
base64	284	2.2
os	17	0.1
hex	9	0.1

Table 1: Distribution of standard library usage from an 30,000 Tingo samples dataset generated with the SelfCodeAlign baseline. The statistics are based on the 12,971 samples that involved API calls.

082 leads to slow inference and a low pass rate, as
083 shown in Figure 1, making it impractical for large-
084 scale data generation. Furthermore, even with this
085 syntax guidance, the generated data still exhibits
086 a severe long-tail distribution, clustering around
087 basic language features, as shown in Table 1. Such
088 long-tail coverage issues have also been observed
089 for self-generated instruction data in high-resource
090 languages (Wang et al., 2025). In our ULRPL set-
091 ting, we find that the effect is further amplified
092 because the model defaults to highly transferable
093 “universal programming tasks” (e.g., string manipu-
094 lation).

095 To overcome these challenges, we propose
096 **CodeRise**, a two-stage framework for ULRPL
097 code generation. Stage 1 (*Syntax-Guided Module-*
098 *Balanced Generation*) bootstraps from near-zero
099 ability by providing syntax guidance and balanc-
100 ing sampling across language-specific libraries
101 to obtain a syntactically valid and diverse seed
102 dataset. Stage 2 (*Progressive Self-Refined Genera-*
103 *tion*) removes the long syntax prompt and improves
104 both difficulty and quality via an adaptive cur-
105 riculum and multi-turn self-debugging. Together,
106 CodeRise progressively strengthens the model and
107 produces high-quality ULRPL code under extreme
108 data scarcity.

109 We conduct extensive experiments on two rep-
110 resentative ULRPLs, **Tengo** and **Janet**. For each
111 language, we use CodeRise to synthesize a high-
112 quality dataset and fine-tune LLMs on the corre-
113 sponding synthetic data. For Tengo, we evaluate
114 on migrated versions of HumanEval-Tengo and
115 MBPP-Tengo, along with our newly constructed
116 **TengoEval** benchmark. For Janet, we evaluate

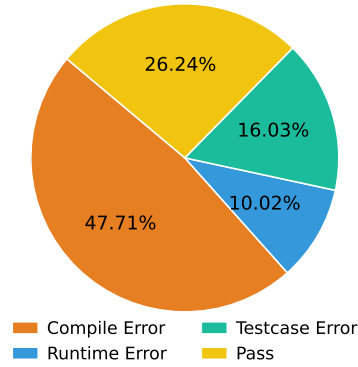


Figure 1: Failure analysis for the Qwen-2.5-Coder-32B-Instruct using in-context learning with the language’s full documentation, the majority of failures are execution errors.

on our newly constructed **JanetEval** benchmark. Results show that models fine-tuned on CodeRise data consistently and significantly outperform a comprehensive suite of baselines, including both training-free and training-based approaches, underscoring the effectiveness of our progressive data generation strategy.

Our contributions are as follows:

- We propose **CodeRise**, a self-improvement framework designed to overcome the cold-start and long-tail challenges in ULRPLs.
- We release benchmark suites for **Tengo** and **Janet**. For Tengo, we provide migrated versions of HumanEval and MBPP together with our newly constructed **TengoEval**; for Janet, we release our newly constructed **JanetEval**.

2 Related Work

LLMs for Coding In recent years, Large Language Models (Guo et al., 2024; Rozière et al., 2024) have achieved remarkable progress in code-related tasks, such as code generation (Zheng et al., 2024; Han et al., 2024), code completion (Yu et al., 2024a; Hayes et al., 2025) and code translation (Szafraniec et al., 2023; Xue et al., 2024; Pan et al., 2024), and the impact is expanding across the entire software engineering (Jimenez et al., 2023; Xie et al., 2025; Yang et al., 2024). Besides, reasoning models like OpenAI-o1 (OpenAI et al., 2024a) and DeepSeek-R1 (DeepSeek-AI et al., 2025) have also demonstrated impressive coding performance due to their long reasoning ability. However, de-

148 spite these advancements, current LLMs still struggle
149 to perform well on ULRPLs due to the scarcity
150 of training data, limited exposure to syntax pat-
151 terns, and the lack of task-specific supervision.
152 While some prior work has addressed challenges
153 in low-resource programming languages (Giagnorio
154 et al., 2025; Mora et al., 2024; Cassano et al.,
155 2022, 2024), such as R and Perl, their methods rely
156 on the model’s existing initial capability in those
157 languages, this foundation is absent for the ULR-
158 PLs setting. In addition, many of these approaches
159 synthesize training data via code translation from
160 high-resource languages, which can introduce task-
161 domain shift.

Data-Driven Code Tuning Instruction tuning
162 further trains a pretrained LLM on instruction-
163 response pairs (Ouyang et al., 2022; Wei et al.;
164 Li et al., 2024) and has proven effective for code
165 tasks (Muennighoff et al., 2024; Yu et al., 2024b;
166 Liu et al., 2024). However, high-quality instruction
167 data is costly to curate. Consequently, the *Self-*
168 *Instruct* paradigm was proposed to generate syn-
169 thetic instruction data from seed sets (Wang et al.,
170 2023; Xu et al., 2024) and has been widely adopted
171 in code (Luo et al., 2024; Wei et al., 2024b), as well
172 as frameworks such as SelfCodeAlign (Wei et al.,
173 2024a). Moreover, self-generated instruction data
174 can be overly concentrated in a few domains, moti-
175 vating feature-aware rebalancing methods such
176 as Epicoder (Wang et al., 2025). However, these
177 approaches typically assume sufficient initial com-
178 petence in the target language, an assumption that
179 often fails in ULRPLs.
180

181 3 Methodology

182 To address the challenges of data scarcity and the
183 pervasive long-tail distribution problem inherent in
184 ULRPLs, we propose an iterative self-improvement
185 framework. This framework adheres to a curricu-
186 lum learning paradigm, enabling an LLM to au-
187 tonomously generate diverse, correct, and progres-
188 sively complex fine-tuning data through two dis-
189 tinct stages. Stage 1, termed Syntax-Guided Data
190 Generation, primarily aims to create an initial high-
191 quality dataset for the preliminary fine-tuning of
192 the LLM’s foundational understanding. Stage 2,
193 Advanced Data Generation via Self-Debugging, is
194 meticulously engineered to achieve two key objec-
195 tives: generating more complex and challenging
196 tasks efficiently, and substantially improving data
197 quality through robust self-correction mechanisms.

198 3.1 Syntax-Guided Module-Balanced 199 Generation

200 The data generation process in Stage 1 is executed
201 as a multi-step pipeline designed to create a founda-
202 tional and diverse dataset. The pipeline begins by
203 generating a composite “seed concept set” derived
204 from both code snippet analysis and a diversity-
205 focused module sampling strategy. This seed set
206 is then used to formulate a specific programming
207 instruction. Following this, the model leverages the
208 complete language syntax to autonomously gener-
209 ate both a code implementation and its correspond-
210 ing test suite. Finally, the entire script is rigorously
211 validated in a sandbox environment to ensure its
212 correctness.

Concept Inspiration from Code Snippets Simi-
213 lar to previous works like OSS-Instruct (Wei et al.,
214 2024b) and SelfCodeAlign (Wei et al., 2024a), we
215 select a code snippet from a meticulously curated
216 open-source code corpus, typically ranging from
217 1 to 15 lines. This chosen snippet serves as an in-
218 spiration for the LLM. Prompted with this snippet,
219 the LLM is tasked with extracting and formulating
220 a set of related “programming concept” that repre-
221 sent the core functionality, algorithms, or common
222 patterns found in the snippet. This approach lever-
223 ages the LLM’s understanding of diverse coding
224 patterns to derive high-level, actionable concepts,
225 which form the building blocks for subsequent task
226 instruction generation.
227

Module Sampling Relying solely on concepts
228 derived from code snippets leads to a long-tail dis-
229 tribution in the generated dataset, with tasks cluster-
230 ing around common language features like string
231 manipulation. To counteract this bias, we employ a
232 weighted sampling strategy to select a module from
233 the language’s standard library, then this module
234 is injected as an additional concept to guide task
235 generation. The sampling weight for each module
236 is determined by our logarithmic inverse frequency
237 formula:
238

$$239 w_i = \frac{1}{\log(c_i + 2)^p} \quad (1)$$

240 where c_i is the frequency that module i appears in
241 correct solution. p is a hyperparameter to control
242 the intensity of this weighting.

Code Generation and Validation The seed con-
243 cept set from the previous step serves as the in-
244 put to a sequential pipeline that generates and val-
245 idates each data point. First, this set is embedded
246

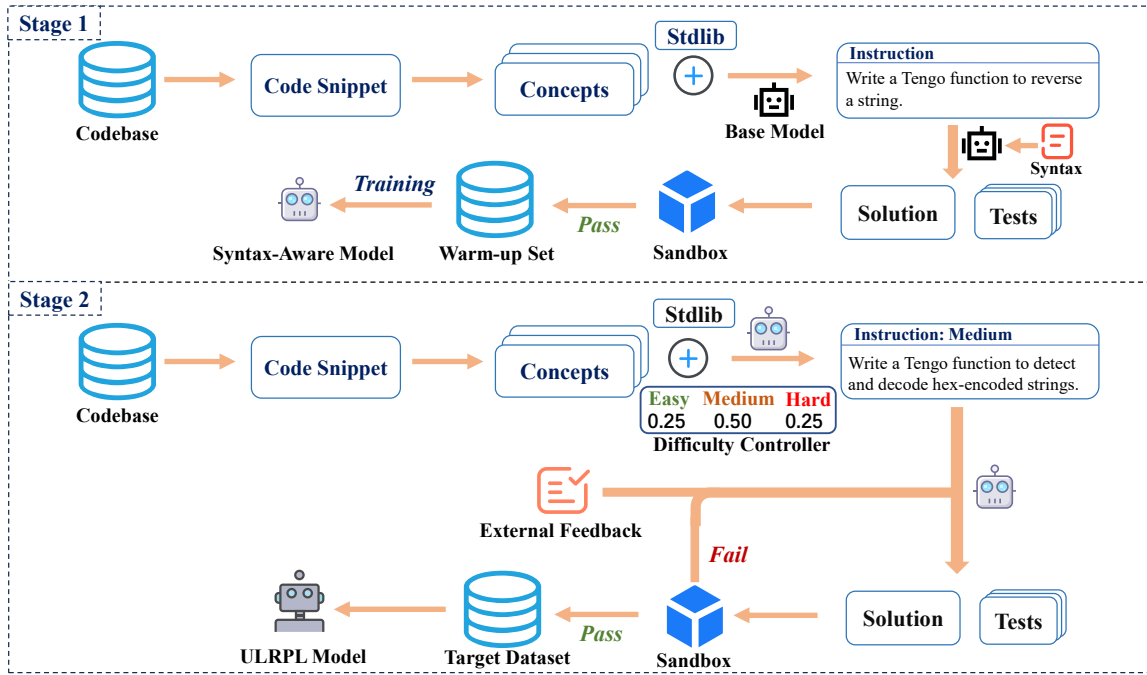


Figure 2: The overview of our two-stage CodeRise framework. **Stage 1: Syntax-Guided Module-Balanced Generation** (Top) addresses the cold-start problem by using explicit syntax guidance and diversity-focused module sampling to create a Warm-up Dataset. **Stage 2: Progressive Self-Refined Generation** (Bottom) then builds upon the fine-tuned model, employing a curriculum-based difficulty controller and a multi-turn self-debugging guided by external feedback to efficiently generate complex, high-quality instruction data for the final ULRPL model.

247 into a prompt that instructs the LLM to formulate
 248 a specific programming task in natural language.
 249 Next, this newly created instruction, along with the
 250 complete language syntax as a guiding scaffold, is
 251 provided to the LLM in a second prompt. In this
 252 step, the model’s task is to generate both the code
 253 solution and a corresponding set of unit tests. Fi-
 254 nally, each resulting (instruction, code, test) triplet
 255 is subjected to rigorous automated validation in a
 256 secure sandbox. An instruction-code pair is only
 257 accepted into our cold-start dataset if the code exe-
 258 cutes flawlessly and passes all associated tests, all
 259 other attempts are discarded.

260 3.2 Progressive Self-Refined Generation

261 Building upon the model fine-tuned in Stage 1,
 262 this advanced stage elevates the data generation
 263 process along two key dimensions: generation effi-
 264 ciency and data quality. A significant leap in
 265 efficiency is achieved by no longer providing the
 266 complete language syntax in the prompt. This is
 267 viable because the model has already internalized
 268 foundational syntactic knowledge, which drasti-
 269 cally reduces context length and improves through-
 270 put. To enhance quality and generation yield, this
 271 stage replaces Stage 1’s “discard-on-failure” strat-

272 egy with the multi-turn self-debugging mechanism.
 273 This strategy enables the model to methodically
 274 correct its errors by leveraging reliable external
 275 feedback, facilitating the successful generation of
 276 more complex instruction-code pairs and elevating
 277 the dataset’s overall quality.

Curriculum-Based Difficulty Control To sys-
 278 tematically increase the complexity of generated
 279 tasks, Stage 2 introduces a structured curriculum
 280 learning approach. While the initial seed gener-
 281 ation process remains consistent with Stage 1 com-
 282 bining concepts from code snippets with diversity-
 283 focused module sampling, we now introduce a cru-
 284 cial new dimension: explicit difficulty levels.
 285

286 Crucially, we categorize code tasks into three dif-
 287 ficulty tiers: Easy, Medium, and Hard, each with
 288 a tailored prompt template. The curriculum’s pro-
 289 gression is adaptive, driven by the model’s mastery
 290 as quantified by its Pass@10 score p ($0 \leq p \leq 1$)
 291 on recent tasks. This score dynamically adjusts the
 292 sampling probability for each tier. The unnormal-
 293 ized weights are calculated as follows:

$$\begin{aligned}
 \text{Weight}_{\text{hard}} &= \alpha \cdot p \\
 \text{Weight}_{\text{medium}} &= C \\
 \text{Weight}_{\text{easy}} &= \alpha \cdot (1 - p)
 \end{aligned}
 \tag{2}$$

where α is a hyperparameter balancing the influence of Pass@10 on hard/easy tasks, and p represents the current Pass@10. The Medium tier’s weight is a fixed constant C that acts as a stable baseline, ensuring the model always has a default difficulty to fall back on.

These unnormalized weights are then converted into probabilities using a softmax function with temperature T .

$$P_i = \frac{\exp(\text{Weight}_i/T)}{\sum_j \exp(\text{Weight}_j/T)} \quad (3)$$

where P_i is the probability of selecting difficulty tier i , temperature T controls the randomness of the tier selection. This adaptive strategy ensures a progressive increase in task complexity, maintaining appropriate control over the curriculum.

Multi-Turn Self-Debugging A cornerstone of Stage 2 is the Multi-Turn Self-Debugging Mechanism, designed to significantly increase both the success rate and quality of data generation. In a crucial departure from the “discard-on-failure” policy in Stage 1, the process begins after an initial solution fails validation in the sandbox. Instead of being discarded, the system first analyzes the specific type of error to trigger one of two tailored debugging strategies, which are detailed as follows:

Debugging Compilation or Runtime Errors. When a failure is caused by a compilation or runtime error, we employ a feedback-driven correction strategy. A new debugging prompt is constructed for the LLM, containing three key sources of information: 1) the verbatim error message from the compiler or runtime; 2) a pre-authored Debugging Guideline with expert strategies for common language-specific issues; 3) relevant chunks from the documentation using BM25 retrieval. By synthesizing these combined sources of feedback, the LLM is tasked with diagnosing the root cause of the error and generating a revised code solution.

Debugging Unit Test Failures via Pseudo Ground Truth. For logical flaws, we employ a distinct cross-lingual debugging strategy with the core objective of transferring the model’s strong reasoning capabilities from a high-resource language (Python) to our target language. To achieve this, the model is prompted to generate a Python solution that serves as a “Pseudo Ground Truth”. This Python solution is generated only upon the first occurrence of a logical error for a given task. The LLM then performs a comparative analysis to

identify and correct logical discrepancies in its low-resource code, with the process concluding once a revised version successfully passes all unit tests.

Iterative Refinement. The two debugging strategies described above both operate within an iterative refinement loop. After each correction attempt, the revised code is resubmitted to the sandbox for re-validation, creating a “debug-correct-revalidate” cycle. This cycle continues until either the code successfully passes all tests or a predefined maximum number of attempts is reached. Stage 2 enables the successful creation of more complex instruction-code pairs that would otherwise be filtered out, thereby elevating the overall quality and challenge level of the final dataset.

4 Experiments

4.1 Experimental Setup

This section details the models, parameters, and evaluation benchmarks used to validate our proposed framework. Our research primarily focuses on enhancing code generation capabilities for ULRPLs. For this purpose, we selected **Tengo** and **Janet** as the target ULRPLs. Tengo follows a paradigm similar to high-resource languages, whereas Janet is a stack-based language with a substantially different programming style, they represent two contrasting ULRPL settings under extreme data scarcity.

Model We use the Qwen-2.5-Coder-32B (Hui et al., 2024) as the data generator to create the training datasets in both Stage 1 and Stage 2. The models selected for fine-tuning are Qwen-2.5-Coder-32B, Qwen-2.5-Coder-7B-Instruct (Hui et al., 2024) and Llama-3.1-8B-Instruct (Grattafiori et al., 2024).

Baseline We compare against a comprehensive set of training-free and training-based baselines. The training-free methods include standard Zero-Shot and Few-Shot prompting, as well as richer-context strategies: (i) **Full-Syntax Prompting**, which provides the entire language manual in context; (ii) **Translation Rules** (Giagnorio et al., 2025); and (iii) **Translation Example** (Giagnorio et al., 2025). For training-based comparisons, we include three representative baselines: (i) **Self-CodeAlign** (Wei et al., 2024a), a general self-alignment approach that fine-tunes a model on self-generated instruction data; (ii) **MultiPL-T** (Casano et al., 2024), which synthesizes training data

Model	Method	HumanEval	MBPP	TengoEval	JanetEval
		Pass@1	Pass@1	Pass@1	Pass@1
Qwen-Coder-2.5-32B-Instruct	Zero-Shot	0.61	2.34	1.00	0.00
	Full-Syntax Prompting	33.5	36.8	23.0	6.00
	Few-Shot Examples	22.6	41.8	15.0	4.00
	Translation Examples	24.4	42.6	18.0	4.00
	Translation Rules	15.9	32.0	11.0	2.00
	SelfCodeAlign	61.6	65.1	49.0	24.0
	MultiPL-T	60.4	64.6	47.0	22.0
	SPEAC	60.4	63.5	45.0	-
	CodeRise (Ours)	71.3	75.5	60.0	36.0
Qwen-Coder-2.5-7B-Instruct	Zero-Shot	0.61	0.00	1.00	0.00
	Full-Syntax Prompting	12.2	19.0	14.0	4.00
	Few-Shot Examples	9.75	25.7	9.00	2.00
	Translation Examples	4.27	6.80	6.00	2.00
	Translation Rules	2.44	10.4	7.00	2.00
	SelfCodeAlign	39.6	41.7	31.0	12.0
	MultiPL-T	43.3	44.5	31.0	14.0
	SPEAC	45.1	45.3	30.0	-
	CodeRise (Ours)	50.3	57.9	39.0	20.0
Llama-3.1-8B-Instruct	Zero-Shot	2.44	3.13	2.00	0.00
	Full-Syntax Prompting	9.76	10.4	10.0	2.00
	Few-Shot Examples	15.2	29.0	11.0	2.00
	Translation Examples	14.6	30.4	7.00	2.00
	Translation Rules	3.67	2.86	6.00	0.00
	SelfCodeAlign	32.9	33.6	29.0	10.0
	MultiPL-T	25.0	28.6	21.0	8.00
	SPEAC	27.4	31.5	20.0	-
	CodeRise (Ours)	35.4	36.7	31.0	18.0
DeepSeek-R1-0528	Zero-Shot	18.3	27.8	16.0	6.00
	Full-Syntax Prompting	67.1	75.5	56.0	22.0

Table 2: Pass@1 comparison of CodeRise against training-free baselines and training-based baselines on migrated HumanEval-Tengo and MBPP-Tengo, as well as ULRPL-specific benchmarks, TengoEval and JanetEval.

for low-resource languages via code translation and test-based filtering; and (iii) **SPEAC** (Mora et al., 2024), which first prompts the LLM to write code in an intermediate language and then compiles it into the target language to construct training data.

Datasets and Metric To evaluate our framework, we employ several benchmarks. For Tengo, since no standard benchmark exists, we migrate HumanEval (Chen et al., 2021) and MBPP (Austin et al., 2021), and construct a new benchmark, **TengoEval**. For Janet, we use our newly annotated benchmark, **JanetEval**. Detailed benchmark construction and statistics are provided in Appendix B. We report Pass@1 (Chen et al., 2021) as the primary metric, which measures the percentage of problems solved in a single attempt.

4.2 Results

To evaluate the effectiveness of CodeRise, we conducted a comprehensive set of experiments, with

Module	Frequency	Proportion(%)
text	6746	43.7 / -4.1
math	1562	10.1 / -1.1
rand	1373	8.9 / -1.1
fmt	1676	10.9 / +0.5
json	1057	6.8 / -0.4
enum	820	5.3 / -0.8
times	692	4.5 / -0.5
base64	504	3.3 / +1.1
os	652	4.2 / +4.1
hex	352	2.3 / +2.2

Table 3: Distribution of standard module usage from our 30,000 dataset, the statistics are based on the 15,434 samples that involved API calls. The last column showing the change compared to the baseline in Table 1.

the main results summarized in Table 2. From these results, we can derive the following conclusions:

CodeRise significantly outperforms both training-free and strong training-based baselines across all tested open-source models. As

shown in Table 2, when applied to the powerful Qwen-Coder-2.5-32B-Instruct model, our method outperforms all training-free and strong training-based baselines. Furthermore, its performance is highly competitive with the powerful DeepSeek-R1-0528 model, especially considering the latter relies on a Full-Syntax Prompting strategy that is exceptionally time-consuming due to its long context requirements. This trend of superior performance extends to the smaller models as well, on both Qwen-Coder-2.5-7B-Instruct and Llama-3.1-8B-Instruct, CodeRise consistently outperforms all baselines, demonstrating the stability.

Our framework effectively mitigates the long-tail distribution while preserving natural usage patterns. As demonstrated in Table 3, our method significantly rebalances the module usage distribution compared to the baseline (e.g., *os*, *hex*). Crucially, this is achieved without artificially distorting the overall distribution. This demonstrates that our framework addresses the most extreme aspects of the long-tail problem, resulting in a more diverse yet practical training dataset.

CodeRise significantly improves the efficiency and overall yield of the data generation process. As detailed in Table 4, CodeRise (Stage 2) achieves a final effective pass rate of 62.86%, a $2.4\times$ increase in yield over the SelfCodeAlign baseline. **Pass Rate** measures generation yield, namely the fraction of generation attempts that produce a syntactically valid solution and pass all unit tests after the method’s refinement and filtering procedure. It is worth noting that the initial pass rate of CodeRise (Stage 1) is slightly lower than the baseline, which is an expected trade-off because our diversity sampling deliberately explores more challenging, long-tail tasks. In contrast, the translation-based baseline MultiPL-T exhibits a substantially lower yield in this setting, since the domain shift between the source language and the target ULRPL causes many translated samples to be invalid or not faithfully expressible, and thus fail compilation. Ultimately, CodeRise solves these harder problems through self-debugging, leading to a superior outcome in both data quality and overall efficiency.

4.3 Analysis

4.3.1 Ablation

To better understand the contribution of each component in our CodeRise framework, we conduct a series of ablation studies on the HumanEval dataset.

Method	Syntax	Pass Rate (%)
MultiPL-T	✓	12.03
SelfCodeAlign	✓	26.24
CodeRise (Stage 1)	✓	23.64
CodeRise (Stage 2)	×	62.86

Table 4: Data generation efficiency and yield on Qwen-2.5-Coder-32B-Instruct for Tengo. We run 30,000 generation attempts for each method and report Pass Rate. The **Syntax** column indicates whether the method includes a full language syntax prompt during generation.

Method	HumanEval	TengoEval
CodeRise	71.3	60.0
w/o Diversity Strategy	70.7	49.0
w/o Curriculum Learning	64.0	52.0
w/o Self-Debugging	67.1	55.0
w/o Retrieval	69.5	53.0
Stage 1 Only (30k data)	59.8	46.0
Stage 1 Only (2k data)	51.8	36.0

Table 5: Ablation study results (Pass@1) on HumanEval and TengoEval using the fine-tuned Qwen-2.5-Coder-32B-Instruct model.

The results are presented in Table 5. We define the following ablation settings: a) Stage 1 Only: The model is fine-tuned on the dataset generated through the Stage 1 method; b) w/o Diversity Strategy: Remove the module sampling component in CodeRise; c) w/o Curriculum Learning: The full CodeRise framework is used, but the adaptive curriculum in Stage 2 is disabled. Tasks are generated without any explicit difficulty guidance or tier-specific prompting; d) w/o Self-Debugging: The full CodeRise framework is used, but the multi-turn self-debugging mechanism in Stage 2 is disabled, and failed attempts are discarded; e) w/o Retrieval: The full CodeRise framework is used, but remove the retrieval component in Stage 2. For a fair comparison, all models were fine-tuned on 30,000 samples. Based on the above ablation study, we can derive the following conclusions:

The Stage 2 self-refinement process is the primary driver of performance. The results in Table 5 clearly show a massive performance gap between our full framework and the Stage 1 Only model. We attribute this significant difference to the quality of the training data. The Stage 1 Only dataset is generated using a simple “discard-on-failure” policy without any curriculum or difficulty control. This inherently biases the dataset towards

simpler problems that the model can solve in a single pass. In contrast, our Stage 2 process can generate and successfully solve much more complex problems, and training on this higher-quality data is what ultimately leads to the superior performance of our final model.

Adaptive curriculum and multi-turn self-debugging are all effective components. The adaptive curriculum proves to be the most impactful component, its removal (w/o Curriculum Learning) causes a substantial performance degradation. This highlights that systematically generating more challenging problems is essential for pushing the model beyond its initial capabilities. The self-debugging mechanism is the crucial counterpart to this process. Its removal (w/o Self-Debugging) leads to a 4.2% decrease in performance on HumanEval, demonstrating that providing the model with multiple attempts and reliable external feedback is what enables it to successfully solve the difficult tasks proposed by the curriculum. This aligns with our core philosophy: for ULRPLs where no powerful teacher model exists, it is more effective to introduce external, verifiable knowledge to guide a model’s own refinement process. For details on self-debugging, please refer to the case study in Appendix C.

The diversity strategy and retrieval component mainly contribute to long-tail coverage rather than algorithmic problem solving. Without the module-balanced sampling strategy in Stage 1, performance remains nearly unchanged on HumanEval but drops sharply on TengeEval, demonstrating that balancing over language-specific modules is critical for improving standard-library coverage and mitigating long-tail bias. Similarly, removing retrieval only slightly affects HumanEval but leads to a larger degradation on TengeEval, suggesting that external documentation is particularly important for resolving library-related details in practical ULRPL usage. Overall, these ablations confirm that CodeRise relies on complementary mechanisms: Stage 1 diversifies and bootstraps generation, while Stage 2 progressively increases difficulty and refines solutions to produce higher-quality and more practical training data.

4.3.2 Sweet Spot of Warm-up

An important consideration for CodeRise is determining the optimal amount of Stage 1 data for the initial fine-tuning. The objective is to achieve two goals: (1) effectively reducing execution errors in

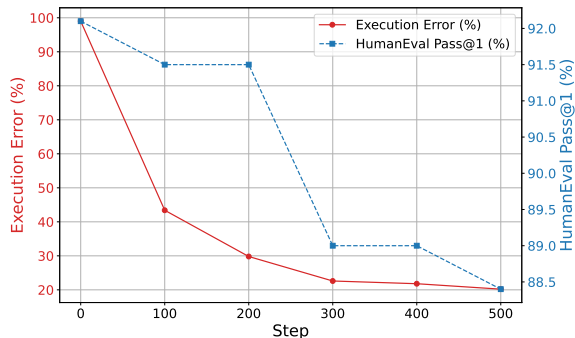


Figure 3: Identifying the sweet spot for the Stage 1 syntax warm-up. We plot the Execution Error Rate (%) on Tenge (left) and the Pass@1 on HumanEval-Python (right) against the number of fine-tuning steps.

the target language; (2) preserving the model’s general ability, it supports our aim of using a single model to save significant computational resources. To find this balance, we tracked the model’s performance at 100-step intervals (batch size = 8). At each checkpoint, we evaluated its Execution Error Rate on a unified set of 5k instructions, and the performance on HumanEval-Python. As shown in Figure 3, the execution error rate on Tenge drops from nearly 100% to under 30% in the first 200 training steps. The intersection of these trends reveals a "sweet spot" around 200-300 steps. Therefore, we identify this range as the optimal warm-up point, as it achieves a substantial reduction in execution errors at a minimal cost to the model’s general capabilities.

5 Conclusion

In this paper, we addressed the critical challenge of applying Large Language Models to ULRPLs, where progress has been hindered by a lack of high-quality training data. We identified two primary obstacles for existing self-generation methods: a severe cold-start problem and the tendency to produce datasets with a long-tail distribution. We introduced CodeRise, a novel two-stage framework that overcomes these issues. By first using syntax-guided generation and a diversity-focused sampling strategy, and then refining the model with an adaptive curriculum and a multi-turn self-debugging mechanism, CodeRise generates a diverse, complex, and high-quality dataset. Our experiments demonstrate that CodeRise significantly outperforms a wide range of baselines, proving the effectiveness of our methodology.

579
580
581
582
583
584
585
586
587
588
589
590

591

592
593
594
595
596

597
598
599
600
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

Limitations

Our study has a key limitation that points to promising directions for future research. Our current instruction generation process does not account for the domain-specific of many ULRPLs. While our generated data is syntactically valid, it may not reflect the idiomatic or practical use cases of the target language. Future work should focus on developing mechanisms to align instruction generation with the intended application scope, thereby improving the real-world relevance of the synthetic dataset.

References

Jacob Austin, Augustus Odena, Maxwell I. Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie J. Cai, Michael Terry, Quoc V. Le, and Charles Sutton. 2021. [Program synthesis with large language models](#). *CoRR*, abs/2108.07732.

Mohammad Bavarian, Heewoo Jun, Nikolas Tezak, John Schulman, Christine McLeavey, Jerry Tworek, and Mark Chen. 2022. [Efficient training of language models to fill in the middle](#). *Preprint*, arXiv:2207.14255.

Federico Cassano, John Gouwar, Francesca Lucchetti, Claire Schlesinger, Anders Freeman, Carolyn Jane Anderson, Molly Q. Feldman, Michael Greenberg, Abhinav Jangda, and Arjun Guha. 2024. [Knowledge transfer from high-resource to low-resource programming languages for code llms](#). *Proc. ACM Program. Lang.*, 8(OOPSLA2):677–708.

Federico Cassano, John Gouwar, Daniel Nguyen, Sydney Nguyen, Luna Phipps-Costin, Donald Pinckney, Ming-Ho Yee, Yangtian Zi, Carolyn Jane Anderson, Molly Q Feldman, Arjun Guha, Michael Greenberg, and Abhinav Jangda. 2022. [Multipl-e: A scalable and extensible approach to benchmarking neural code generation](#). *Preprint*, arXiv:2208.08227.

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Pondé de Oliveira Pinto, Jared Kaplan, and 1 others. 2021. [Evaluating large language models trained on code](#). *CoRR*, abs/2107.03374.

Xinyun Chen, Maxwell Lin, Nathanael Schärli, and Denny Zhou. 2024. [Teaching large language models to self-debug](#). In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net.

DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, and 1 others. 2025. [Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning](#). *Preprint*, arXiv:2501.12948.

Yifeng Ding, Hantian Ding, Shiqi Wang, Qing Sun, Varun Kumar, and Zijian Wang. 2024. [Horizon-length prediction: Advancing fill-in-the-middle capabilities for code generation with lookahead planning](#). *CoRR*, abs/2410.03103.

Emily First, Markus N. Rabe, Talia Ringer, and Yuriy Brun. 2023. [Baldur: Whole-proof generation and repair with large language models](#). In *Proceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering, ESEC/FSE 2023, San Francisco, CA, USA, December 3-9, 2023*, pages 1229–1241. ACM.

Alessandro Giagnorio, Alberto Martin-Lopez, and Gabriele Bavota. 2025. [Enhancing code generation for low-resource languages: No silver bullet](#). *Preprint*, arXiv:2501.19085.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, and 1 others. 2024. [The llama 3 herd of models](#). *Preprint*, arXiv:2407.21783.

Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, and 1 others. 2024. [Deepseek-coder: When the large language model meets programming – the rise of code intelligence](#). *Preprint*, arXiv:2401.14196.

Hojae Han, Jaejin Kim, Jaeseok Yoo, Youngwon Lee, and Seung-won Hwang. 2024. [Archcode: Incorporating software requirements in code generation with large language models](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, pages 13520–13552. Association for Computational Linguistics.

Jamie Hayes, Iliia Shumailov, William P. Porter, and Aneesh Pappu. 2025. [Measuring memorization in RLHF for code completion](#). In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025*. OpenReview.net.

Binyuan Hui, Jian Yang, Zeyu Cui, Jiayi Yang, and Others. 2024. [Qwen2.5-coder technical report](#). *CoRR*, abs/2409.12186.

Carlos E. Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. 2023. [Swe-bench: Can language models resolve real-world github issues?](#) *CoRR*, abs/2310.06770.

Sathvik Joel, Jie JW Wu, and Fatemeh H. Fard. 2024. [A survey on llm-based code generation for low-resource and domain-specific programming languages](#). *Preprint*, arXiv:2410.03981.

Denis Kocetkov, Raymond Li, Loubna Ben Allal, Jia Li, Chenghao Mou, Yacine Jernite, Margaret Mitchell, Carlos Muñoz Ferrandis, Sean Hughes, Thomas Wolf, Dzmitry Bahdanau, Leandro von Werra, and Harm

685	de Vries. 2023. The stack: 3 TB of permissively licensed source code . <i>Trans. Mach. Learn. Res.</i> , 2023.	
686		
687	Yuhang Lai, Chengxi Li, Yiming Wang, Tianyi Zhang,	
688	Ruiqi Zhong, Luke Zettlemoyer, Wen-Tau Yih,	
689	Daniel Fried, Sida I. Wang, and Tao Yu. 2023. DS-1000: A natural and reliable benchmark for data science code generation . In <i>International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA</i> , volume 202 of <i>Proceedings of Machine Learning Research</i> , pages 18319–18345. PMLR.	
690		
691		
692		
693		
694		
695		
696	Kaixin Li, Qisheng Hu, James Xu Zhao, Hui Chen,	
697	Yuxi Xie, Tiedong Liu, Michael Shieh, and Junxian	
698	He. 2024. InstructCoder: Instruction tuning large language models for code editing . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics, ACL 2024 - Student Research Workshop, Bangkok, Thailand, August 11-16, 2024</i> , pages 50–70. Association for Computational Linguistics.	
699		
700		
701		
702		
703		
704		
705	Jimmy Lin, Xueguang Ma, Sheng-Chieh Lin, Jheng-	
706	Hong Yang, Ronak Pradeep, and Rodrigo Nogueira.	
707	2021. Pyserini: A Python toolkit for reproducible	
708	information retrieval research with sparse and dense	
709	representations. In <i>Proceedings of the 44th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2021)</i> , pages 2356–2362.	
710		
711		
712		
713	Bingchang Liu, Chaoyu Chen, Zi Gong, Cong Liao,	
714	Huan Wang, Zhichao Lei, Ming Liang, Dajun Chen,	
715	Min Shen, Hailian Zhou, Wei Jiang, Hang Yu, and	
716	Jianguo Li. 2024. MftCoder: Boosting code llms with multitask fine-tuning . In <i>Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, KDD 2024, Barcelona, Spain, August 25-29, 2024</i> , pages 5430–5441. ACM.	
717		
718		
719		
720		
721	Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and	
722	LINGMING ZHANG. 2023. Is your code generated by chatGPT really correct? rigorous evaluation of large language models for code generation . In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> .	
723		
724		
725		
726		
727	Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xi-	
728	ubo Geng, Wenxiang Hu, Chongyang Tao, Jing Ma,	
729	Qingwei Lin, and Daxin Jiang. 2024. WizardCoder: Empowering code large language models with evolve-instruct . In <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . OpenReview.net.	
730		
731		
732		
733		
734	Federico Mora, Justin Wong, Haley Lepe, Sahil Bha-	
735	tia, Karim Elmaaroufi, George Varghese, Joseph E.	
736	Gonzalez, Elizabeth Polgreen, and Sanjit A. Seshia.	
737	2024. Synthetic programming elicitation for text-to-code in very low-resource programming and formal languages . In <i>The Thirty-eighth Annual Conference on Neural Information Processing Systems</i> .	
738		
739		
740		
741	Niklas Muennighoff, Qian Liu, Armel Randy Ze-	
742	baze, Qinkai Zheng, Binyuan Hui, Terry Yue Zhuo,	
	Swayam Singh, Xiangru Tang, Leandro von Werra,	743
	and Shayne Longpre. 2024. Octopack: Instruction tuning code large language models . In <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . OpenReview.net.	744
		745
		746
		747
		748
	OpenAI, :, Aaron Jaech, Adam Kalai, Adam Lerer, and	749
	1 others. 2024a. Openai o1 system card . <i>Preprint</i> , arXiv:2412.16720.	750
		751
	OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal,	752
	Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,	753
	and 1 others. 2024b. Gpt-4 technical report . <i>Preprint</i> , arXiv:2303.08774.	754
		755
	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida,	756
	Carroll L. Wainwright, Pamela Mishkin, Chong	757
	Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray,	758
	John Schulman, Jacob Hilton, Fraser Kelton, Luke	759
	Miller, Maddie Simens, Amanda Askell, Peter Welin-	760
	der, Paul F. Christiano, Jan Leike, and Ryan Lowe.	761
	2022. Training language models to follow instructions with human feedback . In <i>Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022</i> .	762
		763
		764
		765
		766
		767
	Rangeet Pan, Ali Reza Ibrahimzada, Rahul Krishna,	768
	Divya Sankar, Lambert Pougues Wassi, Michele	769
	Merler, Boris Sobolev, Raju Pavuluri, Saurabh Sinha,	770
	and Reyhaneh Jabbarvand. 2024. Lost in translation: A study of bugs introduced by large language models while translating code . In <i>Proceedings of the 46th IEEE/ACM International Conference on Software Engineering, ICSE 2024, Lisbon, Portugal, April 14-20, 2024</i> , pages 82:1–82:13. ACM.	771
		772
		773
		774
		775
		776
		777
	Stephen Robertson and Hugo Zaragoza. 2009. The probabilistic relevance framework: Bm25 and beyond . <i>Foundations and Trends® in Information Retrieval</i> , 3(4):333–389.	778
		779
		780
	Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten	781
	Sootla, Itai Gat, Xiaoqing Ellen Tan, and 1 others.	782
	2024. Code llama: Open foundation models for code . <i>Preprint</i> , arXiv:2308.12950.	783
		784
	Marc Szafraniec, Baptiste Rozière, Hugh Leather,	785
	Patrick Labatut, François Charton, and Gabriel Syn-	786
	naeve. 2023. Code translation with compiler representations . In <i>The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net.	787
		788
		789
		790
	Runchu Tian, Yining Ye, Yujia Qin, Xin Cong, Yankai	791
	Lin, Yinxu Pan, Yesai Wu, Hui Haotian, Liu We-	792
	ichuan, Zhiyuan Liu, and Maosong Sun. 2024. DebugBench: Evaluating debugging capability of large language models . In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 4173–4198, Bangkok, Thailand. Association for Computational Linguistics.	793
		794
		795
		796
		797
		798

799	Yaoxiang Wang, Haoling Li, Xin Zhang, Jie Wu, Xiao Liu, Wenxiang Hu, Zhongxin, and 1 others. 2025. Epicoder: Encompassing diversity and complexity in code generation . In <i>Forty-second International Conference on Machine Learning, ICML 2025, Vancouver, BC, Canada, July 13-19, 2025</i> . OpenReview.net.	Xinran Yu, Chun Li, Minxue Pan, and Xuandong Li. 2024a. Droidcoder: Enhanced android code completion with context-enriched retrieval-augmented generation . In <i>Proceedings of the 39th IEEE/ACM International Conference on Automated Software Engineering, ASE 2024, Sacramento, CA, USA, October 27 - November 1, 2024</i> , pages 681–693. ACM.	856 857 858 859 860 861 862
805	Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2023. Self-instruct: Aligning language models with self-generated instructions . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 13484–13508. Association for Computational Linguistics.	Zhaojian Yu, Xin Zhang, Ning Shang, Yangyu Huang, Can Xu, Yishujie Zhao, Wenxiang Hu, and Qiufeng Yin. 2024b. Wavocoder: Widespread and versatile enhancement for code large language models by instruction tuning . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024</i> , pages 5140–5153. Association for Computational Linguistics.	863 864 865 866 867 868 869 870 871
814	Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. Finetuned language models are zero-shot learners . In <i>International Conference on Learning Representations</i> .	Lin Zheng, Jianbo Yuan, Zhi Zhang, Hongxia Yang, and Lingpeng Kong. 2024. Self-infilling code generation . <i>Preprint</i> , arXiv:2311.17972.	872 873 874
819	Yuxiang Wei, Federico Cassano, Jiawei Liu, Yifeng Ding, Naman Jain, Zachary Mueller, Harm de Vries, Leandro von Werra, Arjun Guha, and Lingming Zhang. 2024a. Selfcodealign: Self-alignment for code generation . In <i>Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024</i> .	Li Zhong, Zilong Wang, and Jingbo Shang. 2024. Debug like a human: A large language model debugger via verifying runtime execution step by step . In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 851–870, Bangkok, Thailand. Association for Computational Linguistics.	875 876 877 878 879 880
828	Yuxiang Wei, Zhe Wang, Jiawei Liu, Yifeng Ding, and Lingming Zhang. 2024b. Magicoder: Empowering code generation with oss-instruct . In <i>Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024</i> . OpenReview.net.		
834	Chengxing Xie, Bowen Li, Chang Gao, He Du, Wai Lam, Difan Zou, and Kai Chen. 2025. Swe-fixer: Training open-source llms for effective and efficient github issue resolution . <i>CoRR</i> , abs/2501.05040.		
838	Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, Qingwei Lin, and Daxin Jiang. 2024. Wizardlm: Empowering large pre-trained language models to follow complex instructions . In <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . OpenReview.net.		
845	Min Xue, Artur Andrzejak, and Marla Leuther. 2024. An interpretable error correction method for enhancing code-to-code translation . In <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . OpenReview.net.		
851	John Yang, Carlos E. Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik Narasimhan, and Ofir Press. 2024. Swe-agent: Agent-computer interfaces enable automated software engineering . <i>Preprint</i> , arXiv:2405.15793.		

881	A Implementation Details	
882	This section provides the implementation details of	
883	CodeRise to ensure full reproducibility. We report	
884	the key hyperparameters for our data generation	
885	pipeline. For all inference processes conducted	
886	during testing and evaluation, the sampling temper-	
887	ature was consistently set to $T = 0$.	
888	A.1 Syntax Warm-up	
889	For the initial syntax warm-up, we fine-tuned the	
890	Qwen-2.5-Coder-32B-Instruct model on a dataset	
891	of 2k samples generated from Stage 1. The training	
892	was conducted for 1 epoch, with a learning rate of	
893	2×10^{-5} , and the batch size is set to 8. We utilized	
894	a cosine scheduler and AdamW optimizer for the	
895	training process.	
896	A.2 Main Model Training	
897	We fine-tuned three models on the 30,000 samples	
898	dataset generated by CodeRise. All training runs	
899	were conducted for 3 epochs using the AdamW	
900	optimizer with a cosine learning rate scheduler and	
901	a maximum sequence length of 2048 tokens. No	
902	warm-up steps were used. For Qwen-2.5-Coder-	
903	32B-Instruct, we used a learning rate of 2×10^{-5}	
904	with a batch size of 8. For Qwen-2.5-Coder-7B-	
905	Instruct and Llama-3.1-8B-Instruct, we used a	
906	learning rate of 4×10^{-5} with a batch size of 32.	
907	A.3 Seed CodeBase	
908	To ensure a fair comparison with prior work, we	
909	use the same seed codebase ² (Wei et al., 2024a) as	
910	the SelfCodeAlign framework, which consists of	
911	574k functions. This dataset is a curated collec-	
912	tion of functions originally sourced from The Stack	
913	V1 ³ (Kocetkov et al., 2023) corpus, a 3.1 TB dataset	
914	consisting of permissively licensed source code in	
915	30 programming languages.	
916	B Benchmark Details	
917	We evaluate CodeRise on migrated benchmarks	
918	for Tingo (HumanEval and MBPP) and our newly	
919	constructed native benchmark suite for ultra low-	
920	resource programming languages, consisting of	
921	TengoEval and JanetEval .	
	² https://huggingface.co/datasets/bigcode/python-stack-v1-functions-filtered	
	³ https://huggingface.co/datasets/bigcode/the-stack	
	B.1 TengoEval and JanetEval	922
	Scope and format TengoEval and JanetEval	923
	are file-level benchmarks, each containing 100	924
	problems. Unlike translation-based benchmarks,	925
	both are natively designed for their target lan-	926
	guages and grounded in each language’s standard	927
	library, rather than translated from high-resource	928
	languages.	929
	Module coverage TengoEval is designed to	930
	cover all native modules in the Tingo standard li-	931
	brary, and each problem uses 2.3 modules on av-	932
	erage. JanetEval contains 100 problems and each	933
	problem uses 1.6 modules on average.	934
	Unit tests and semantics All test cases are writ-	935
	ten and executed under the native semantics of the	936
	target language. Each task includes multiple unit	937
	tests and at least one explicit edge case to reduce	938
	ambiguity and improve robustness.	939
	Quality control To ensure correctness and fea-	940
	sibility, two senior engineers reviewed every prob-	941
	lem for task feasibility and test correctness. All	942
	problems, reference solutions, and unit tests are ex-	943
	ecuted and validated in a sandboxed environment	944
	before inclusion.	945
	B.2 Migrated Benchmarks	946
	We evaluate on migrated versions of HumanEval	947
	and MBPP for Tingo. We use the EvalPlus (Liu	948
	et al., 2023) problem sets as the source, but due to	949
	limited staffing, we only migrate the original unit	950
	tests provided by HumanEval and MBPP rather	951
	than the additional tests introduced by EvalPlus.	952
	All migrated problems and tests are manually veri-	953
	fied by two engineers to ensure correctness under	954
	native Tingo semantics.	955
	C Case Study	956
	Figure 4 presents a case study where the model was	957
	tasked with calculating Euclidean distance. The ini-	958
	tial code contained a compilation error, it used tuple	959
	assignment to unpack coordinates, which is a com-	960
	mon feature in Go but unsupported in Tingo. Upon	961
	failure, it receives three complementary pieces of	962
	information: the direct compiler error identifies	963
	what is wrong; the pre-authored Debugging Guide-	964
	line explains why it is wrong in Tingo and provides	965
	the correct pattern; and the retrieved documentation	966
	confirms that tuple assignment is a feature explic-	967
	itly absent from the language. This multi-faceted	968

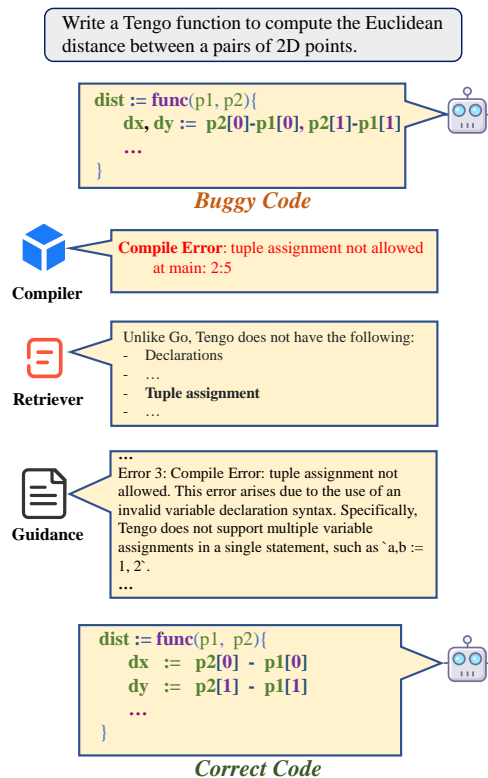


Figure 4: A case study of the CodeRise mechanism. An Compile Error caused by tuple assignment is resolved after the model receives synergistic feedback from the compiler, a documentation retriever, and our pre-authored guidance.

feedback enables the model to effectively diagnose and resolve the language-specific syntactic issue.

D Document Retrieval Strategy

The document was divided into three structurally parts: **(i) Basic Syntax**; **(ii) Built-in Functions**; and **(iii) Standard Libraries**. The chunking strategy was tailored to the characteristics of each section. For **Basic Syntax** and **Built-in Functions**, we leveraged the native Markdown hierarchy, treating each top-level heading as an individual chunk. In contrast, the **Standard Libraries** section, which consists of numerous concise function descriptions, was uniformly partitioned by grouping every three consecutive function definitions into one chunk. This design ensures semantic coherence within chunks while maintaining consistent retrieval granularity across heterogeneous documentation types.

Since retrieval is only triggered for *compilation* or *runtime* errors, BM25(Robertson and Zaragoza, 2009) proves more effective than dense retrieval. We utilized the BM25 algorithm from

Iterations	Samples	Proportion (%)
1	11,090	53.1
2	4,476	21.4
3	2,150	10.3
4	1,194	5.7
5	751	3.6
6	510	2.4
7	395	1.9
8	319	1.5
9	257	1.2

Table 6: Iteration statistics of Stage 2 self-debugging on Tengo. We report the distribution over the number of refinement iterations for 20,885 successfully salvaged samples.

the Pyserini(Lin et al., 2021) with its default settings. The feedback messages from compilers often contain explicit lexical cues that directly correspond to documentation text, such as “Runtime Error: wrong number of arguments in call to user-function:has_prefix”. In such cases, keyword-based retrieval better exploits these literal overlaps, whereas dense embeddings tend to smooth away the precise token-level distinctions that are crucial for error localization and resolution. Consequently, BM25 provides more reliable and interpretable matches for this type of debugging-oriented retrieval.

E Iteration statistics

To characterize the efficiency of Stage 2, we analyze 20,885 samples that were successfully salvaged by our self-debugging pipeline on Tengo. Table 6 reports the distribution over the number of refinement iterations needed to pass compilation/execution and all unit tests. More than half of the samples succeed in a single iteration; 74.5% succeed within two iterations, and 90.5% within four iterations. Only a small long tail (10.5%) requires five or more iterations, indicating that Stage 2 is typically efficient while still being able to rescue hard cases when needed.