TAROT: TEST-DRIVEN AND CAPABILITY-ADAPTIVE CURRICULUM REINFORCEMENT FINE-TUNING FOR CODE GENERATION

Anonymous authorsPaper under double-blind review

000

001

002

004

006

008 009 010

011 012

013

014

015

016

017

018

019

021

023

025

026

027

028

029

031

033

034

037

040

041

042

043

044

046

047

048

051

052

ABSTRACT

Large Language Models (LLMs) are fundamentally changing the coding paradigm, known as vibe coding, yet synthesizing algorithmically sophisticated and robust code still remains a critical challenge. Incentivizing the deep reasoning capabilities of LLMs is essential to overcome this hurdle. Reinforcement Fine-Tuning (RFT) has emerged as a promising strategy to address this need. However, most existing approaches overlook the heterogeneous difficulty and granularity inherent in test cases, leading to an imbalanced distribution of reward signals and consequently biased gradient updates during training. To address this, we propose "TAROT", Test-driven and cApability-adaptive cuRriculum reinfOrcement fine-Tuning. TAROT systematically constructs, for each problem, a four-tier test suite (basic, intermediate, complex, edge), providing a controlled difficulty landscape for curriculum design and evaluation. Crucially, TAROT decouples curriculum progression from raw reward scores, enabling capability-conditioned evaluation and principled selection from a portfolio of curriculum policies rather than incidental test-case difficulty composition. This design fosters stable optimization and more efficient competency acquisition. Extensive experimental results reveal that the optimal curriculum for reinforcement fine-tuning in code generation is closely tied to a model's inherent capability, with less capable models achieving greater gains with an easy-to-hard progression, whereas more competent models excel under a hard-first curriculum. TAROT provides a reproducible method that adaptively tailors curriculum design to a model's capability, thereby consistently improving the functional correctness and robustness of the generated code. All code and data are released to foster reproducibility and advance community research at https://anonymous.4open.science/r/TAROT-B675/.

1 Introduction

Large Language Models (LLMs) are driving significant changes in software engineering, with automated code generation emerging as a pivotal application (Du et al., 2024; Jiang et al., 2024). Foundational models exhibit a strong capacity to translate natural language specifications into functional code, promising significant enhancements in developer productivity (Weber et al., 2024). Nevertheless, advancing the frontier toward synthesizing algorithmically sophisticated and highly robust solutions remains a critical challenge (Zhuo et al., 2025). The next key step hinges on significantly augmenting the deep reasoning and problem-solving faculties of these models.

Curriculum Learning (CL), a methodology that structures training data by difficulty (Bengio et al., 2009), presents a promising avenue for cultivating these capabilities and improving training efficiency. However, existing applications of CL in code generation primarily focus on sequencing entire problems based on coarse difficulty metrics (Naïr et al., 2024; Khant et al., 2025). While this inter-problem curriculum is intuitive, it neglects the nuanced, intra-problem difficulty gradient inherent in software verification. Human developers naturally employ practices like Test-Driven Development (TDD) (Beck, 2003), incrementally refining a solution against increasingly complex test cases to ensure robustness. Yet, this natural curriculum axis remains largely untapped in LLM training. Furthermore, reliance on problem-level sequencing often leads to flat reward landscapes when integrated with Reinforcement Fine-Tuning (RFT), dampening the learning signal. This oversight

of heterogeneous test-case difficulty results in imbalanced reward signals and consequently biased gradient updates during training, hindering the model's ability to acquire robust, sophisticated reasoning skills.

Furthermore, while the trend in curriculum learning for LLMs is shifting towards more dynamic approaches that progressively increase task complexity (Xu et al., 2024; Cheng et al., 2025), these methods predominantly define difficulty based on the intrinsic properties of the data or the task structure. For instance, curricula are often structured using automated metrics of the source code itself, such as cyclomatic complexity (Naïr et al., 2024), or by decomposing a problem into a fixed sequence of simpler subtasks (Dou et al., 2024). This prevailing focus on the data, rather than the learner, overlooks the crucial variable of the model's own evolving and multi-faceted capability. A curriculum tailored to an early-stage model may cause learning stagnation for a more advanced one, while a curriculum designed for experts can overwhelm a nascent model and hinder its convergence. Therefore, for a more holistic approach to effective learning, curriculum design should consider not only the intrinsic properties of the data but also the evolving capabilities of the model itself, leading to a capability-adaptive framework.

To address these limitations, we introduce **TAROT**, a novel framework for **T**est-driven and c**A**pability-adaptive cu**R**riculum reinf**O**rcement fine-**T**uning. Crucially, TAROT decouples curriculum progression from raw reward scores, enabling capability-conditioned evaluation and principled selection from a portfolio of curriculum policies rather than incidental test-case difficulty composition. This design fosters stable optimization and more efficient competency acquisition. The framework's novelty is twofold. First, TAROT operationalizes the concept of an intra-problem difficulty gradient through a novel, test-driven curriculum. To instantiate this gradient, which is absent in standard coding datasets, we constructed the TAROT dataset. Each coding problem is systematically augmented with a test suite built upon four tiers of difficulty including basic, intermediate, complex, and edge cases. This structure defines difficulty as a spectrum of functional correctness. This engineered gradient directly counteracts the flat reward landscape common in RL, providing a structured and nuanced signal for learning robust solutions.

Second, we study capability-adaptive curriculum design. Given the TAROT dataset, we instantiate a portfolio of curriculum policies that vary along three axes, namely allocation across tiers, the sequence and proportion of tiers, and reward weighting across tiers. This setup enables capability-conditioned evaluation and principled selection among policies for models differing in effective capability which is defined via instruction-following fidelity and baseline coding proficiency, both influenced by model scale, specialization, and architecture. Our central thesis is that the optimal learning path is capability dependent. In particular, Nascent models learn best with basic to complex progression, whereas stronger models learn best by focusing on complex tiers.

To validate this thesis, this paper introduces the TAROT framework and demonstrates its effectiveness through comprehensive experiments, making the following primary contributions:

- A novel intra-problem, test driven curriculum for code generation, embodied in the new TAROT dataset. Each problem in our dataset is augmented with a four-tiered test suite (basic, intermediate, complex, edge cases) to provide a granular difficulty landscape and enable nuanced reward modeling.
- A capability conditioned study and guideline for curriculum design in code generation.
 We present a portfolio of curriculum policies in allocation sequence and reward weighting together with a reproducible evaluation protocol for capability conditioned comparison and principled selection informed by the characteristics of model capability such as scale and specialization.
- Comprehensive empirical validation demonstrating the effectiveness of TAROT. Our experiments show that our capability-adaptive approach significantly improves model performance and training efficiency compared to baseline curriculum strategies.

2 Related works

Our work is positioned at the intersection of two key research domains: curriculum learning for structuring pedagogical data, and reinforcement learning for policy optimization in code generation. We review prior work in these areas and highlight how TAROT offers a novel synthesis of both.

2.1 Curriculum Learning for Code

109 110 111

112

113

114

115

116

117

118

119

120

121

122

123

124

108

Curriculum Learning (CL) is a training strategy inspired by human cognition that presents data to a model in a structured order, typically from simple to complex examples (Bengio et al., 2009). This method has been shown to accelerate convergence and improve generalization by guiding optimization toward better solutions. In the context of Large Language Models (LLMs), curricula have been implemented in various ways, such as using a teacher model to progressively generate more complex instructions, as seen in the Evol-Instruct method (Xu et al., 2024), or by fine-tuning on a small set of meticulously curated, high-quality examples as demonstrated by LIMA (Zhou et al., 2023). For code generation, where task complexity varies widely, CL is a particularly promising but challenging area. While many code datasets rely on manual difficulty labels, recent research has focused on more systematic approaches. A notable example is the use of automatic difficulty metrics, combining measures like cyclomatic complexity and Halstead difficulty, to sort problems into a multi-stage curriculum (Naïr et al., 2024). Training with this structured approach yielded significant gains, demonstrating the value of CL in the code domain. Other methods, like StepCoder, create an implicit curriculum by breaking a complex problem into a sequence of simpler code-completion subtasks (Dou et al., 2024). These efforts show a clear trend towards leveraging curricula to organize the training process for code generation. Our work contributes to this line of research by proposing a novel method to generate a tiered test suite that serves as the basis for our curriculum.

125 126 127

2.2 Reinforcement Fine-Tuning for Code LLMs

128129130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

Reinforcement Learning (RL) is a dominant paradigm for aligning LLMs with desired behaviors, using techniques like RLHF (Ouyang et al., 2022), DPO (Rafailov et al., 2023) (Rafailov et al., 2023), PPO (Schulman et al., 2017) (Schulman et al., 2017), GRPO (Shao et al., 2024), and GSPO (Zheng et al., 2025). In code generation, RL is adapted to optimize for functional correctness, typically using unit test outcomes as a reward signal. This "RL from unit test feedback" approach, while effective, often suffers from two key limitations: a sparse and "flat" reward landscape. The reward is sparse because a model gets no learning signal on complete failure, and it is flat because all successful solutions receive the same reward, regardless of the problem's difficulty. This flatness, where all successful solutions receive a similar reward regardless of the challenge, generates imbalanced reward signals that can lead to biased gradient updates. Recent work has begun to address these shortcomings. To combat sparse rewards, Process Reward Models have been introduced to provide dense, line-level feedback, guiding the model even when the final code is incorrect (Dai et al., 2025). To address the flat reward landscape, researchers are exploring ways to incorporate a sense of difficulty into the learning process. The idea of combining RL with a curriculum is gaining traction. For instance, some approaches use RL to guide a model through a curriculum of subtasks, while others dynamically adjust the curriculum during RL training using techniques like the Self-Evolving Curriculum, which treats problem selection as a multi-armed bandit problem to maximize learning progress (Chen et al., 2025). Our TAROT framework specifically addresses the flat reward problem by making the reward signal itself curriculum-aware. Instead of treating all successes equally, we modulate the reward based on the difficulty of the solved test tier, a concept inspired by curriculum design. By integrating this tiered reward scheme directly into a stable policy optimization algorithm, the TAROT framework provides a more nuanced learning gradient that encourages the model to master harder problems. This approach of infusing a static curriculum structure directly into the RL reward mechanism is a novel contribution that complements other recent innovations in the field.

153 154 155

3 TAROT FRAMEWORK

160

161

In this section, we elaborate on the details of the proposed TAROT framework which enhances the code generation capability of language models by having them solve test cases of varying difficulty in appropriate order and rewarding weights that are adaptively determined based on the model's baseline capability.

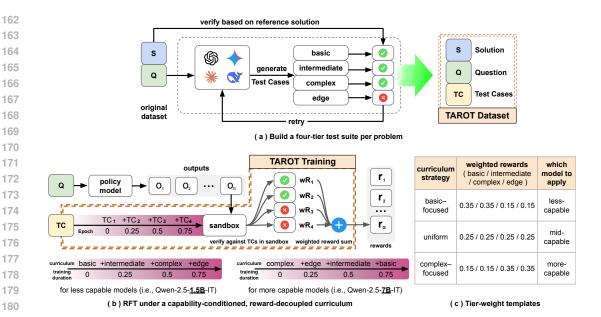


Figure 1: Overview of TAROT framework. (a) Build a four-tier test suite (basic/intermediate/complex/edge) per problem using frontier LLMs and verify them against the reference solution. (b) Reinforcement fine-tuning under a capability-conditioned, reward-decoupled curriculum. Less capable models perform best with basic \rightarrow complex, whereas more capable models perform best with complex \rightarrow basic. (c) Tier-weight templates specifying reward weights for basic, intermediate, complex, and edge, with suggested use by capability buckets.

3.1 TAROT DATASET

A coding problem, denoted as P, is formally defined as a tuple consisting of three core components: a problem statement (S), a reference solution (R), and a set of test suite (T). This relationship can be expressed as:

$$\mathcal{P} = (\mathcal{S}, \mathcal{R}, \mathcal{T}) \tag{1}$$

In this structure, the problem statement S outlines the task, and the reference solution R provides a correct implementation. The primary purpose of the test suite, \mathcal{T} , is to serve as a final validation mechanism to verify the correctness of a solver's proposed solution, but they are not designed to facilitate users' incremental learning processes. Consequently, the number and nature of test cases can vary significantly. For instance, a problem's test suite might consist of a single, simple case to verify the primary logic, or conversely, focus exclusively on complex edge cases, neither of which is structured to support a step-by-step learning process.

From a software engineering perspective, development is commonly test-driven. It begins with simple tests and progressively adds more complex and edge cases; implementations are refactored along the way, strengthening correctness and design. This staged expansion of the test suite mirrors the intuition behind curriculum learning. However, the test suite accompanying typical coding problems are not authored with this incremental pedagogy in mind. They are primarily designed for summative verification rather than stepwise scaffolding, so both their cardinality and difficulty mix vary widely and arbitrarily across problems.

To address the absence of a pedagogical structure, we introduce the TAROT dataset \mathcal{D}_{TAROT} , constructed by the procedure depicted in Figure 1 (a). Each problem in the dataset is augmented by incorporating a tiered test suite organized into four predefined difficulty levels (basic, intermediate, complex, and edge), without modifying the original statement or the reference solution. This structure is formally defined as follows:

221 222

224

225

226

227

228 229

230 231

232

233

234

235

236

237

238

239

240

241 242

243

244

245

246

247

248

249

250 251 252

253

254

255

256

257 258

259

260

261

262 263

$$\mathcal{D}_{\text{TAROT}} = \left\{ \left(\mathcal{S}_i, \ \mathcal{R}_i, \ \{ \mathcal{T}_{i,l} \}_{l \in L} \right) \right\}_{i=1}^N, \tag{2}$$

219 $L = \{ \text{basic, intermediate, complex, edge} \},$ 220

$$\mathcal{T}_{i} = \bigcup_{l \in L} \mathcal{T}_{i,l},$$
s.t. $\forall i \in [N], \ \forall l \in L, \ \forall t \in \mathcal{T}_{i,l} : \operatorname{Pass}(\mathcal{R}_{i}, t).$ (5)

s.t.
$$\forall i \in [N], \forall l \in L, \forall t \in \mathcal{T}_{i,l} : \operatorname{Pass}(\mathcal{R}_i, t).$$
 (5)

(3)

Here, L is the set of difficulty levels, and the full suite for each problem is the union of its per-level subsets equation 4. By construction (see equation 5), every test case is validated against the reference solution, ensuring data quality. Any curriculum order (e.g., basic \rightarrow complex \rightarrow basic) is imposed at training time and is not part of the dataset definition.

3.2 TAROT TRAINING MECHANISM

The TAROT training mechanism is designed to decouple the curriculum from raw test scores. It achieves this by utilizing two pre-defined components: a curriculum allocator that defines a fixed proportion of training focus for each difficulty tier $l \in L$, and tailored reward weights that prioritize tiers by placing greater value where the learning signal is most beneficial. During the training loop, the model first generates candidate solutions for a given problem. These solutions are then executed against the tiered test cases, and the resulting pass/fail outcomes are used to calculate and accumulate a tier-weighted return. Both the curriculum allocation and reward weights are aligned with the model's capability. The guiding principle is to concentrate the training signal within a zone of optimal difficulty, which is unique to each model's effective capability, a composite of instruction following fidelity and baseline coding proficiency. Therefore, the entire training schedule is preconfigured to match this profile, creating a fixed yet highly customized learning path.

Specifically, this means that models with lower baseline capability receive a larger share of basic and intermediate cases, whereas models with higher capability are assigned more complex and edge cases to push their frontier. The reward weights mirror this design, ensuring that successes on capability-appropriate tiers contribute more to the final objective.

We now formulate the training objective in the reinforcement learning setting. For each problem P_i and difficulty level $l \in L$, we define the tier-level success of a policy π as the average pass rate over that tier's tests.

$$r_{i,l}(\pi) = \frac{1}{|\mathcal{T}_{i,l}|} \sum_{t \in \mathcal{T}_{i,l}} \mathbf{1} \{ \operatorname{Pass}(\pi, t) \}.$$
 (6)

Here, $Pass(\pi, t)$ indicates that the solution produced under π satisfies test case t.

Given a curriculum allocation $\alpha = (\alpha_l)_{l \in L}$ and reward weights $w = (w_l)_{l \in L}$, we define the TAROT return for P_i as a weighted sum over tiers.

$$R_{\text{TAROT}}(P_i, \pi; \boldsymbol{\alpha}, \boldsymbol{w}) = \sum_{l \in L} \alpha_l \, w_l \, r_{i,l}(\pi), \qquad \sum_{l \in L} \alpha_l = 1, \quad w_l \ge 0. \tag{7}$$

We interpret α_l as the share of training effort assigned to tier l, and w_l as how much a success on tier l contributes given the model's baseline capability. Here, each α_l not only weights the contribution of tier l to the return but also specifies the fraction of training updates allocated to that tier.

Training then maximizes the expected TAROT return over problems.

$$J_{\text{TAROT}}(\theta) = \mathbb{E}_{P_i \sim \mathcal{D}_{\text{TAROT}}} \Big[R_{\text{TAROT}} \Big(P_i, \pi_{\theta}; \boldsymbol{\alpha}, \boldsymbol{w} \Big) \Big].$$
 (8)

268

269

This formulation provides a simple yet powerful objective function. By decoupling the allocation of training effort (α) from the valuation of success (w), the TAROT framework enables a fine-grained, capability-matched curriculum. This approach moves beyond imposing a single, model-agnostic learning path, instead concentrating the training signal on the most productive difficulty tiers for any given model.

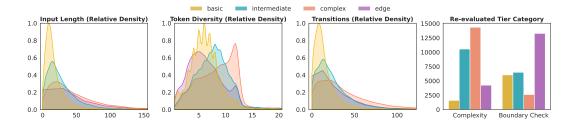


Figure 2: Quantitative and qualitative validation of the TAROT dataset's four-tiered structure. The KDE plots show a systematic increase in structural complexity metrics across tiers. On the right, a GPT-40 based re-validation qualitatively confirms that complex tiers test for complexity and edge tiers for boundary conditions.

4 EXPERIMENTS

4.1 EXPERIMENTAL SETTINGS

We construct a TAROT dataset based on 15k Python coding interview problems with validated basic/intermediate/complex/edge test suites. We prepare curricula couple allocation order (Forward/Reversed/Static with 0.2/0.4/0.6 staged transitions) and tier weights (Uniform, B&I, C&E). We evaluate TAROT training mechanism on a diverse suite of models, including Qwen2.5-Instruct, Qwen2.5-Coder-Instruct (1.5B, 3B, 7B) (Qwen et al., 2025; Hui et al., 2024), Gemma-2-IT (2B, 9B) (Team et al., 2024), and Qwen3-4B-Instruct-2507 (Yang et al., 2025) on well-known benchmarks including HumanEval (Chen et al., 2021), MBPP (Austin et al., 2021), HumanEval+, MBPP+ (Liu et al., 2024), LiveCodeBench v5 (Jain et al., 2024), CodeForces (Penedo et al., 2025), and CruxEval (Gu et al., 2024). This selection allows us to assess the framework's effectiveness across a wide spectrum of model scales, architectures, coding specializations, and performance tiers, including those at the frontier. Full implementation details and the curriculum schedules appear in Appendix C.

4.2 EXPERIMENTAL RESULTS

To validate the empirical integrity of the TAROT dataset's tiering, we analyzed its structure using quantitative and qualitative metrics, as illustrated in Figure 2. The three KDE plots demonstrate a clear progression: as the tiers advance from basic to complex, the distributions for input length, token diversity, and character transitions all exhibit a consistent rightward shift, signifying a systematic increase in structural complexity. Furthermore, the qualitative bar chart reveals a crucial distinction between the two hardest tiers. It shows that test cases designed to probe complexity peak in the complex tier, while those targeting boundary checks are overwhelmingly concentrated in the edge tier. These complementary findings confirm that our four-level taxonomy not only stratifies overall difficulty but also effectively separates different types of challenge, establishing a robust foundation for the subsequent experiments.

Experimental results, illustrated in Figure 3, reveal a nuanced relationship between model scale, specialization, and optimal curriculum design. The general-purpose Qwen2.5-Instruct models exhibit a straightforward, scale-dependent trend; the largest model (7B) performs best with complex-focused strategies, while the smallest (1.5B) benefits from a conventional basic-focused approach. However, this correlation with scale does not fully explain performance. The coding specialized Qwen2.5-Coder models introduce a critical insight, as the mid-scale Qwen2.5-Coder-3B model displays a learning preference akin to the much larger Instruct-7B model despite its smaller parameter count. It achieves its peak HumanEval score using the same complex-focused strategy and decisively outperforms its general-purpose 3B counterpart. This finding strongly suggests that a model's prior specialization enhances its effective capability, making it a more critical determinant of the ideal learning path than parameter count alone. The evidence supports the thesis that the optimal curriculum is dictated not by raw scale, but by a more holistic measure of a model's intrinsic abilities

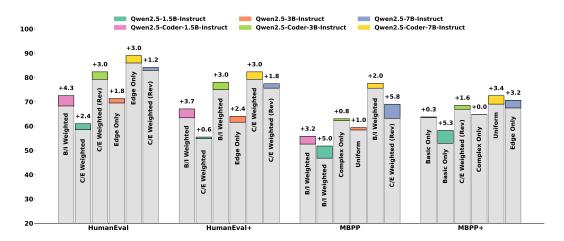


Figure 3: Experimental results for Qwen2.5-Instruct and Qwen2.5-Coder-Instruct (1.5B, 3B, 7B) on HumanEval, HumanEval+, MBPP, and MBPP+. Scores are pass@1. Numbers above bars indicate gains in percentage points (pp) relative to each model's base checkpoint without RFT. Labels inside bars indicate the best performing curriculum strategy.

shaped by its training history. Complete per-strategy results for Qwen2.5-Instruct and Qwen2.5-Coder-Instruct models are provided in Appendix G, Table 6.

To generalize these findings, the investigation was extended to the more recent Qwen3-4B-Instruct-2507. As detailed in Table 1, this newer model reinforces the core thesis. The optimal curriculum strategy, C/E Weighted, consistently outperforms the base model across all evaluated benchmarks, yielding substantial gains ranging from +2.12 to +4.26 percentage points. Notably, these improvements were achieved on a model that already possesses a strong performance baseline, confirming that curriculum learning is an effective method for eliciting further gains even from highly capable models. This result is significant because it demonstrates that a model's preference for advanced, complex-focused curricula is not merely a function of parameter count but is strongly tied to its overall capability. The success of this strategy on a powerful, state-of-the-art Qwen3-4B-Instruct-2507 further solidifies the argument that the most capable models benefit most from curricula that prioritize challenging examples. Full per-curriculum results for Qwen3-4B-Instruct-2507 are in Appendix G, Table 7. We also study sensitivity to training hyper-parameters (temperature, GRPO β) and to the inference max-token limit; see Appendix E and F for details.

Table 1: Performance comparison of curriculum strategies for Qwen3-4B-Instruct-2507, highlighting the top-performing *C/E Weighted* curriculum strategy against the base model. The performance delta is shown in percentage points (pp).

Strategy	HumanEval	HumanEval+	MBPP	MBPP+
Base	89.02%	78.66%	52.60%	56.61%
C/E Weighted	91.46 % (+2.44pp)	82.92 % (+4.26pp)	55.20 % (+2.60pp)	58.73 % (+2.12pp)

4.3 IN-DEPTH ANALYSIS

Performance on Out-of-Distribution Benchmarks. Evaluating the TAROT framework on out-of-distribution benchmarks like CodeForces, LiveCodeBench v5, and CruxEval revealed a crucial insight: the optimal curriculum strategy is not universal but is instead highly task-dependent. While curriculum learning consistently surpassed the baseline across all OOD (Out-of-Distribution) tasks (Figure 4), the most effective learning path varied significantly. For example, Qwen2.5-7B model achieved its best performance on LiveCodeBench v5 with the *Basic Only* curriculum (+5.6 pp), whereas the *C/E Weighted* strategy proved most effective for the same model on CruxEval (+5.0 pp)

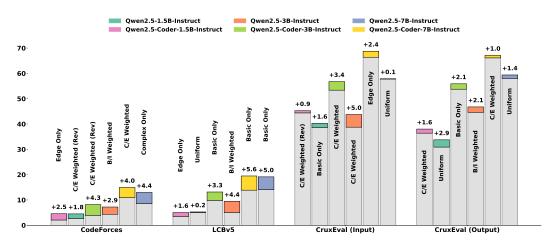


Figure 4: Experimental results for Qwen2.5-Instruct and Qwen2.5-Coder-Instruct models (1.5B, 3B, 7B) on CodeForces, LiveCodeBench v5(LCBv5), and CruxEval. Scores are the overall accuracy across easy, medium, and hard problems. Numbers above bars indicate gains in percentage points (pp) relative to each model's base checkpoint without reinforcement fine-tuning. Labels inside bars indicate the best performing curriculum strategy.

and CodeForces (+4.0 pp). This divergence arises because these benchmarks test different skills than our coding interview style training data. This finding underscores the limitations of a one-size-fits-all curriculum and suggests that different learning strategies may be required across tasks. It further highlights the value of task specific intra-problem test suites that reflect the computational structure of the target domain, calling for future research on curriculum policy selection and intra-problem test design.

Applying TAROT to Other Model Architectures. Applying TAROT to Gemma2 architecture validated our hypothesis that baseline proficiency, not parameter count, dictates the optimal learning path. For the larger Gemma2-9B-IT, *Complex Only* curriculum offered no decisive advantage, as simpler strategies like *Basic Only* often proved superior on key benchmarks shown in Table 2. This principle was even more starkly visible with the smaller Gemma2-2B-IT. As detailed in Appendx H, most curricula were actively harmful, leading to a performance collapse from a sparse reward signal. In stark contrast, a *Basic Only* strategy focused on fundamentals yielded substantial improvements. This demonstrates that for less-capable models, a fundamentals-first curriculum is a prerequisite for successful fine-tuning, whereas unstructured approaches can be severely detrimental.

Table 2: Performance comparison for Gemma2-9B-IT across key curriculum strategies. Highest scores on each benchmark are highlighted in **bold**.

Strategy	HumanEval	HumanEval+	MBPP	MBPP+	CodeForces	LCBv5	CruxEval
Base	60.37%	54.88%	59.60%	65.08%	8.61%	11.83%	45.63/47.63%
Uniform	65.85%	57.93%	59.20%	64.55%	10.82%	13.62%	51.63 /47.13%
Basic Only	63.41%	56.10%	60.40%	65.08%	9.49%	14.70%	51.00/48.00%
Complex Only	65.85%	60.37%	58.60%	64.55%	9.93%	12.54%	48.25/48.00%

The Limits of the Reward Signal. Figure 5 (a) shows that the training reward increases steadily and is clearly separated by model capacity, indicating a stable optimization process. This pattern suggests that the policy learns the training distribution well and that stronger models achieve higher reward levels under the same curriculum. However, the reward observed during training does not reliably anticipate downstream benchmark outcomes. As shown in Figure 5 (c), the final reward has only a weak Pearson correlation coefficient with benchmark scores, which means that runs with similar rewards can still deliver very different levels of task performance.

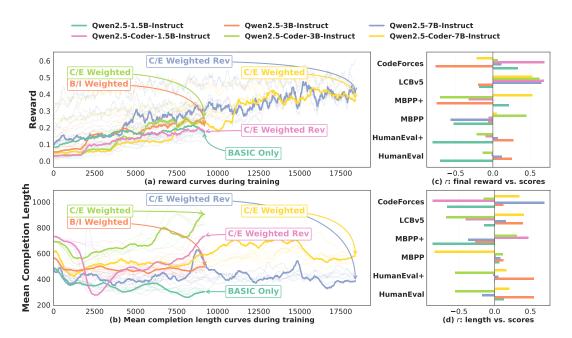


Figure 5: Training dynamics vs. downstream performance. (a) and (b) show the reward and the mean completion length curves during reinforcement fine-tuning, and the annotations mark the curriculum strategy with the best average downstream performance. (c) and (d) show the Pearson correlation coefficient r of the final rewards vs. benchmark scores and the mean completion length vs. benchmark scores, respectively.

Conciseness as a Proxy for Advanced Reasoning. A different perspective comes from analyzing completion length. Figure 5 (b) shows that models with greater capability tend to produce shorter solutions as training progresses, and this tendency becomes more pronounced for stronger configurations. Importantly, Figure 5 (d) indicates that mean completion length exhibits a stronger negative correlation with benchmark scores than the reward does, implying that conciseness aligns better with final solution quality. Shorter programs are more likely to capture the essential reasoning steps without unnecessary detours, whereas longer outputs often reflect uncertainty or inefficient search. These observations support using solution conciseness as a practical secondary proxy for advanced reasoning quality, complementing the reward based perspective and providing a more informative early indicator of downstream performance.

5 CONCLUSION

We introduced TAROT, a test-driven and capability-adaptive framework for curriculum reinforcement fine-tuning in code generation. TAROT challenges the prevailing one-size-fits-all approach by constructing a four-tier, intra-problem test suite that allows curriculum design to be tailored to a model's unique abilities. Experiments confirmed our central thesis that the optimal learning path is capability-dependent: less capable models benefit most from an basic-focused progression, while more competent models excel with curricula that prioritize complex-focused challenges. We found that the most critical factor is not parameter count alone but a more holistic effective capability, which accounts for a model's prior specialization. Ultimately, TAROT provides a practical and effective recipe for enhancing the code generation capabilities of Large Language Models. Our framework proved its value across a wide spectrum of models, from less-capable base models to highly proficient code-specialized and state-of-the-art foundation models, confirming its broad applicability. While significant, our results also show that the optimal curriculum is task-dependent, pointing toward future work in developing domain-specific test suites and automated policy selection methods. Continuing to refine this capability-aware approach is a crucial step toward developing more powerful and reliable code generation models.

REFERENCES

- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, and Charles Sutton. Program synthesis with large language models, 2021. URL https://arxiv.org/abs/2108.07732.
- K. Beck. Test-driven Development: By Example. Addison-Wesley signature series. Addison-Wesley, 2003. ISBN 9780321146533. URL https://books.google.co.kr/books?id=CUlsAQAAQBAJ.
- Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. Curriculum learning. In *Proceedings of the 26th Annual International Conference on Machine Learning ICML '09*, pp. 1–8, Montreal, Quebec, Canada, 2009. ACM Press. ISBN 978-1-60558-516-1. doi: 10.1145/1553374.1553380. URL http://portal.acm.org/citation.cfm?doid=1553374.1553380.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large language models trained on code. 2021.
- Xiaoyin Chen, Jiarui Lu, Minsu Kim, Dinghuai Zhang, Jian Tang, Alexandre Piché, Nicolas Gontier, Yoshua Bengio, and Ehsan Kamalloo. Self-evolving curriculum for llm reasoning, 2025. URL https://arxiv.org/abs/2505.14970.
- Yang Cheng, Zilai Wang, Weiyu Ma, Wenhui Zhu, Yue Deng, and Jian Zhao. Evocurr: Self-evolving curriculum with behavior code generation for complex decision-making, 2025. URL https://arxiv.org/abs/2508.09586.
- Ning Dai, Zheng Wu, Renjie Zheng, Ziyun Wei, Wenlei Shi, Xing Jin, Guanlin Liu, Chen Dun, Liang Huang, and Lin Yan. Process supervision-guided policy optimization for code generation, 2025. URL https://arxiv.org/abs/2410.17621.
- Shihan Dou, Yan Liu, Haoxiang Jia, Enyu Zhou, Limao Xiong, Junjie Shan, Caishuang Huang, Xiao Wang, Xiaoran Fan, Zhiheng Xi, Yuhao Zhou, Tao Ji, Rui Zheng, Qi Zhang, Tao Gui, and Xuanjing Huang. StepCoder: Improving code generation with reinforcement learning from compiler feedback. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 4571–4585, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.251. URL https://aclanthology.org/2024.acl-long.251/.
- Xueying Du, Mingwei Liu, Kaixin Wang, Hanlin Wang, Junwei Liu, Yixuan Chen, Jiayi Feng, Chaofeng Sha, Xin Peng, and Yiling Lou. Evaluating large language models in class-level code generation. In *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering*, ICSE '24, New York, NY, USA, 2024. Association for Computing Machinery. ISBN 9798400702174. doi: 10.1145/3597503.3639219. URL https://doi.org/10.1145/3597503.3639219.
- Alex Gu, Baptiste Rozière, Hugh Leather, Armando Solar-Lezama, Gabriel Synnaeve, and Sida I. Wang. Cruxeval: A benchmark for code reasoning, understanding and execution. *arXiv preprint arXiv:2401.03065*, 2024.
 - Hugging Face. Open r1: A fully open reproduction of deepseek-r1, January 2025. URL https://github.com/huggingface/open-r1.

- Binyuan Hui, Jian Yang, Zeyu Cui, Jiaxi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Keming Lu, Kai Dang, Yang Fan, Yichang Zhang, An Yang, Rui Men, Fei Huang, Bo Zheng, Yibo Miao, Shanghaoran Quan, Yunlong Feng, Xingzhang Ren, Xuancheng Ren, Jingren Zhou, and Junyang Lin. Qwen2.5-coder technical report, 2024. URL https://arxiv.org/abs/2409.12186.
- Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free evaluation of large language models for code, 2024. URL https://arxiv.org/abs/2403.07974.
- Juyong Jiang, Fan Wang, Jiasi Shen, Sungju Kim, and Sunghun Kim. A survey on large language models for code generation, 2024. URL https://arxiv.org/abs/2406.00515.
- Kyi Shin Khant, Hong Yi Lin, and Patanamon Thongtanunam. Should code models learn pedagogically? a preliminary evaluation of curriculum learning for real-world software engineering tasks. In 2025 IEEE/ACM 22nd International Conference on Mining Software Repositories (MSR), pp. 249–254, 2025. doi: 10.1109/MSR66628.2025.00044.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*, 2023.
- Jiawei Liu, Songrun Xie, Junhao Wang, Yuxiang Wei, Yifeng Ding, and Lingming Zhang. Evaluating language models for efficient code generation. In *First Conference on Language Modeling*, 2024. URL https://openreview.net/forum?id=IBCBMeAhmC.
- Marwa Naïr, Kamel Yamani, Lynda Lhadj, and Riyadh Baghdadi. Curriculum learning for small code language models. In Xiyan Fu and Eve Fleisig (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 4: Student Research Workshop)*, pp. 390–401, Bangkok, Thailand, August 2024. Association for Computational Linguistics. ISBN 979-8-89176-097-4. doi: 10.18653/v1/2024.acl-srw.44. URL https://aclanthology.org/2024.acl-srw.44/.
- 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 Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, NIPS '22, Red Hook, NY, USA, 2022. Curran Associates Inc. ISBN 9781713871088.
- Guilherme Penedo, Anton Lozhkov, Hynek Kydlíček, Loubna Ben Allal, Edward Beeching, Agustín Piqueres Lajarín, Quentin Gallouédec, Nathan Habib, Lewis Tunstall, and Leandro von Werra. Codeforces. https://huggingface.co/datasets/open-r1/codeforces, 2025.
- Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report, 2025. URL https://arxiv.org/abs/2412.15115.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. Direct preference optimization: your language model is secretly a reward model. In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, NIPS '23, Red Hook, NY, USA, 2023. Curran Associates Inc.
- Negin Raoof, Etash Kumar Guha, Ryan Marten, Jean Mercat, Eric Frankel, Sedrick Keh, Hritik Bansal, Georgios Smyrnis, Marianna Nezhurina, Trung Vu, Zayne Rea Sprague, Mike A

595

596

597

598

600 601

602

603 604

605

607

608 609

610

612

613

614

615

616

617

618

619

620

621

622

623

625

626

627

630

631

632

633

634

635

636

637

638

639

640

641 642

643

644

645

646

647

Merrill, Liangyu Chen, Caroline Choi, Zaid Khan, Sachin Grover, Benjamin Feuer, Ashima Suvarna, Shiye Su, Wanjia Zhao, Kartik Sharma, Charlie Cheng-Jie Ji, Kushal Arora, Jeffrey Li, Aaron Gokaslan, Sarah M Pratt, Niklas Muennighoff, Jon Saad-Falcon, John Yang, Asad Aali, Shreyas Pimpalgaonkar, Alon Albalak, Achal Dave, Hadi Pouransari, Greg Durrett, Sewoong Oh, Tatsunori Hashimoto, Vaishaal Shankar, Yejin Choi, Mohit Bansal, Chinmay Hegde, Reinhard Heckel, Jenia Jitsev, Maheswaran Sathiamoorthy, Alex Dimakis, and Ludwig Schmidt. Evalchemy, 2025.

John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language models, 2024. URL https://arxiv.org/abs/2402.03300.

Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, Johan Ferret, Peter Liu, Pouya Tafti, Abe Friesen, Michelle Casbon, Sabela Ramos, Ravin Kumar, Charline Le Lan, Sammy Jerome, Anton Tsitsulin, Nino Vieillard, Piotr Stanczyk, Sertan Girgin, Nikola Momchev, Matt Hoffman, Shantanu Thakoor, Jean-Bastien Grill, Behnam Neyshabur, Olivier Bachem, Alanna Walton, Aliaksei Severyn, Alicia Parrish, Aliya Ahmad, Allen Hutchison, Alvin Abdagic, Amanda Carl, Amy Shen, Andy Brock, Andy Coenen, Anthony Laforge, Antonia Paterson, Ben Bastian, Bilal Piot, Bo Wu, Brandon Royal, Charlie Chen, Chintu Kumar, Chris Perry, Chris Welty, Christopher A. Choquette-Choo, Danila Sinopalnikov, David Weinberger, Dimple Vijaykumar, Dominika Rogozińska, Dustin Herbison, Elisa Bandy, Emma Wang, Eric Noland, Erica Moreira, Evan Senter, Evgenii Eltyshev, Francesco Visin, Gabriel Rasskin, Gary Wei, Glenn Cameron, Gus Martins, Hadi Hashemi, Hanna Klimczak-Plucińska, Harleen Batra, Harsh Dhand, Ivan Nardini, Jacinda Mein, Jack Zhou, James Svensson, Jeff Stanway, Jetha Chan, Jin Peng Zhou, Joana Carrasqueira, Joana Iljazi, Jocelyn Becker, Joe Fernandez, Joost van Amersfoort, Josh Gordon, Josh Lipschultz, Josh Newlan, Ju yeong Ji, Kareem Mohamed, Kartikeya Badola, Kat Black, Katie Millican, Keelin McDonell, Kelvin Nguyen, Kiranbir Sodhia, Kish Greene, Lars Lowe Sjoesund, Lauren Usui, Laurent Sifre, Lena Heuermann, Leticia Lago, Lilly McNealus, Livio Baldini Soares, Logan Kilpatrick, Lucas Dixon, Luciano Martins, Machel Reid, Manvinder Singh, Mark Iverson, Martin Görner, Mat Velloso, Mateo Wirth, Matt Davidow, Matt Miller, Matthew Rahtz, Matthew Watson, Meg Risdal, Mehran Kazemi, Michael Moynihan, Ming Zhang, Minsuk Kahng, Minwoo Park, Mofi Rahman, Mohit Khatwani, Natalie Dao, Nenshad Bardoliwalla, Nesh Devanathan, Neta Dumai, Nilay Chauhan, Oscar Wahltinez, Pankil Botarda, Parker Barnes, Paul Barham, Paul Michel, Pengchong Jin, Petko Georgiev, Phil Culliton, Pradeep Kuppala, Ramona Comanescu, Ramona Merhej, Reena Jana, Reza Ardeshir Rokni, Rishabh Agarwal, Ryan Mullins, Samaneh Saadat, Sara Mc Carthy, Sarah Cogan, Sarah Perrin, Sébastien M. R. Arnold, Sebastian Krause, Shengyang Dai, Shruti Garg, Shruti Sheth, Sue Ronstrom, Susan Chan, Timothy Jordan, Ting Yu, Tom Eccles, Tom Hennigan, Tomas Kocisky, Tulsee Doshi, Vihan Jain, Vikas Yadav, Vilobh Meshram, Vishal Dharmadhikari, Warren Barkley, Wei Wei, Wenming Ye, Woohyun Han, Woosuk Kwon, Xiang Xu, Zhe Shen, Zhitao Gong, Zichuan Wei, Victor Cotruta, Phoebe Kirk, Anand Rao, Minh Giang, Ludovic Peran, Tris Warkentin, Eli Collins, Joelle Barral, Zoubin Ghahramani, Raia Hadsell, D. Sculley, Jeanine Banks, Anca Dragan, Slav Petrov, Oriol Vinyals, Jeff Dean, Demis Hassabis, Koray Kavukcuoglu, Clement Farabet, Elena Buchatskaya, Sebastian Borgeaud, Noah Fiedel, Armand Joulin, Kathleen Kenealy, Robert Dadashi, and Alek Andreev. Gemma 2: Improving open language models at a practical size, 2024. URL https://arxiv.org/abs/2408.00118.

Leandro von Werra, Younes Belkada, Lewis Tunstall, Edward Beeching, Tristan Thrush, Nathan Lambert, Shengyi Huang, Kashif Rasul, and Quentin Gallouédec. Trl: Transformer reinforcement learning. https://github.com/huggingface/trl, 2020.

Thomas Weber, Maximilian Brandmaier, Albrecht Schmidt, and Sven Mayer. Significant productivity gains through programming with large language models. *Proc. ACM Hum.-Comput. Interact.*, 8(EICS), June 2024. doi: 10.1145/3661145. URL https://doi.org/10.1145/3661145.

 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pp. 38–45, Online, October 2020. Association for Computational Linguistics. URL https://www.aclweb.org/anthology/2020.emnlp-demos.6.

- Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, Qingwei Lin, and Daxin Jiang. WizardLM: Empowering large pre-trained language models to follow complex instructions. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=CfXh93NDgH.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. Qwen3 technical report, 2025. URL https://arxiv.org/abs/2505.09388.
- Chujie Zheng, Shixuan Liu, Mingze Li, Xiong-Hui Chen, Bowen Yu, Chang Gao, Kai Dang, Yuqiong Liu, Rui Men, An Yang, Jingren Zhou, and Junyang Lin. Group sequence policy optimization, 2025. URL https://arxiv.org/abs/2507.18071.
- Chunting Zhou, Pengfei Liu, Puxin Xu, Srini Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, Susan Zhang, Gargi Ghosh, Mike Lewis, Luke Zettlemoyer, and Omer Levy. Lima: less is more for alignment. In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, NIPS '23, Red Hook, NY, USA, 2023. Curran Associates Inc.
- Terry Yue Zhuo, Minh Chien Vu, Jenny Chim, Han Hu, Wenhao Yu, Ratnadira Widyasari, Imam Nur Bani Yusuf, Haolan Zhan, Junda He, Indraneil Paul, Simon Brunner, Chen Gong, Thong Hoang, Armel Randy Zebaze, Xiaoheng Hong, Wen-Ding Li, Jean Kaddour, Ming Xu, Zhihan Zhang, Prateek Yadav, Naman Jain, Alex Gu, Zhoujun Cheng, Jiawei Liu, Qian Liu, Zijian Wang, Binyuan Hui, Niklas Muennighoff, David Lo, Daniel Fried, Xiaoning Du, Harm de Vries, and Leandro Von Werra. Bigcodebench: Benchmarking code generation with diverse function calls and complex instructions, 2025. URL https://arxiv.org/abs/2406.15877.

A SAMPLE TIERED TEST CASES

Table 3 presents concrete examples of the four-tiered test cases generated for several problems in the TAROT dataset. These samples illustrate a clear and intentional progression in difficulty and scope, which is a cornerstone of our framework.

The tiers are generally designed to validate different aspects of a solution. Basic tiers focus on the core logic of a problem with simple, straightforward inputs. Following this, intermediate and complex tiers introduce greater difficulty through larger inputs, more intricate scenarios, or patterns requiring more sophisticated algorithmic reasoning. Finally, edge tiers are designed to test for robustness by probing boundary conditions, constraints, and performance-intensive cases such as large numbers or long strings. This tiered structure exemplifies the intra-problem difficulty gradient that forms the basis of our capability-adaptive curriculum.

B TEST CASE GENERATION PROMPTS

To ensure the consistent generation of high-quality, four-tiered test cases, we designed a detailed prompt template. This template, listed in Table 4, guides the language model to act as an expert software engineer and produce test cases that adhere to our specific difficulty criteria.

C IMPLEMENTATION DETAILS

TAROT Dataset. Our experiments utilize the TAROT dataset, which we constructed by augmenting approximately 15,000 problems from the verifiable-coding-problems-python dataset 1 . For each problem, we employed OpenAI's the most powerful o3 and o4 models 2 with the highest reasoning effort to generate a four-tiered test suite with distinct levels: basic, intermediate, complex, and edge. The specific prompts used for this generation process are detailed in Appendix B. To ensure high quality, every generated test case was validated against the reference solution, and any problem with even one failing tier was discarded. This rigorous curation process yields a final dataset of approximately 60,000 tiered test suites (15,000 problems \times 4 tiers). Samples of these generated tiered test cases can be found in Appendix A.

Model Selection. To validate our methodology, we selected a diverse set of models to investigate four key research questions: (1) the effect of model scale, to test our hypothesis that the optimal curriculum is capability-dependent, using three Qwen2.5 models of varying sizes (1.5B, 3B, 7B) (Qwen et al., 2025); (2) the impact of specialization, to determine if TAROT can further enhance models already proficient in coding, using their code-specialized counterparts (Hui et al., 2024); (3) architectural generalizability, to test if our findings apply beyond a single model family, by incorporating two instruction-tuned Gemma2 models (2B, 9B) (Team et al., 2024); and (4) pushing performance frontiers, to assess if our framework can improve even state-of-the-art models with strong baselines, by fine-tuning the recent Qwen3-4B-Instruct-2507 (Yang et al., 2025). For all models, we used their instruction-tuned variants to ensure a foundational code-generation capability, a prerequisite for effective RL-based fine-tuning.

Curriculum Settings. Our curriculum experiments are designed around two independent variables detailed in Table 5: the allocation strategy and the reward scheme. We test three primary allocation strategies: Forward, which progressively introduces tiers from basic to complex; Reversed, which inverts this sequence; and Static, which dedicates all training effort to a single tier. These strategies are combined with three distinct reward weighting schemes: Uniform, which assigns equal value to all tiers; Basic & Intermediate Weighted, which emphasizes easier tiers; and Complex & Edge Weighted, which prioritizes the most difficult ones. This diverse set of configurations allows us to empirically test our central hypothesis: that the optimal learning path is not universal, but is instead contingent on a model's effective capability.

https://huggingface.co/datasets/open-r1/verifiable-coding-problems-python
https://platform.openai.com/docs/models

Table 3: Samples of Python coding problems with synthetically generated four tiers of test cases.

Problem	Test Case Difficulty	Input	Expected Output
Color de Cillerin e d'acceptance de la constitución de	basic	2 1 zebra 3 apple	azebr aelpp
Solve the following coding problem using the programming lan- guage python: You are given a string S and an integer L. A op- eration is described as: "You are given a string S, and you have to perform the given L operations. In each operation, you will be	intermediate	3 1 dbca 4 zzxyyxzx 5 hellohello	adbc xxxyyzzz eehhlllloo
given an integer 'p' and a character 'c'. You have to replace the p-th character of the string with c. After performing all the operations, you have to sort the string in ascending order and print it."	complex	2 3 mississippi 1 baab	iiiimppssss aabb
	basic 2 1 zebra 3 apple az lan-lan-land op-have iintermediate 3 1 dbca 4 zzxyyxzx 5 hellohello ad iill be e the peratit." edge 2 1 qwertyuiopasdfghjklz 20 qwe rtyuiopasdfghjklz 20 qwe rtyuiopasdf	asdfghjklzqwertyuiop adefghijk opqrstuwyz	
	basic	3 BA BABA ABAB	yes yes yes
Solve the following coding problem using the programming lan- guage python: Two sisters, A and B, play the piano every day. During the day, they can play in any order. For a given sequence of A and B, determine if it is possible to divide it into several parts so	intermediate		no yes no
that each part is a concatenation of two identical strings.	complex	AB BAABBAABBAABBAAB BAAA ABBAABBAABBAAB	yes no yes
	edge	AAA BBBBBBBBBBBBBBB	no no
Solve the following coding problem using the programming lan- guage python: Valya and Tolya are an ideal pair, but they quarrel sometimes. Recently, Valya took offense at Tolya and left home. Now Tolya is very sad and wants to reconcile with Valya. For this,	basic	5 abcde bcdea	4 a b b c c d d e
he is going to make her a gift — a necklace. A necklace is a sequence of beads on a string. Valya will be happy if she can read her name on the necklace. Her name is a string S. The necklace is a	intermediate	6 abcdef bcdefa	5 a b b c c d d e e f
string T. Valya can read her name if S is a subsequence of T. Tolya has two strings of beads, A and B. He wants to create a necklace T by choosing some subsequence of A and some subsequence of B	complex	10 abcdefghij bcdeaghijf	8 a b b c c d d e f g g h h i i j
and concatenating them. Moreover, he wants to use as many beads as possible, i.e. the length of T must be maximum. Find the maximum possible length of T such that S is a subsequence of T.	intermediate 3 1 dbca 4 zzxyyxzx 5 hellohello complex 2 3 mississippi 1 baab edge 21 qwertyuiopasdfghjklz 20 qwe rtyuiopasdfghjklz basic 3 BA BABA ABAB intermediate 3 AABB BAABABA AB- BAAAAB complex 3 BAABBAABBAABBAABBAABBAABBAAABBAAABBA	25 abbccddeeffgghhii jkkllmmnnooppqqrrs ttuuvvwwxxyyz	
	basic	7 23	23
Solve the following coding problem using the programming lan- guage python: In a recent [breakthrough] in mathematics, the proof utilized a concept called Height Balanced Tree. A Height Balanced Tree is a binary tree in which the height of the left and right sub-	intermediate	50 99	99
trees of any node differ by not more than 1. You are given an integer N. You have to find the maximum number of nodes a Height Balanced Tree of height N can have.	complex	97 89	97
	basic 2 1 zebra 3 apple intermediate 3 1 dbca 4 zzxyyxzx 5 helloho complex 2 3 mississippi 1 baab edge 2 1 qwertyuiopasdfghjklz 20 or rtyuiopasdfghjklz basic 3 BA BABA ABAB intermediate 3 AABB BAABAABAABAABAABBAABBAABBAABBAAB	11	1
Solve the following coding problem using the programming lan-	basic	3 1 2 3	4
guage python: Alex doesn't like boredom. That's why whenever he gets bored, he comes up with games. One day he came up with a game with numbers. He has a sequence of n numbers a1, a2,,	intermediate	6122334	8
an. He is allowed to perform the following operation any number of times: choose an index i $(1 = i = n)$ and replace the number ai with ai $+ 1$. He wants to make the sequence strictly increasing. What is the minimum number of operations he needs to perform?	complex	13 2 3 2 6 7 2 8 7 3 6 4 7 10	41
	•		200002
This problem is actually a subproblem of problem G from the same contest. There are n bone piles on the ground. You have to choose m piles and paint them black. Then, you will need to arrange the n	basic	1512345	1
m piles and paint them black. Then, you will need to arrange the n piles as a sequence with some order (you may place the piles at the same location). Each time you choose some piles with the same maximum height (larger than 0), take away the top block of each	intermediate	1 10 1 2 1 2 1 2 1 2 1 2	9
maximum neight (larger man 0), take away the top block of each of these piles. If all of the chosen piles are black, this operation is called a "good move"; otherwise, it is a "bad move". The score of this sequence is defined as the number of "good moves". You want to maximize the score over all m-blackened and all possible	complex		25
sequences. You need to output the maximum score.	edge	3114111141234	1 4 1

Table 4: The Prompt template used to generate a tiered test suite per a given coding problem. The problem statement and the default test case from the original source is injected into {problem_statement} and {baseline_test_case} placeholders respectively.

You are an expert software engineer with extensive experience in designing comprehensive unit tests. Your task is to generate four distinct unit tests for a given code implementation based solely on the provided problem statement. Treat this as a black-box testing exercise—focus exclusively on the inputs and expected outputs without assuming any details about the internal implementation.

Important: A baseline test case will be provided separately. Each test case you generate must be more challenging than the baseline test case.

Please generate four unit tests with the following guidelines:

1. Basic Complexity Test (label as "basic"):

- · Use simple, straightforward inputs.
- · Validate the core behavior under normal conditions.
- · Focus on the happy path scenario.
- · This should be the least challenging test case relative to the others.

2. Medium Complexity Test (label as "intermediate"):

- · Include moderately complex inputs with some edge conditions.
- · Test with mixed data types or larger inputs.
- · Incorporate common edge cases and boundary values.
- · Ensure this test is more challenging than the basic test.

3. High Complexity Test (label as "complex"):

- Use complex, nested, or structured inputs.
- · Validate advanced functionality and complex logic paths.
- · Stress test the implementation with challenging scenarios.
- · This test should be more intricate than both the basic and intermediate tests.

4. Edge Case Test (label as "edge"):

- · Use extreme boundary conditions and special cases.
- · Validate behavior with empty, null, or invalid inputs.
- Focus on error handling and exception scenarios.
- · This should be the most challenging test case among the four.

For each test case, follow the JSON format provided in the example below (include only the input and expected output):

Remember:

- Do not assume any knowledge about the internal code; base your tests purely on the input-output behavior described in the problem statement.
- Ensure that each of your test cases is incrementally more challenging than the baseline test case provided.

Problem Statement: {problem_statement}
Baseline Test Case: {baseline_test_case}

Training Details. We fine-tune all selected models for a single epoch using the TAROT framework. For policy optimization, we employ GRPO (Shao et al., 2024). All models are trained using the AdamW optimizer with a constant learning rate of 1×10^{-6} . We set the global batch size to 8, reducing it to 4 for larger models (Qwen2.5-7B-Instruct, Qwen2.5-Coder-7B-Instruct, and Gemma2-9B-IT) to accommodate memory constraints. The maximum input and completion token lengths

Table 5: Overview of the experimental schedules for curriculum learning. Each strategy varies in reward distribution and the sequence of difficulties presented to the model. The abbreviations B, I, C, and E correspond to basic, intermediate, complex, and edge difficulty tiers, respectively. For staged curricula, transitions occur at 0.2, 0.4, and 0.6 of the total epoch.

Strategy	Reward Weights (B, I, C, E)	Curriculum Schedule Progression
Forward (Uniform)	(0.25, 0.25, 0.25, 0.25)	$\overline{ B \to (B,I) \to (B,I,C) \to All }$
Forward (B & I Weighted)	(0.35, 0.35, 0.15, 0.15)	$B \to (B,I) \to (B,I,C) \to All$
Forward (C & E Weighted)	(0.15, 0.15, 0.35, 0.35)	$B \to (B,I) \to (B,I,C) \to All$
Reversed (C & E Weighted)	(0.15, 0.15, 0.35, 0.35)	$C \to (C,\!E) \to (C,\!E,\!I) \to All$
Basic Only	(1.0,,,)	Static
Complex Only	(,, 1.0,)	Static
Edge Only	(,,, 1.0)	Static

were set to 1,024 and 4,096, respectively. For GRPO-specific settings, we generated 8 candidate completions per prompt to estimate the policy advantage, with the core hyperparameter β set to 0.01 in our main experiments. We provide an ablation study on key training hyperparameters, including the GRPO β value (0.1, 0.05, 0.01) and the sampling temperature during training (1.0, 0.7, 0.5), in Appendix E.

The GRPO hyperparameter β controls the strength of the Kullback-Leibler (KL) divergence regularization term, which penalizes the policy for deviating too far from the original base model's behavior. The training temperature, in turn, manages the exploration-exploitation trade-off; higher values encourage the model to sample a wider variety of solutions (exploration), while lower values cause it to refine high-probability ones (exploitation). Our ablation studies were designed to identify the optimal settings for these crucial parameters within our code generation task.

All fine-tuning experiments were conducted on a server with 8 x NVIDIA A100 (80 GB) GPUs, running CUDA 12.4 and PyTorch 2.6. Our implementation is based on open-source libraries including Transformers (Wolf et al., 2020), TRL (von Werra et al., 2020), vLLM (Kwon et al., 2023), and Open-R1 (Hugging Face, 2025).

Evaluation Metrics. We evaluate the efficacy of the TAROT framework on a comprehensive suite of code generation benchmarks. For functional correctness, we measure the pass@1 metric on HumanEval (Chen et al., 2021), MBPP (Austin et al., 2021), and their more challenging variants, HumanEval+ and MBPP+ (Liu et al., 2024). To assess competitive problem-solving skills, we use the overall accuracy on LiveCodeBench v5 (Jain et al., 2024) and CodeForces (Penedo et al., 2025), averaged across their difficulty tiers. Finally, the model's code reasoning capability is evaluated using the input and output prediction accuracy on CruxEval (Gu et al., 2024). The detailed generation parameters and execution environment are described in Appendix D.

D GENERATION AND EXECUTION ENVIRONMENT

The entire evaluation pipeline is managed by the EvalChemy framework (Raoof et al., 2025). We follow the benchmark-specific generation configurations predefined within the framework—such as temperature, top-p, prompt formatting, and stopping criteria—to ensure consistency with established evaluation protocols. By default, the maximum completion tokens for each benchmark adhered to its standard setting; however, for an ablation study on generation length (Appendix F), we systematically increased this limit to 4,096, 8,192, and 16,384 tokens to observe performance trends.

All code generation for evaluation was conducted by serving the fine-tuned models via the vLLM framework (Kwon et al., 2023) on servers equipped with 4 x NVIDIA A100 (80 GB) GPUs, using a batch size of 64. The resulting code is executed in a secure, sandboxed Python 3.11 environment,

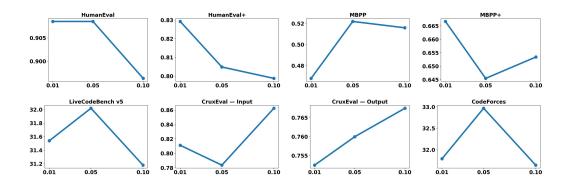


Figure 6: Performance sensitivity to the GRPO hyperparameter β . The plots show the final pass@1 or accuracy scores on various benchmarks as β is varied. The optimal value is task-dependent; for instance, HumanEval and HumanEval+ benefit from a smaller β (0.01) that allows greater policy exploration, whereas MBPP and CodeForces achieve peak performance with a larger β (0.05) that enforces stronger regularization.

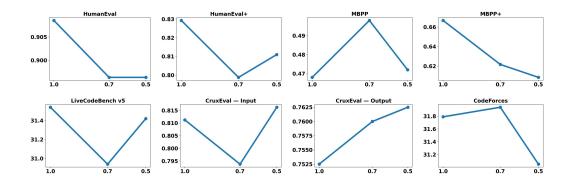


Figure 7: Performance sensitivity to the sampling temperature during training. The plots illustrate the final benchmark scores for different training temperatures. A higher temperature of 1.0, which encourages greater exploration, is optimal for benchmarks like HumanEval and HumanEval+. In contrast, other benchmarks such as MBPP show a preference for a more moderate temperature of 0.7, highlighting that the ideal exploration-exploitation balance is task-specific.

where a strict 10-second timeout is enforced for each test case to prevent infinite loops and manage evaluation time.

E HYPERPARAMETER SENSITIVITY ANALYSIS

This section provides ablation studies on two key training hyperparameters to analyze their impact on final benchmark performance: the GRPO regularization coefficient β and the sampling temperature during training.

Impact of GRPO's β The hyperparameter β in GRPO controls the strength of the Kullback-Leibler (KL) divergence regularization, which prevents the fine-tuned policy from deviating excessively from the original base model. The results of varying β are shown in Figure 6. Performance sensitivity to β is not uniform across benchmarks. For function-synthesis tasks like HumanEval and HumanEval+, a small β of 0.01, which allows for greater policy exploration, yields the best results. Conversely, benchmarks like MBPP and CodeForces appear to benefit from slightly stronger regularization ($\beta = 0.05$). This variance suggests that the optimal regularization strength is task-

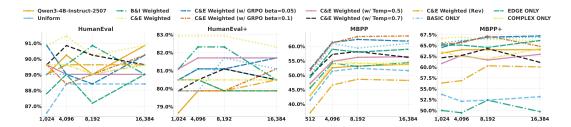


Figure 8: Performance sensitivity to the maximum completion token limit at inference time for Qwen3-4B-Instruct-2507 fine-tuned on various curriculum strategies. The results reveal a clear, benchmark-dependent dichotomy. For function-completion tasks like HumanEval and HumanEval+, performance tends to degrade as the token limit increases beyond 4,096, suggesting that a larger generation space may encourage verbose, error-prone solutions. Conversely, for benchmarks like MBPP and MBPP+, a larger token limit is generally beneficial, indicating that their problem structures may require more extensive code to solve correctly.

dependent. We selected $\beta=0.01$ for our main experiments as it proved most effective on our primary evaluation benchmarks.

Impact of Training Temperature The sampling temperature manages the exploration-exploitation trade-off during training. The results, presented in Figure 7, indicate that a higher temperature of 1.0, which encourages greater exploration of diverse solutions, is optimal for HumanEval and HumanEval+. However, other benchmarks show different trends; MBPP, for example, peaks at a more conservative temperature of 0.7. This highlights that the optimal degree of exploration is also task-specific, and suggests that task-adaptive temperature scheduling could be a potential area for future work.

F IMPACT OF MAXIMUM COMPLETION TOKENS AT INFERENCE TIME

We analyzed the impact of the maximum completion token limit during inference on the fine-tuned Qwen3-4B model, with results presented in Figure 8. The findings reveal a clear, benchmark-dependent dichotomy. On function-completion tasks like HumanEval and HumanEval+, performance generally degrades as the token limit increases beyond 4,096. In stark contrast, benchmarks like MBPP and MBPP+ benefit from a larger generation space, with optimal results often found at 8,192 or 16,384 tokens.

This divergence suggests that for tasks requiring concise solutions, such as those in HumanEval, a larger token limit may encourage verbose and error-prone code. Conversely, the nature of MBPP problems may necessitate a longer generation process to fully develop the correct logic. This analysis underscores a critical point for standardized evaluation: the ideal setting for maximum completion tokens is highly contingent on the characteristics of the target benchmark.

G FULL BENCHMARK TABLES (QWEN2.5 & QWEN3-4B)

We report the complete benchmark results for all curriculum strategies on Qwen2.5 family (1.5B/3B/7B, including Coder variants) and Qwen3-4B-Instruct-2507. These tables expand the main figures by listing pass@1 on HumanEval/HumanEval+, MBPP/MBPP+, and average accuracy of CodeForces, LiveCodeBench v5, and CruxEval for every strategy in Table 6 and 7.

Consistent with the main text, the *C/E Weighted* strategy tends to be the top performer for the stronger Qwen3-4B model, improving over the base across all four code-function benchmarks (see Table 1 in the main paper for deltas). The full per-strategy breakdowns here allow exact comparison across OOD tasks as well.

Table 6: Comprehensive performance evaluation of all curriculum strategies across Qwen2.5-Instruct and Qwen2.5-Coder-Instruct models (1.5B, 3B, 7B). For each model size, the highest score on each benchmark is highlighted in **bold**.

Model	Strategy	HumanEval	HumanEval+	MBPP	MBPP+	CodeForces	LCBv5	CruxEval
Qwen2.	5-7B-Instruct							
	Base	82.93%	75.61%	63.20%	67.46%	8.54%	14.10%	57.75/58.009
	Uniform	82.93%	73.78%	67.40%	67.46%	12.36%	14.93%	57.88/59.389
	B/I Weighted	78.05%	76.83%	67.60%	69.58%	11.56%	15.89%	57.25/59.259
	C/E Weighted	79.27%	73.78%	66.20%	70.37%	10.89%	15.77%	57.13/55.509
	C/E Weighted (Rev)	84.15%	77.44%	69.00%	70.11%	8.24%	15.41%	57.88 /56.389
	Basic Only	82.32%	75.61%	66.20%	68.52%	12.29%	19.12%	55.63/57.509
	Edge Only	83.54%	76.22%	67.60%	70.63%	11.11%	17.08%	56.13/57.75
	Complex Only	84.15%	75.61%	69.00%	69.05%	12.95%	17.80%	57.25/56.38
Qwen2.	5-3B-Instruct							
_	Base	69.51%	61.59%	58.40%	64.81%	4.34%	5.02%	38.75/44.639
	Uniform	71.34%	63.41%	59.40%	63.49%	6.92%	8.72%	42.00/42.50
	B/I Weighted	69.51%	62.20%	59.00%	63.49%	7.21%	9.44%	42.38/ 46.75
	C/E Weighted	69.51%	62.80%	56.60%	63.76%	7.21%	7.17%	43.75 /44.50
	C/E Weighted (Rev)	70.12%	62.80%	57.00%	63.49%	6.92%	8.00%	43.63/42.50
	Basic Only	66.46%	59.15%	59.40%	64.02%	6.33%	6.09%	40.50/44.13
	Edge Only	71.34%	64.02%	58.20%	62.70%	6.11%	7.05%	43.13/42.63
	Complex Only	67.68%	60.37%	59.00%	64.81%	6.84%	6.33%	41.25/42.88
Owen2.	5-1.5B-Instruct							
	Base	58.54%	54.88%	46.80%	52.91%	2.65%	5.02%	38.63/30.88
	Uniform	60.98%	54.88%	50.00%	57.14%	3.68%	5.26%	37.13/ 33.75
	B/I Weighted	59.15%	54.27%	51.80%	57.94%	3.83%	4.54%	36.00/29.75
	C/E Weighted	60.98%	55.49%	49.40%	56.08%	3.61%	5.02%	34.75/32.38
	C/E Weighted (Rev)	56.71%	52.44%	50.40%	58.20%	4.49%	4.90%	34.00/31.75
	Basic Only	57.32%	53.05%	50.60%	58.20%	4.05%	4.66%	40.25 /33.00
	Edge Only	55.49%	50.61%	50.20%	56.08%	3.75%	4.42%	35.50/31.50
	Complex Only	59.76%	54.88%	51.80%	55.29%	3.46%	4.54%	36.13/33.38
02								
Qwenz.	5-Coder-7B-Instruct Base	85.98%	79.27%	75.60%	69.05%	10.89%	13.86%	66.38/66.13
	Uniform	85.76%	79.27%	77.20%	72.49%	13.98%	17.68%	66.50/66.38
	B/I Weighted	84.76%	78.66%	77.60%	71.96%	13.32%	17.44%	68.38/65.88
	C/E Weighted	87.80%	82.32%	76.20%	70.90%	14.94%	19.24%	66.25/67.13
	C/E Weighted (Rev)	88.41%	81.10%	75.00%	70.90%	13.98%	19.24%	68.63/65.00
					71.42%		19.12% 19.47%	
	Basic Only	85.98%	79.88%	76.20%		14.86%		67.50/66.38
	Edge Only	79.02%	81.07%	77.20%	71.96%	12.14%	19.12%	68.75 /66.00
	Complex Only	87.80%	80.49%	76.60%	70.90%	14.35%	18.16%	67.75/66.50
Qwen2.	5-Coder-3B-Instruct	70.270	75.000/	62.200	66 0201	2 000	0.900	E2 20/E2 7E
	Base	79.27%	75.00%	62.20%	66.93%	3.90%	9.80%	53.38/53.75
	Uniform	81.10%	76.83%	62.00%	67.20%	7.21%	10.75%	54.00/54.75
	B/I Weighted	81.71%	78.05%	61.40%	66.93%	6.70%	9.80%	54.25/53.38
	C/E Weighted	79.88%	76.83%	61.00%	67.46%	8.02%	10.51%	56.75 /53.50
	C/E Weighted (Rev)	82.32%	77.44%	62.00%	68.52%	8.17%	10.75%	52.63/55.13
	Basic Only	80.49%	76.22%	62.80%	66.67%	7.95%	13.14%	55.88/ 55.88
	Edge Only	79.27%	75.00%	62.60%	66.14%	7.21%	10.63%	53.75/53.25
	Complex Only	78.05%	73.78%	63.00%	67.72%	7.65%	10.63%	53.13/55.25
Qwen2.	5-Coder-1.5B-Instruc							
	Base	68.29%	63.41%	52.60%	63.49%	2.06%	3.46%	44.38/36.38
	Uniform	71.34%	65.24%	52.80%	62.96%	4.56%	4.42%	44.75/36.38
	B/I Weighted	72.56%	64.02%	55.80%	62.70%	4.19%	4.66%	45.13/35.75
	C/E Weighted	71.34%	66.65%	54.60%	62.96%	3.46%	4.42%	45.13/38.00
	C/E Weighted (Rev)	72.56%	64.20%	54.20%	62.96%	3.38%	4.18%	45.25 /37.00
	Basic Only	70.12%	64.02%	54.00%	64.76%	4.49%	4.66%	43.25/36.00
	Edge Only	72.56%	67.10%	53.60%	62.17%	4.56%	5.02%	44.86/35.63
	Complex Only	71.34%	66.46%	53.20%	63.49%	3.31%	4.54%	43.75/37.36

H ADDITIONAL RESULTS ON GEMMA2-2B-IT

This appendix provides the full curriculum comparison for Gemma2-2B-IT as in Table 8. Unlike larger or stronger models, Gemma2-2B-IT exhibits curriculum fragility: most curricula depress performance, consistent with the observation in the main text that sparse reward signals can cause collapse for less-capable models. In contrast, *Basic Only*—a fundamentals-first schedule—yields the most reliable gains among the tested strategies.

Table 7: Comprehensive performance evaluation of all curriculum strategies on Qwen3-4B-Instruct-2507. The highest score on each benchmark is highlighted in **bold**. The performance of Qwen3-Coder-30B-A3B-Instruct is included to enable comparison against a leading code-specialized model.

Model	Strategy	HumanEval	HumanEval+	MBPP	MBPP+	CodeForces	LCBv5	CruxEval
Qwen3-	-Coder-30B-A3B-Insti	ruct						
	Base	94.51%	86.59%	73.80%	75.13%	29.65%	37.63%	81.75/79.25%
Qwen3-	-4B-Instruct-2507							
•	Base	89.02%	78.66%	52.60%	56.61%	33.63%	32.02%	78.25/ 77.75 %
	Uniform	88.41%	80.09%	35.30%	53.70%	31.86%	31.96%	79.37/75.75%
	B/I Weighted	89.63%	81.09%	28.00%	52.38%	33.04%	33.81%	79.50/75.38%
	C/E Weighted	91.46%	82.92%	55.20%	58.73%	31.79%	31.54%	81.12 /75.25%
	C/E Weighted (Rev)	89.63%	80.48%	36.20%	35.98%	34.66%	31.66%	79.50/76.00%
	Basic Only	89.63%	79.87%	39.80%	56.34%	33.11%	31.90%	78.50/75.00%
	Edge Only	89.63%	79.88%	47.20%	56.61%	31.86%	30.59%	80.25/74.00%
	Complex Only	90.85%	80.48%	28.60%	51.85%	30.61%	31.30%	80.37/76.37%

These results reinforce our capability-dependent view of curriculum design: for weaker models, emphasizing simpler tiers is a prerequisite for successful fine-tuning, whereas complex-focused or mixed curricula can be harmful.

Table 8: Performance comparison for Gemma2-2B-IT across all curriculum strategies. Scores are colored and bolded based on their deviation from the Base strategy (blue for higher, red for lower).

Strategy	HumanEval	HumanEval+	MBPP	MBPP+	CodeForces	LCBv5	CruxEval
Base	42.07%	34.76%	41.20%	47.09%	2.21%	4.30%	37.50/26.88%
Uniform	39.02%	31.09%	33.80%	39.95%	0.22%	3.58%	33.00/26.63%
B/I Weighted	35.98%	32.32%	35.60%	42.06%	0.22%	4.30%	38.63/27.25%
C/E Weighted	41.46%	34.15%	39.20%	48.41%	0.22%	3.94%	36.63/26.75%
C/E Weighted (Rev)	40.86%	35.37%	40.20%	44.44%	0.44%	3.94%	35.63/26.75%
Basic Only	44.51%	37.20%	38.60%	46.83%	1.77%	3.94%	39.88/27.88%
Edge Only	39.63%	35.37%	38.00%	46.03%	0.22%	4.30%	39.00/28.13%
Complex Only	42.07%	36.59%	37.00%	45.77%	2.21%	2.87%	35.63/27.55%