

REVEAL: SELF-EVOLVING CODE AGENTS VIA RELIABLE SELF-VERIFICATION

Anonymous authors
 Paper under double-blind review

ABSTRACT

Reinforcement learning with verifiable rewards (RLVR) has advanced the reasoning capabilities of large language models. However, existing methods rely solely on outcome rewards, without explicitly optimizing verification or leveraging reliable signals from realistic environments, leading to unreliable self-verification and limited test-time scaling. To address this, we widen the verification-generation asymmetry by explicitly optimizing self-verification, making it a reliable driver of deeper test-time scaling. We introduce *ReVeal*, a multi-turn Reinforcement learning framework that evolves code generation through self-Verification and tool-based evaluation. ReVeal structures long-horizon reasoning as iterative generation-verification turns and incorporates TAPO for turn-level credit assignment, fostering the co-evolution of code and test generation. At inference, this strengthened self-verification enables the model to use self-constructed tests and tool feedback to continuously evolve code for **20+** turns on LiveCodeBench despite training on only three. It also significantly improves *Pass@k*, indicating stronger exploration that expands the reasoning boundaries of the base model. These findings highlight the promise of ReVeal as a scalable paradigm for RL training and test-time scaling, paving the way for more robust and autonomous AI agents.

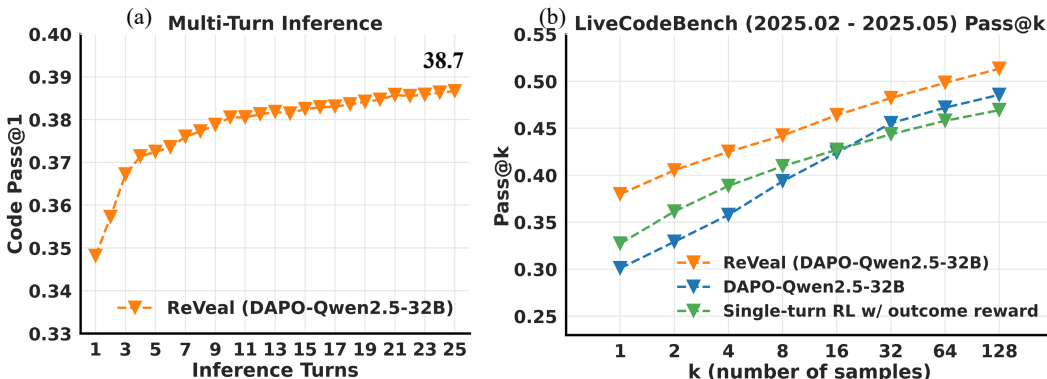


Figure 1: Performance of ReVeal on LiveCodeBench V6. (a) ReVeal enables effective test-time scaling, with Pass@1 accuracy improving from 34.8% at turn 1 to 38.7% at turn 25. (b) ReVeal (max_turn=10) consistently outperforms both the base model and the RL baseline in Pass@k, expanding the base model’s reasoning boundaries, which the RL baseline fails to achieve.

1 INTRODUCTION

Reinforcement learning with verifiable rewards (RLVR) has recently shown strong potential to enhance the reasoning abilities of large language models (LLMs) (DeepSeek-AI et al., 2025; OpenAI). A key factor behind this success is the emergence of reflection and self-verification, which allow models to iteratively refine their reasoning. Recent analyses identify the *verification-generation asymmetry* (i.e., easier to verify than to solve) as the underlying mechanism for these improvements and a key driver of test-time scaling (Wei, 2025; Setlur et al., 2025). However, current RLVR methods rely solely on outcome rewards without explicitly optimizing verification. This leads to unreliable

self-verification, where models often produce verbose, uninformative reflections or random guessing on hard problems, and limits the effectiveness of test-time scaling: prior studies show that reasoning performance plateaus once test-time compute exceeds the training horizon (Setlur et al., 2025).

Complex problem-solving, such as competitive programming, typically requires multiple iterations of verification and revision rather than being solved in a single attempt, making accurate feedback essential to guide refinement. This highlights the need for verification-driven multi-turn reasoning. Prior work has attempted this either by training a separate critic model to assess each attempt—without leveraging tool feedback and at the cost of added inference-time complexity (Xie et al., 2025)—or by relying on execution feedback against pre-existing public tests, which are rarely available in real-world scenarios (Gehring et al., 2025). As a result, these methods provide limited and non-generalizable verification, leaving self-verification unreliable and limiting sustained improvement.

To address these limitations, we propose **ReVeal**, a multi-turn RL framework that *explicitly optimizes self-verification*, thereby widening the verification-generation (V-G) asymmetry and fostering co-evolution of both capabilities during training. This enables models at inference to obtain reliable verification signals from realistic environments and iteratively refine their solutions, without needing to rely on pre-existing tests. The widened V-G gap allows verification to drive sustained improvements in generation, ultimately enabling deeper test-time scaling (Figure 2).

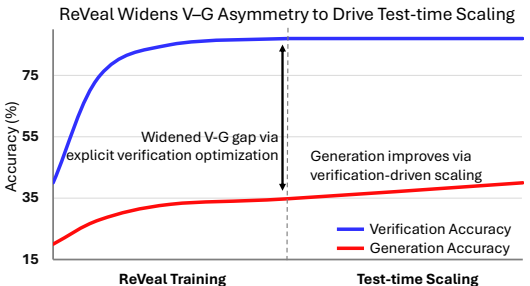


Figure 2: ReVeal expands the V-G gap.

Concretely, ReVeal structures long-horizon reasoning into iterative generation and verification turns. At each turn, the model generates candidate code and *self-verifies* its correctness by constructing test cases and invoking external tools (e.g., a Python interpreter) for execution. This closed loop yields actionable verification signals and fine-grained feedback, allowing the model to identify errors, revise strategies, and progressively refine its output across turns. For training, we attach *dense, turn-level rewards* that directly supervise both code quality and verification accuracy. To ensure robustness, ReVeal employs a *Turn-Aware Policy Optimization (TAPO)* tailored for the generation-verification interplay, assigning credit at the turn granularity and preventing reward gaming (e.g., generating trivial code to hack verification rewards). Unlike outcome-only RL methods, ReVeal makes verification itself an optimization target, turning verification signals into reliable drivers of improvement.

We evaluate ReVeal on the challenging LiveCodeBench benchmark (Jain et al., 2024). Notably, despite being trained on only three reasoning turns, ReVeal sustains continuous refinement for over 20+ inference turns, showing robust extrapolation beyond its training horizon and tackling problems previously unsolved. Furthermore, ReVeal significantly outperforms the base model in Pass@k by leveraging verification signals and tool feedback to guide more effective exploration, achieving an expansion of the underlying model’s reasoning boundaries that standard RL methods fail to reach. These results validate ReVeal as not only a practical framework for self-evolving code agents, but also as a general RL paradigm for tasks with verification-generation asymmetry, where explicitly optimizing verification unlocks reliable long-horizon reasoning.

2 METHODS

2.1 REVEAL FRAMEWORK

2.1.1 ITERATIVE GENERATION-VERIFICATION LOOP

ReVeal organizes long-horizon reasoning into an interleaved *generation-verification* loop with tool execution feedback, where verification itself is explicitly optimized to provide reliable signals for multi-turn refinement. As illustrated in Figure 3, we use a single policy for both generation and verification to reduce system complexity and cost and to enable cross-capability transfer, so that solutions and their verification strategies co-evolve under a shared training scheme. In the code-generation setting, *generation* produces candidate code, whereas *verification* synthesizes and executes

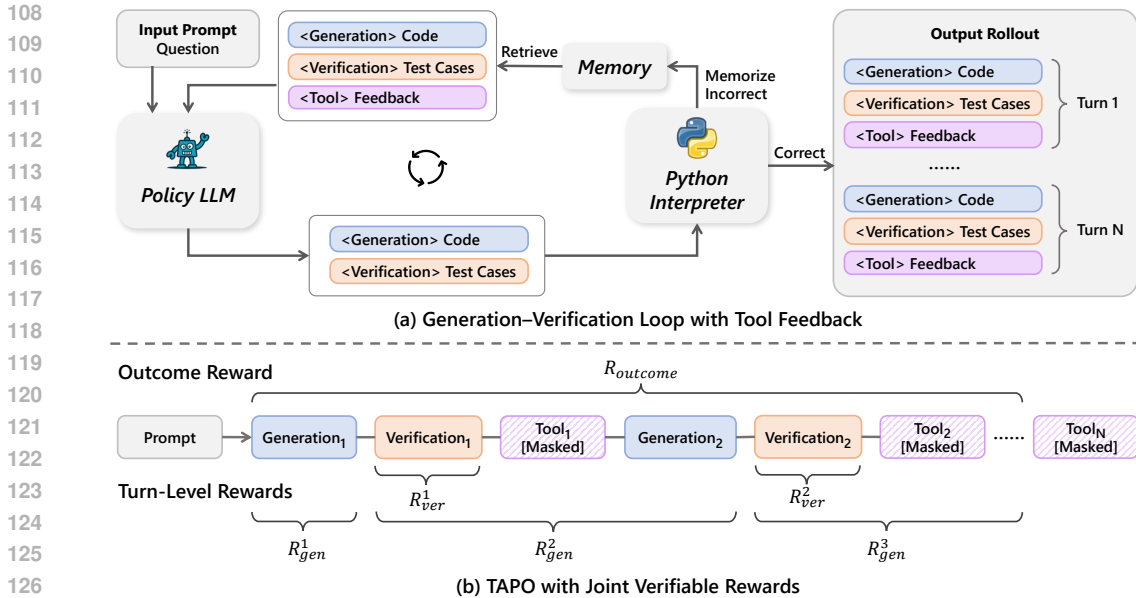


Figure 3: Illustration of ReVeal. (a) Iterative generation-verification loop with tool feedback. (b) TAPO with joint verifiable rewards: outcome, generation, and verification rewards.

tests to assess correctness. Fine-grained feedback from tool execution (e.g., Python interpreter) is appended to the rollout and conditions the next turn. The loop continues until a valid solution is found or a turn budget K is reached, enabling progressive refinement without external critics or predefined test cases.

Table 4 illustrates a multi-turn rollout under ReVeal’s structured prompting, which decouples generation, verification, and tool feedback into distinct segments. At each turn, the policy first reasons thoroughly and explores diverse reasoning patterns freely, then emits structured outputs: executable code in `<generation-answer>` and executable tests in `<verification-answer>`. As shown in the case study, after producing candidate code the model begins verification: it hypothesizes potential failure modes and edge conditions to propose diverse test cases. The `<tool-feedback>` section then records execution results, including runtime errors, invalid test cases, as well as the expected output, actual output, and pass/fail judgment for each valid test case. Based on this feedback, the model interprets traces and error messages, diagnoses underlying causes, and adjusts both its candidate code and its verification plan in the next turn. Full prompting and feedback templates are provided in Tables 5 and 8.

2.1.2 TOOL-AUGMENTED VERIFICATION

The interaction with external tools provides reliable, fine-grained supervisory signals that condition subsequent reasoning and enable systematic refinement of both code and verification strategies across turns. More importantly, tool interaction broadens exploration during reinforcement learning by revealing concrete failure modes, steering the policy into promising regions of the search space beyond a single attempt and helping it escape local optima. Empirically (see §3.3), this yields consistently higher pass@k than the base model.

During RL training, the `<tool-feedback>` section is excluded from the loss and used only as contextual input, which stabilizes optimization while preserving coherent multi-turn rollouts. To ensure feedback quality during training, we adopt a filtering mechanism: model-generated test cases are executed on candidate code only if they are verified against a golden solution. This guarantees that execution traces provide legitimate supervision, thereby improving feedback precision and guiding exploration toward correct solutions. At test time, no golden reference is available; all generated test cases are executed, making verification fully autonomous. This places a strong demand on the model’s ability to generate high-quality tests. To meet this demand, ReVeal adopts a novel RL algorithm that incentivize diverse and reliable test construction.

2.2 TURN-AWARE RL FOR THE GENERATION-VERIFICATION PARADIGM

Prior RLVR methods rely on outcome-only signals to optimize an entire long reasoning trace, but this provides imprecise credit to intermediate verification and often degenerates into blind reflection. Yet one may ask: *is the current paradigm fully sustain reliable verification and deeper test-time scaling?* Verification, however, is a non-trivial task: with well-designed verifiable rewards, a task can often be solved effectively. This motivates ReVeal to explicitly optimize verification with hard-to-hack rewards, which widen the verification-generation gap. At test time, this asymmetry becomes an asset: easier and more reliable verification signals can effectively guide the harder generation process to evolve over many turns.

2.2.1 JOINT VERIFIABLE REWARDS

To jointly train generation and verification, ReVeal decomposes the reward into three complementary components (Fig. 3b): an *outcome reward* supervising the final solution, a *generation reward* capturing improvements across generation turns, and a *verification reward* evaluating the quality of generated tests. This design naturally links the two roles in a co-evolutionary loop.

Outcome reward. The outcome reward shapes the entire reasoning chain by the final solution quality:

$$r_{\text{outcome}} = r_{\text{format}} + r_{\text{passrate}}, \quad (1)$$

where the format reward r_{format} ensures that the model produces well-formed generation and verification blocks,

$$r_{\text{format}} = \begin{cases} 1, & \text{if the output format is correct,} \\ -1, & \text{otherwise,} \end{cases} \quad (2)$$

and $r_{\text{passrate}} = 5 \times \text{passrate}$ measures final code accuracy with $\text{passrate} \in [0, 1]$, giving $r_{\text{outcome}} \in [-1, 6]$.

Generation reward. For each generation turn k (odd), we compute the pass rate r_{passrate}^k of the code produced and define:

$$r_{\text{gen}}^k = \begin{cases} r_{\text{passrate}}^1, & k = 1, \\ \text{abs} \cdot r_{\text{passrate}}^k + \text{imp} \cdot (r_{\text{passrate}}^k - r_{\text{passrate}}^{k-2}), & k \geq 3, \end{cases}$$

where abs and imp weight absolute accuracy and iterative improvement. We set $\text{abs} = 0$, $\text{imp} = 1$ so that the reward encourages real improvements in code accuracy across turns.

Verification reward. For each verification turn k (even), we reward the proportion of generated tests that succeed when executed on a golden code:

$$r_{\text{ver}}^k = \frac{\#\{\text{test cases in turn } k \text{ that pass}\}}{\#\{\text{test cases generated in turn } k\}}. \quad (3)$$

2.2.2 TURN-AWARE POLICY OPTIMIZATION

Preliminaries. Our algorithm builds on the Proximal Policy Optimization (PPO) framework (Schulman et al., 2017), an on-policy actor-critic method that optimizes a clipped surrogate objective for stable updates. PPO typically estimates token-level advantages using Generalized Advantage Estimation (GAE) (Schulman et al., 2018):

$$\hat{A}_t^{\text{GAE}(\gamma, \lambda)} = \sum_{l=0}^{\infty} (\gamma \lambda)^l (r_{t+l} + \gamma V_{t+l+1} - V_{t+l}), \quad (4)$$

where $\gamma \in [0, 1]$ is the discount factor and $\lambda \in [0, 1]$ controls the bias-variance trade-off.

Turn-Aware Policy Optimization. Building on our structured reward design, we introduce *Turn-Aware Policy Optimization* (TAPO), which preserves the PPO actor-critic framework but replaces GAE-based advantages with a *turn-aware* return. TAPO leverages the critic to efficiently bootstrap from both token-level Monte Carlo returns and turn-level returns, enabling stable learning across these two reward granularities.

216 1. **Token-level return.** We set $\lambda = 1$ and $\gamma = 1$ (pure Monte Carlo). For token step t with final step
217 T :

$$218 R_t = \sum_{l=0}^{T-t} r_{t+l} = r_t + R_{t+1}, \quad R_{T+1} = 0. \quad (5)$$

221 2. **Turn-level return.** To mitigate adversarial reward gaming (e.g., generating trivial code that
222 hacks the verification reward), we introduce a *turn-level* return tailored to the generation-verification
223 interplay. Specifically, (i) each generation reward is assigned both to its own generation turn and to
224 the immediately preceding verification turn, and (ii) each verification reward is confined strictly to
225 its own verification turn. This design prevents reward hacking by ensuring that generation turns are
226 rewarded solely based on code quality, rather than verification success. Let $\{t_1, \dots, t_K\}$ denote the
227 token indices at which each turn ends (alternating generation and verification), and define:

$$228 R^{\text{turn}}(t_k) = \begin{cases} r_{\text{gen}}^k, & \text{if turn } k \text{ is generation,} \\ r_{\text{ver}}^k + R^{\text{turn}}(t_{k+1}), & \text{if turn } k \text{ is verification,} \end{cases} \quad R^{\text{turn}}(t_{K+1}) = 0. \quad (6)$$

231 For token t , let $\tau(t) = \min\{t_k \mid t_k \geq t\}$ and define

$$232 R_t^{\text{turn}} = \begin{cases} R^{\text{turn}}(\tau(t)), & \text{if } \tau(t) \text{ exists,} \\ 0, & \text{otherwise.} \end{cases} \quad (7)$$

235 3. **Turn-aware return.** The final return combines the two levels:

$$236 \tilde{R}_t = R_t + R_t^{\text{turn}}, \quad A_t = \tilde{R}_t - V_t, \quad (8)$$

238 where V_t is the critic model’s estimate at step t . These advantages A_t then replace the standard GAE
239 estimates in the PPO objective, completing the TAPO update.

241 **Discussion.** TAPO provides sharper supervision than outcome-only methods by explicitly assigning
242 credit at both token and turn levels. It integrates outcome rewards, which keep the process aligned
243 with final correctness, and turn-level signals, which provide dense supervision for progressive
244 refinement. This structure establishes a feedback loop: stronger tests expose errors that drive code
245 improvements, which are then reinforced by the generation reward, while improved code raises the
246 bar for verification, pushing the model to generate richer and more challenging tests. By design,
247 TAPO prevents reward gaming and turns this loop into stable co-evolution of code and tests. Crucially,
248 TAPO is a *general* credit-assignment algorithm, applicable to any reasoning task with verifiable
249 rewards for both generation and verification.

250 3 EXPERIMENTS

251 3.1 SETTINGS

254 **Dataset** We construct our training dataset from TACO (Li et al., 2023), a large-scale corpus
255 comprising 26,443 algorithmic programming problems sourced from competitive programming
256 platforms such as LeetCode (LLC, 2015) and Codeforces (Codeforces, 2025). Each problem consists
257 of a natural language description, golden solutions, and multiple test cases.

258 To address noise in the raw dataset, we first filter out problems containing unsupported content types,
259 specifically those tagged with interactive or image elements. To ensure testability and correctness,
260 we process two types of test case format, function-based tests and standard input/output tests, into a
261 unified structure compatible with our code execution environment. We then execute each test case
262 against the first available golden solution in our execution environment. Problems where the golden
263 code fails to pass all associated test cases are discarded. After preprocessing, we retain a high-quality
264 dataset of 11,151 problems for training and 509 problems for testing.

266 **Models and Training Details** We adopt DAPO-Qwen-32B (Yu et al., 2025) as our base model,
267 which is reinforced with mathematical data, and we continue RL training on code datasets to adapt
268 its reasoning capabilities to coding tasks. Our models are trained using Verl (Sheng et al., 2024)
269 framework on 8/16 AMD Mi300x GPUs. The RL training process follows the hyperparameter
settings listed in Table 2. We set maximum turns to 3 during RL training.

Table 1: Performance comparison of ReVeal with baseline methods on LiveCodeBench V6 and CodeContests. Pass@1 indicates the success rate; Δ_{\uparrow} and Δ_{\downarrow} represent the percentages of incorrect solutions corrected and correct solutions degraded after revision, respectively.

Model	LiveCodeBench V6			CodeContests		
	Pass@1	Δ_{\uparrow}	Δ_{\downarrow}	Pass@1	Δ_{\uparrow}	Δ_{\downarrow}
<i>Existing Baselines</i>						
Qwen2.5-32B-Instruct	24.8	-	-	13.3	-	-
DAPO-Qwen2.5-32B	31.1	-	-	18.5	-	-
Qwen2.5-Coder-32B-Instruct	29.5	-	-	14.6	-	-
w/ critic×5 Qwen2.5-Coder	29.6	2.14	3.04	-	-	-
w/ critic×5 GPT-4o	32.9	4.82	2.50	-	-	-
w/ critic×5 CTRL	33.4	3.75	0.89	-	-	-
<i>RL based on DAPO-Qwen2.5-32B</i>						
Single-turn RL	32.8	-	-	21.0	-	-
ReVeal×25	38.7	7.50	0.0	33.6	15.69	0.0
<i>Ablation Study: TAPO with Joint Verifiable Rewards</i>						
ReVeal×8 w/ outcome reward	36.1	4.69	1.32	27.4	9.24	2.36
ReVeal×8 w/ TAPO with joint rewards	37.7	5.62	0.0	30.4	12.30	0.0

Evaluation We evaluate ReVeal on two code-generation benchmarks: LiveCodeBench (LCB) V6 (2025.02–2025.05) (Jain et al., 2024) and CodeContests (Li et al., 2022). The evaluation process follows the hyperparameter configuration specified in Table 3. Although training is performed with a maximum of 3 turns, we evaluate the model under extended turn settings (8 and 25 turns) to assess its generalization to longer reasoning horizons and test-time scaling performance.

We use Pass@1 to measure the success rate of the model’s final code solutions. To evaluate the model’s verification and self-correction capabilities, we introduce two additional metrics: Δ_{\uparrow} denotes the fraction of initially incorrect solutions that become correct after revision, and Δ_{\downarrow} denotes the fraction of initially correct solutions that become incorrect after revision. In line with recent work (Yue et al., 2025), we use Pass@k up to $k = 128$ to assess whether ReVeal can push the reasoning boundaries beyond the base model, with at most 10 generation–verification turns per example.

Memory Mechanism for Context Management To improve inference efficiency under extended multi-turn rollouts, we use a short-term memory mechanism that retains only the last three turns as context, which prevents excessive context growth without hurting accuracy (detail in Appendix D).

Code Execution Tool We use Code Judge¹ as our code execution environment. Code Judge supports both function-based and standard input-output test case formats through a consistent interface. Designed for scalability and robustness, it enables efficient long-batch execution through multi-processing and provides reliable code evaluation.

Baselines We compare ReVeal against following baselines: (1) *Base*: base models without code-specific RL training; (2) *CTRL* (Xie et al., 2025) + *Qwen2.5-Coder-32B-Instruct*: five-turn critic–revision with a dedicated critic model; results cited from the original paper (evaluated on LCB 24.08–24.11); (3) *Single-turn RL with outcome reward*: RL with outcome-only rewards under standard <think>-<answer> prompting template without any external tool calls.

3.2 MAIN RESULTS

Table 1 shows that single-turn RL (outcome-only, no explicit optimization of self-verification or tool use) improves Pass@1 over the base models. ReVeal goes further by explicitly optimizing verification and enabling deeper inference, it surpasses the single-turn RL baseline by a wide margin. Beyond deeper-turn gains, ReVeal also achieves higher Pass@1 at turn 1 (34.8%) than the single-turn

¹<https://github.com/0xWJ/code-judge>

RL baseline under equal inference budget on LCB V6, indicating that multi-turn training (3 turns) transfers exploration benefits into a stronger policy and that increasing training depth may further amplify gains.

ReVeal significantly outperforms critic-based methods such as CTRL. While critic models tailored for code tasks can be paired with policy models for multi-turn critique and revision, ReVeal employs a single policy model that self-verifies and iteratively refines its own outputs, yet achieves superior results, highlighting the benefit of jointly optimizing generation and verification. Specifically, ReVeal attains larger correction rates with near-zero degradation, demonstrating highly robust and reliable capabilities in self-verification, critique, and revision. (CTRL numbers are cited from earlier LCB version; see Table 9 for a V5 comparison on Qwen2.5-32B-Instruct.)

Ablation studies confirm the benefit of TAPO with joint verifiable rewards: at the same turn budget it yields higher Pass@1, increases Δ_{\uparrow} , and suppresses Δ_{\downarrow} compared to outcome-only training. In contrast, outcome-only rewards exhibit higher Δ_{\downarrow} , indicating that insufficiently optimized verification can drive incorrect revisions.

More Experiments To validate effectiveness and scalability across models, we evaluate ReVeal on another base model, Qwen2.5-32B-Instruct. The detailed results are in Appendix Table 9 and Figure 5. ReVeal outperforms single-turn RL baseline by 4.1%, which demonstrates the effectiveness of ReVeal. To demonstrate that the performance improvements are statistically significant, we repeated the experiment 8 times and reported the *mean* \pm *std* in the Appendix Table 7. Across models of varying capability, ReVeal remains effective. With the stronger DAPO-Qwen2.5-32B backbone, ReVeal unlocks greater headroom: accuracy continues to improve with deeper inference turns and surpasses outcome-only RL by a wider margin. This underscores ReVeal’s potential on stronger backbones.

To verify that the performance improvements originate from explicit optimization rather than mere tool use, we conduct comparative experiments with ReAct prompting (Shinn et al., 2023), showing that ReVeal significantly outperforms ReAct across all turn budgets. Without explicit optimization for verification, ReAct fails to sustain effective deep multi-turn refinement. See Table 6 for details.

3.3 ANALYSIS

ReVeal Enables Test-time Scaling into Deeper Inference Regimes. As shown in Figure 1 (a), ReVeal enables effective test-time scaling through iterative generation and verification. Although the model is trained with a maximum of three reasoning turns, it continues to improve its solutions when more turns are allowed at inference time, leading to progressively higher code accuracy. For instance, Pass@1 increases from 34.8% at turn 1 to 36.7% at turn 3, and further rises to 38.7% by turn 25 for LiveCodeBench. This compellingly demonstrates how reliable self-verification and iterative environment feedback can enable compute scaling into deeper inference regimes, allowing ReVeal to solve previously intractable problems and evolve novel solutions. As a result, ReVeal supports self-improvement beyond the training horizon, enabling strong generalization in long-horizon reasoning during inference. Furthermore, these newly discovered solutions can be distilled back into the code LLM to further enhance its reasoning capabilities through continued training.

ReVeal Pushes Beyond the Reasoning Boundaries of the Base Model. We compare DAPO-Qwen2.5-32B and single-turn RL baseline with ReVeal using Pass@k metrics on LiveCodeBench. As shown in Figure 1 (b), the RL baseline outperforms the base model when $k < 32$, but its performance gain gradually diminishes as k increases. In contrast, ReVeal consistently outperforms both the base model and the RL baseline across all k values from 1 to 128, demonstrating its ability to surpass the reasoning boundaries beyond the base model. We attribute this improvement to ReVeal’s verification-driven exploration: tool-assisted verification provides targeted, execution-based feedback and precise judgments that guide the model to explore better solutions more effectively. With this enhanced exploration capability, the model continually self-evolves and grows beyond its initial reasoning capability during RL training. We believe this approach offers a promising path towards developing self-evolving agents with stronger reasoning capabilities.

ReVeal Co-evolves the Model’s Generation and Verification Capabilities. Figure 4(b, d) illustrates the co-evolution of the model’s code and test case generation capabilities. As shown in

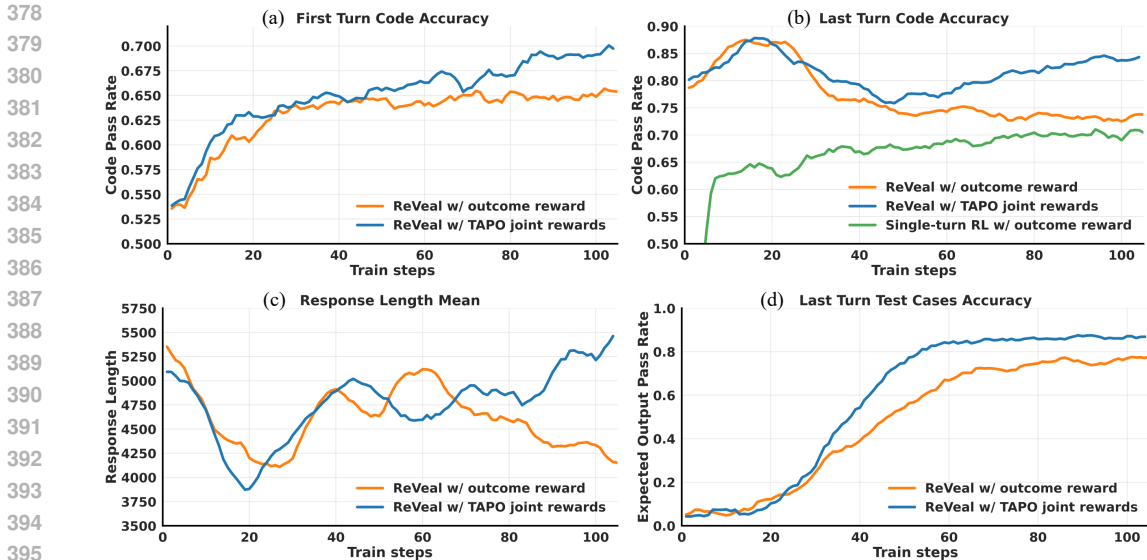


Figure 4: Training curves for (a) first-turn code accuracy, (b) last-turn code accuracy, (c) mean response length, and (d) last-turn test-case accuracy across three training methods. **Note on (b):** the dip before step 40 is due to expanded evaluation coverage: as format score reaches 0.9 around step 40, more problems enter the evaluation set, temporarily lowering accuracy.

Figure 4(b), once the format is learned, final code accuracy steadily improves during training and significantly surpasses the single-turn RL baseline. Moreover, comparing Figure 4(a) and (b) reveals that final solutions consistently outperform those generated at Turn 1, with the performance gap widening over time. This trend indicates that as the model’s verification ability strengthens, multi-turn refinement enables the exploration of better solutions, progressively enhancing its capacity to generate and refine code. After the format is learned, test-case accuracy rises substantially from about 50% at step 40 to nearly 88%, as shown in Figure 4 (d). Additionally, for correctly generated test cases, the model achieves over 85% accuracy in judging code correctness. This demonstrates that during inference, the model can reliably generate valid test cases and effectively leverage tool to produce accurate verification signals, which are critical for continuous improvements in code quality. These results provide strong evidence that ReVeal jointly and effectively optimizes both generation and verification, enabling the model to evolve its reasoning capabilities throughout training.

The Effectiveness of TAPO with Joint Verifiable Rewards. As shown in Table 1, TAPO with joint rewards further enhances multi-turn performance compared to relying solely on outcome rewards. The training curves in Figure 4(a,b) show TAPO with joint rewards achieves more stable and consistent per-turn code gains, and Figure 4(d) shows it achieves higher test-case accuracy, indicating that explicitly optimizing verification yields higher-quality tests and more effective reasoning in code generation task. These benefits amplify in longer-sequence and harder verification scenarios. On the stronger DAPO-Qwen2.5-32B backbone with longer chains, dense turn-level supervision yields larger gains than on Qwen2.5-32B-Instruct with much shorter chains (see Table 9 and Figure 5). This is because outcome-only signals are too coarse for extremely long chains, providing imprecise credit to intermediate verification steps. Furthermore, in more challenging verification scenarios, such fine-grained supervision becomes increasingly essential, offering richer learning signals to enhance the model’s verification capabilities.

4 RELATED WORK

4.1 TOOL-AUGMENTED REASONING

Tool-integrated reasoning enables large language models (LLMs) to leverage external tools, such as search engines or code interpreters, to overcome inherent limitations in domain knowledge and mathematical operations. Early approaches demonstrated the benefits of tool integration via prompt

432 engineering (Yao et al., 2023; Chen et al., 2023; Shinn et al., 2023) and supervised fine-tuning (Gou
433 et al., 2024). ReAct (Yao et al., 2023) and Reflexion (Shinn et al., 2023), which interleave reasoning
434 with acting or verbal self-critique to iteratively refine solutions under tool access. These approaches
435 highlight the value of interactive signals, but they typically rely on prompt heuristics. More recently,
436 multi-turn RL has been adopted to further enhance this capability on various reasoning tasks (Jin
437 et al., 2025; Feng et al., 2025; Li et al., 2025). For example, Search-R1 (Jin et al., 2025) incorporates
438 multi-turn interactions with a search engine to retrieve relevant contextual information during RL
439 training. ReTool (Feng et al., 2025) and ToRL (Li et al., 2025) enable multi-turn code execution to
440 support mathematical reasoning. Building on the promising potential of tool-integrated RL, Agent-
441 R1 (Ouyang et al., 2025) introduces an open-source RL framework capable of supporting multi-turn,
442 customization tool invocations.

443 Despite their effectiveness, most tool-augmented RL methods are predominantly outcome-driven: they
444 rely on task success or failure as the sole training signal and do not explicitly optimize verification or
445 assign credit across turns. Likewise, prompt-only agents lack turn-level, verifiable supervision, which
446 can make self-verification unreliable on harder problems and limit sustained test-time improvement.
447 Unlike prior tool-augmented works, ReVeal treats verification itself as a first-class optimization target
448 alongside generation, and introduces Turn-Aware Policy Optimization rewards (TAPO) to provide
449 fine-grained credit to both generation and verification turns. Our approach and prior tool-augmented
450 methods are orthogonal and complementary, and can be naturally combined to further enhance
451 reasoning capability.

452 4.2 SELF-VERIFICATION OF LLMs

453

454 Enabling LLMs to iteratively refine their outputs is critical for enhancing their reasoning capabilities.
455 However, LLMs typically lack reliable self-judgment (Huang et al., 2024). One common solution
456 is to introduce a separate critic model to verify the output of the policy model (Zhang et al., 2025;
457 Xie et al., 2025). For example, CTRL (Xie et al., 2025) uses RL to train a critic model for code
458 completion tasks. Although effective, these approaches incur the cost and complexity of maintaining
459 and coordinating two distinct models.

460 An alternative strategy is to enable one single model to generate outputs and self-verify them. In
461 mathematical reasoning, (Xiong et al., 2025) synthesizes long chains of thought that incorporate
462 "self-reward" and "self-correction" signals as seed data for supervised fine-tuning, and then further
463 enhances this ability via RL. In the code domain, execution feedback effectively verifies code
464 correctness and provides useful information for fixing errors. RLEF (Gehring et al., 2025) performs
465 multi-turn code generation and verification with an integrated code execution tool; however, it
466 depends on publicly available test cases, limiting its applicability.

467 In contrast, ReVeal advances self-verification by having the model generate its own high-quality test
468 cases on the fly. By explicitly crafting and executing these tests, ReVeal eliminates the dependency
469 on pre-existing test suites and improves applicability to real-world software systems.

471 5 CONCLUSION

472

473 We presented REVEAL, a multi-turn reinforcement learning (RL) framework that makes verification a
474 first-class optimization target alongside generation and organizes long reasoning chains into iterative
475 generation-verification turns with tool feedback. Using TAPO with joint verifiable rewards, REVEAL
476 equips LLMs with strong verification capabilities and demonstrates the surprising power of enabling
477 code LLMs to self-evolve—both during RL training, where it pushes boundaries beyond the base
478 model, and at test time, where multi-turn generation and verification continually refine outputs, even
479 up to 20+ inference turns. This compellingly demonstrates that REVEAL can enable compute scaling
480 into deeper inference regimes, allowing it to solve previously intractable problems and evolve novel
481 solutions. Furthermore, these newly discovered solutions can be distilled back into the code LLM to
482 further enhance its reasoning capabilities through continued training.

483 Although we demonstrate REVEAL on code tasks, its general concept of generation-verification,
484 TAPO, and turn-level reward design can be applied to any domain with verifiable rewards for both
485 generation and verification and that exhibits verification asymmetry, offering a promising blueprint
for future advances in self-improving, more robust, and autonomous AI agents.

REFERENCES

- 486
487
488 Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W. Cohen. Program of thoughts prompting:
489 Disentangling computation from reasoning for numerical reasoning tasks, 2023. URL <https://arxiv.org/abs/2211.12588>.
490
- 491 Codeforces. Codeforces - competitive programming platform, 2025. URL <https://codeforces.com/>. Accessed: 2025-03-18.
492
- 493 DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu,
494 Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu,
495 Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao
496 Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan,
497 Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao,
498 Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding,
499 Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang
500 Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong,
501 Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao,
502 Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang,
503 Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang,
504 Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L.
505 Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang,
506 Shuiping Yu, Shunfeng Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shaoqing Wu, Shengfeng
507 Ye, Tao Yun, Tian Pei, Tianyu Sun, T. Wang, Wangding Zeng, Wanbiao Zhao, Wen Liu, Wenfeng
508 Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan
509 Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xinyu Yang,
510 Xinyuan Li, Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin, Xiaojin Shen, Xiaosha Chen,
511 Xiaowen Sun, Xiaoxiang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia Shan, Y. K. Li,
512 Y. Q. Wang, Y. X. Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Wang,
513 Yi Yu, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang, Yixuan Tan,
514 Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yuduan Wang, Yue Gong, Yuheng Zou, Yujia
515 He, Yunfan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanhong
516 Xu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian Ma, Ying Tang, Yukun Zha,
517 Yuting Yan, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda Xie, Zhengyan Zhang,
518 Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li,
519 Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu Zhang, and Zhen
520 Zhang. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025.
521 URL <https://arxiv.org/abs/2501.12948>.
522
- 523 Jiazhan Feng, Shijue Huang, Xingwei Qu, Ge Zhang, Yujia Qin, Baoquan Zhong, Chengquan Jiang,
524 Jinxin Chi, and Wanjun Zhong. Retool: Reinforcement learning for strategic tool use in llms, 2025.
525 URL <https://arxiv.org/abs/2504.11536>.
526
- 527 Jonas Gehring, Kunhao Zheng, Jade Copet, Vegard Mella, Quentin Carbonneaux, Taco Cohen, and
528 Gabriel Synnaeve. Rlef: Grounding code llms in execution feedback with reinforcement learning,
529 2025. URL <https://arxiv.org/abs/2410.02089>.
530
- 531 Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen, Yujiu Yang, Minlie Huang, Nan Duan, and
532 Weizhu Chen. Tora: A tool-integrated reasoning agent for mathematical problem solving, 2024.
533 URL <https://arxiv.org/abs/2309.17452>.
534
- 535 Jie Huang, Xinyun Chen, Swaroop Mishra, Huaixiu Steven Zheng, Adams Wei Yu, Xinying Song,
536 and Denny Zhou. Large language models cannot self-correct reasoning yet, 2024. URL <https://arxiv.org/abs/2310.01798>.
537
- 538 Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando
539 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>.
- 538 Bowen Jin, Hansi Zeng, Zhenrui Yue, Jinsung Yoon, Sercan Arik, Dong Wang, Hamed Zamani, and
539 Jiawei Han. Search-r1: Training llms to reason and leverage search engines with reinforcement
learning, 2025. URL <https://arxiv.org/abs/2503.09516>.

- 540 Rongao Li, Jie Fu, Bo-Wen Zhang, Tao Huang, Zhihong Sun, Chen Lyu, Guang Liu, Zhi Jin, and
541 Ge Li. Taco: Topics in algorithmic code generation dataset, 2023. URL [https://arxiv.org/
542 abs/2312.14852](https://arxiv.org/abs/2312.14852).
- 543 Xuefeng Li, Haoyang Zou, and Pengfei Liu. Torl: Scaling tool-integrated rl, 2025. URL [https://
544 arxiv.org/abs/2503.23383](https://arxiv.org/abs/2503.23383).
- 546 Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom
547 Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, Thomas Hubert, Peter Choy, Cyprien
548 de Masson d’Autume, Igor Babuschkin, Xinyun Chen, Po-Sen Huang, Johannes Welbl, Sven Gowal,
549 Alexey Cherepanov, James Molloy, Daniel Mankowitz, Esme Sutherland Robson, Pushmeet Kohli,
550 Nando de Freitas, Koray Kavukcuoglu, and Oriol Vinyals. Competition-level code generation with
551 alphacode. *arXiv preprint arXiv:2203.07814*, 2022.
- 552 LeetCode LLC. Leetcode: Online coding platform for technical interview preparation, 2015. URL
553 <https://leetcode.com/>. Accessed: 2025-05-16.
- 554 OpenAI. Learning to reason with llms.
555 [urlhttps://openai.com/index/learning-to-reason-with-llms/](https://openai.com/index/learning-to-reason-with-llms/). Accessed: 15 March 2025.
- 557 Jie Ouyang, Ruiran Yan, Yucong Luo, Mingyue Cheng, Qi Liu, Zirui Liu, Shuo Yu, and Daoyu
558 Wang. Training powerful llm agents with end-to-end reinforcement learning, 2025. URL [https://
559 github.com/0russwest0/Agent-RL](https://github.com/0russwest0/Agent-RL).
- 560 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
561 optimization algorithms, 2017. URL <https://arxiv.org/abs/1707.06347>.
- 562 John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-dimensional
563 continuous control using generalized advantage estimation, 2018. URL [https://arxiv.org/
564 abs/1506.02438](https://arxiv.org/abs/1506.02438).
- 566 Amrith Setlur, Matthew Y. R. Yang, Charlie Snell, Jeremy Greer, Ian Wu, Virginia Smith, Max
567 Simchowitz, and Aviral Kumar. e3: Learning to explore enables extrapolation of test-time compute
568 for llms, 2025. URL <https://arxiv.org/abs/2506.09026>.
- 569 Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng,
570 Haibin Lin, and Chuan Wu. Hybridflow: A flexible and efficient rlhf framework. *arXiv preprint
571 arXiv:2409.19256*, 2024.
- 572 Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik Narasimhan, and Shunyu
573 Yao. Reflexion: Language agents with verbal reinforcement learning, 2023.
- 574 Jason Wei. Asymmetry of verification and the verifier’s law. [https://www.jasonwei.net/
575 blog/asymmetry-of-verification-and-verifiers-law](https://www.jasonwei.net/blog/asymmetry-of-verification-and-verifiers-law), 2025.
- 576 Zhihui Xie, Jie chen, Liyu Chen, Weichao Mao, Jingjing Xu, and Lingpeng Kong. Teaching language
577 models to critique via reinforcement learning, 2025. URL [https://arxiv.org/abs/2502.
578 03492](https://arxiv.org/abs/2502.03492).
- 581 Wei Xiong, Hanning Zhang, Chenlu Ye, Lichang Chen, Nan Jiang, and Tong Zhang. Self-
582 rewarding correction for mathematical reasoning, 2025. URL [https://arxiv.org/abs/
583 2502.19613](https://arxiv.org/abs/2502.19613).
- 584 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao.
585 React: Synergizing reasoning and acting in language models, 2023. URL [https://arxiv.
586 org/abs/2210.03629](https://arxiv.org/abs/2210.03629).
- 587 Qiying Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Tiantian Fan, Gaohong
588 Liu, Lingjun Liu, Xin Liu, Haibin Lin, Zhiqi Lin, Bole Ma, Guangming Sheng, Yuxuan Tong, Chi
589 Zhang, Mofan Zhang, Wang Zhang, Hang Zhu, Jinhua Zhu, Jiase Chen, Jiangjie Chen, Chengyi
590 Wang, Hongli Yu, Weinan Dai, Yuxuan Song, Xiangpeng Wei, Hao Zhou, Jingjing Liu, Wei-Ying
591 Ma, Ya-Qin Zhang, Lin Yan, Mu Qiao, Yonghui Wu, and Mingxuan Wang. Dapo: An open-source
592 llm reinforcement learning system at scale, 2025. URL [https://arxiv.org/abs/2503.
593 14476](https://arxiv.org/abs/2503.14476).

594 Yang Yue, Zhiqi Chen, Rui Lu, Andrew Zhao, Zhaokai Wang, Yang Yue, Shiji Song, and Gao Huang.
595 Does reinforcement learning really incentivize reasoning capacity in llms beyond the base model?,
596 2025. URL <https://arxiv.org/abs/2504.13837>.
597

598 Lunjun Zhang, Arian Hosseini, Hritik Bansal, Mehran Kazemi, Aviral Kumar, and Rishabh Agarwal.
599 Generative verifiers: Reward modeling as next-token prediction, 2025. URL <https://arxiv.org/abs/2408.15240>.
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
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

A THE USE OF LARGE LANGUAGE MODELS

In preparing this manuscript, we used a large language model (LLM) solely for polishing the writing style and improving the clarity of the manuscript. The LLM was not used for generating research ideas, designing experiments, conducting analyses, or deriving results. All scientific contributions, including the conceptualization, methodology, experiments, and conclusions, were developed entirely by the authors.

B IMPLEMENTATION DETAILS

B.1 HYPERPARAMETERS

Table 2 and Table 3 show the detailed hyperparameters we use during training and evaluation.

B.2 PROMPT TEMPLATES

Table 5 shows the comparison between the commonly used Think-Answer prompt and our Generation-Verification Prompt. Our prompt guides the model to continuously alternate between generation and verification until the correct answer is obtained. Additionally, to enable the extraction of code generated by the model for providing accurate training rewards, we instruct the model to enclose the code within python blocks.

B.3 CASE STUDY OF ReVEAL

Table 4 shows the detail case of ReVeal.

B.4 TEMPLATES USED FOR TOOL FEEDBACK

Table 8 shows the mapping between execution results and hint templates: (1) for test cases that are verified as success, we give a [Passed] signal in the judgement area; (2) for test cases that are verified as failed, we give a [Failed] signal in the judgement area; (3) for test cases that are verified as wrong, we give a clear feedback of [Wrong test case] for individual failures, or [No correct test cases generated] if all test cases are invalid; (4) for format error, we will give the feedback of formatting instructions to guide correct generation.

C COMPARE WITH REACT-STYLE PROMPTING FRAMEWORK

Reflexion (Shinn et al., 2023) uses GPT-4 to generate code and tests, call tools for feedback, and refine code on easier benchmarks like HumanEval and MBPP, yet still suffers from inaccurate tests (causing 59% false negatives, 16% false positives) that can even hurt performance. Our work targets harder code competition problems and explicitly optimizes self-verification — a key factor for LongCoT success. As shown in Figure 4(d), test accuracy improves from 50% to 90% after ReVeal RL training.

We tested ReAct-style prompting on Qwen2.5-32B-Instruct with the same tool. Without RL-trained verification, it failed to support effective deep multi-turn refinement, highlighting the necessity of ReVeal training for stronger generation as well as enhanced verification, refinement capabilities to enable deeper scaling.

D THE EFFECTIVENESS OF SHORT-TERM MEMORY

To ensure context understanding in long interaction loops, we adopt short-term memory component, which implements a short-term memory mechanism that leverages in-context learning to store recent generation-verification loop within the contextual window. The system maintains a rolling history of the most recent interactions, including code generations, test verifications, and tool feedback. Critical information such as successful patterns, error types, and effective test structures are preserved in complete formats. This structured memory representation enables the model to quickly identify relevant patterns from recent history and build upon previous attempts, leading to faster convergence and reduced redundant exploration in the solution space.

Table 2: RL Training Hyperparameters for ReVeal-Qwen2.5-32B-Instruct and ReVeal-DAPO-Qwen-32B

Parameter	Qwen2.5-32B-Instruct	DAPO-Qwen-32B
Max Turn	3	3
Training Batch Size	128	1024
Mini-Batch Size	16	256
Learning Rate (Actor)	5×10^{-7}	1×10^{-6}
Learning Rate (Critic)	1×10^{-5}	1×10^{-5}
KL Coefficient	0.0	0.0
Maximum Prompt Length	4,096	4,096
Maximum Response Length	8,192	12,288
Maximum Tool Response Length	4,096	4,096
Temperature	1.0	1.0
Training Epochs	4	10

Table 3: RL Evaluation Hyperparameters for Qwen2.5-32B-Instruct and DAPO-Qwen-32B

Parameter	Qwen2.5-32B-Instruct	DAPO-Qwen-32B
Maximum Prompt Length	4,096	4,096
Maximum Response Length	20,480	51,200
Maximum Tool Response Length	4,096	4,096
Do Sample	True	True
Temperature for Pass@1	0.6	0.2/0.6
Top_p for Pass@1	0.9	0.9
Rollout for Pass@1	10	10
Temperature for Pass@k	0.6	0.6
Top_p for Pass@k	0.9	0.9
Rollout for Pass@k	128	128

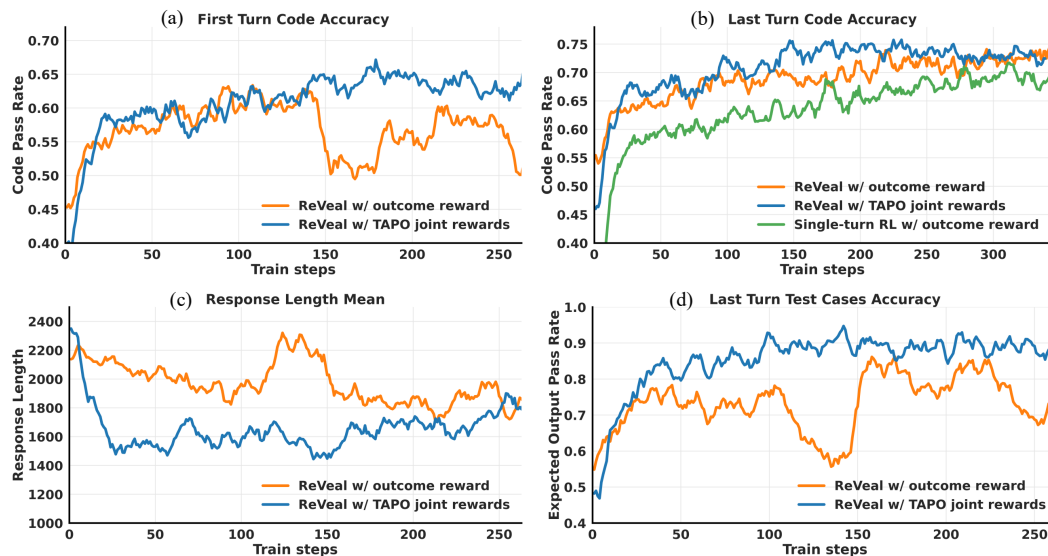


Figure 5: Comparison of code accuracy, test case accuracy, and response length across training for ReVeal (Qwen2.5-32B-Instruct) with turn-level rewards, ReVeal with outcome-only rewards, and single-turn RL without tool integration.

756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

Table 4: Case Study of ReVeal Performing Iterative Generation and Verification

Question:

The Two-dimensional kingdom has a regular army of n people... The Two-dimensional kingdom has m vests at its disposal, the j -th vest's size equals b^j . Help mobilize the Two-dimensional kingdom's army: equip with vests as many soldiers as possible.

Response:

<generation-think>

The problem can be solved using a dual circulation.

Here's a step-by-step approach to the solution:

...

Wait, I can use a greedy approach to solve this problem. Let me try the greedy algorithm by following steps:

...

</generation-think>

<generation-answer>

```
```python
```

```
def equip_soldiers(n, m, x, y, a, b):
```

```
 ...
```

**</generation-answer>**

**<verification-think>**

Now, we will verify the code with the provided test cases.

- Example Input 1 should result in matching 2 pairs: (Soldier at index 1 with Vest at index 1) and (Soldier at index 3 with Vest at index 2). There is no possible matching for the remaining soldiers with the available vests within the given range.

Wait, I could add an extra test case for verification: more vests than soldiers...

**</verification-think>**

**<verification-answer>**

- Input:

```
```1 3 5```
```

- Expected Output:

```
```3 2```
```

...

**</verification-answer>**

**<tool-feedback>**

...

- Actual Output:

```
```3 1```
```

- Judgement:

Failed

...

</tool-feedback>

<generation-think>

The initial code provided did not successfully handle the actual output as expected. There is a discrepancy where the vest index tracked in the loop is not correctly being reset or updated..

Alternatively, I should create a list that keeps track of which vests have already been used and adjust the code logic to ensure each vest is used only once. Here is a revised version of the code:

...

</verification-answer>

<tool-feedback>

...

- Actual Output:

```
```3 2```
```

- Judgement:

Passed

...

**</tool-feedback>**

---

Table 5: Comparison Between Think-Answer Prompt and ReVeal Prompt

Think-Answer Prompt	ReVeal Prompt
<p><b>system</b> You are Qwen, created by Alibaba Cloud. You are a helpful assistant.</p> <p><b>user</b> {question}</p> <p>First think about the reasoning process in the mind and then provides the user with the answer. The reasoning process and answer are enclosed within <code>&lt;think&gt;</code> <code>&lt;/think&gt;</code> and <code>&lt;answer&gt;</code> <code>&lt;/answer&gt;</code> tags, respectively, i.e., <code>&lt;think&gt;</code> reasoning process here <code>&lt;/think&gt;</code> <code>&lt;answer&gt;</code> answer here <code>&lt;/answer&gt;</code>. Enclose your code within delimiters as follows.</p> <pre>```python YOUR CODE HERE ```</pre> <p><b>assistant</b></p>	<p><b>system</b> You are Qwen, created by Alibaba Cloud. You are a helpful assistant.</p> <p><b>user</b> {question}</p> <p>First think through the reasoning process and write Python code to solve the problem, enclose your reasoning process in <code>&lt;generation-think&gt;</code> <code>&lt;/generation-think&gt;</code> and present the code in <code>```python</code> Your code <code>```</code></p> <p>within <code>&lt;generation-answer&gt;</code> <code>&lt;/generation-answer&gt;</code> tags. After that, verify your code by generating test cases:</p> <ol style="list-style-type: none"> <li>1. Extract sample test cases if the problem description includes them. When necessary, generate a small number of additional test cases to validate the correctness of the generated code.</li> <li>2. Enclose your reasoning process in <code>&lt;verification-think&gt;</code> <code>&lt;/verification-think&gt;</code> tags and enclose the final test cases and your verification conclusion within <code>&lt;verification-answer&gt;</code> <code>&lt;/verification-answer&gt;</code> tags and wrap each test case using the following format: - Input: <code>```</code> testcase input <code>```</code> - Expected Output: <code>```</code> expected testcase output <code>```</code></li> <li>3. Note that for "Use Call-Based format" questions, the testcase input should use a function call format, e.g., <code>fn_name(12, 12, 12)</code>.</li> </ol> <p><b>assistant</b></p>

Table 6: Performance comparison of ReVeal with baseline methods on LiveCodeBench. Pass@1 indicates the success rate;  $\Delta_{\uparrow}$  and  $\Delta_{\downarrow}$  represent the percentages of incorrect solutions corrected and correct solutions degraded after revision, respectively.

Model	LiveCodeBench V5		
	Pass@1	$\Delta_{\uparrow}$	$\Delta_{\downarrow}$
ReAct Prompting			
×1 turn	26.3	-	-
×3 turn	27.5	-	-
×6 turn	27.4	2.40	1.31
ReVeal (Qwen2.5-32B-Instruct)			
×1 turn	35.7	-	-
×3 turn	37.5	-	-
×6 turn	38.0	3.41	0.0

Table 7: Significance test of ReVeal (Qwen2.5-32B-Instruct) on LiveCodeBench V5. *mean ± std* indicates the average code pass@1 from 8 repeated experiments.

Model	LiveCodeBench V5
	<i>mean ± std</i>
ReVeal×6 turn w/ turn-level reward	38.02 ± 0.43
ReVeal×6 turn w/ outcome reward	37.09 ± 0.33



Table 8: Tool Feedback Templates for Different Execution Result Types.

Execution Results	Feedback
Success Test cases	- Input: {input}
	- Expected Output: {expected output}
	- Actual Output: {actual output}
	- Judgement Passed
Failed Test cases	- Input: {input}
	- Expected Output: {expected output}
	- Actual Output: {actual output}
	- Judgement Failed
- Failed Reason {failed reason}	
Wrong Test Cases	- Input: {input}
	- Expected Output: {expected output}
	- Actual Output: {actual output}
	- Judgement Wrong test case.
	No correct test cases are generated.
Error Format	No valid code because of the incorrect format. Write Python code again, and present the code in ```python Your code ```
	within <generation-answer> </generation-answer> tags. After that, verify your code by generating test cases:
	1. Extract sample test cases if the problem description includes them...
	2. Wrap each test case using the following format: - Input: ```
	testcase input ```
	- Expected Output: ```
	expected testcase output ```

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

Table 9: Performance comparison of ReVeal (Qwen2.5-32B-Instruct) with baseline methods on LiveCodeBench. Pass@1 indicates the success rate;  $\Delta_{\uparrow}$  and  $\Delta_{\downarrow}$  represent the percentages of incorrect solutions corrected and correct solutions degraded after revision, respectively.

Model	LiveCodeBench V5		
	Pass@1	$\Delta_{\uparrow}$	$\Delta_{\downarrow}$
<i>Existing Baselines</i>			
Qwen2.5-32B-Instruct	26.6	-	-
DAPO-Qwen2.5-32B	29.6	-	-
Qwen2.5-Coder-32B-Instruct	30.5	-	-
w/ critic×5 Qwen2.5-Coder	29.6	2.14	3.04
w/ critic×5 GPT-4o	32.9	4.82	2.50
w/ critic×5 CTRL	33.4	3.75	0.89
<i>RL based on Qwen2.5-32B-Instruct</i>			
Single-turn RL	33.9	-	-
ReVeal×6	<b>38.0</b>	<b>3.41</b>	<b>0.0</b>
<i>Ablation Study: TAPO with Joint Verifiable Rewards</i>			
ReVeal×6 w/ outcome reward	37.1	2.98	0.0
ReVeal×6 w/ TAPO with joint rewards	38.0	3.41	0.0

we compare models with and without memory integration. The baseline model without memory provides complete historical information from all previous turns directly to the model, maintaining full contextual details throughout the interaction sequence. We verified the impact of ReVeal on Pass@1 by LiveCodeBench with and without memory. Test results indicate that introducing short-term memory does not cause a decline in Pass@1 (wo/ memory 38.2% vs. w/ memory 38.3% at turn 15), and may even yield a slight performance boost. This sustained improvement capability highlights the memory mechanism’s effectiveness in enabling continuous learning and adaptation within computational constraints.