
REVES: REvision and VERification–Augmented Training for Test-Time Scaling

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Test-time scaling via sequential revision has emerged as a powerful paradigm for
2 enhancing Large Language Model (LLM) reasoning. However, standard post-
3 training methods primarily optimize single-shot objectives, creating a fundamental
4 misalignment with multi-step inference dynamics. While recent work treats this as
5 multi-turn reinforcement learning (RL), conventional approaches optimize over the
6 multi-step trajectories directly, failing to further exploit the high-quality mistakes
7 in intermediate steps that model can learn from correcting them. We propose a
8 two-stage iterative framework that alternates between online data/prompt augmen-
9 tation and policy optimization. By converting the intermediate steps (“near-miss”
10 answers) in the successful recovery trajectories into decoupled revision and ver-
11 ification prompts, our approach concentrates training on both effective answer
12 transformation and error identification. This approach enables efficient off-policy
13 data generation and reduces the computational overhead of long-horizon sampling
14 compared to standard multi-turn RL. On LiveCodeBench, using publicly avail-
15 able test cases as feedback, we observe gains of +6.5 points over the RL baseline
16 and +4.0 points over standard multi-turn training. Beyond coding, our approach
17 matches the previously reported SOTA result on circle packing. Math results under
18 ground-truth verification further confirm improved correction ability. It also gener-
19 alizes to out-of-distribution constraint-satisfaction puzzles such as `n_queens` and
20 `mini_sudoku`, where correctness is defined entirely by problem constraints.

21 1 Introduction

22 Large language models (LLMs) deployed in challenging settings rarely produce a correct answer
23 on the first attempt; real-world workflows routinely loop back with feedback. A growing family of
24 test-time scaling (TTS) algorithms formalizes this pattern, including sequential revision (Madaan
25 et al., 2023; Shinn et al., 2023), tree search (Inoue et al., 2025), evolutionary refinement (Lee et al.,
26 2025b), and best-of- N with a verifier (Chow et al., 2025). They all share a single primitive: given a
27 problem, a previous attempt, and feedback, produce an improved response.

28 This raises a basic question: *can we design a post-training algorithm that explicitly improves a*
29 *model’s test-time sequential revision capability?* We focus on sequential revision (SR) because gains
30 on SR transfer to the rest of the family: every revision-using TTS algorithm eventually calls the
31 policy on revise-from-prior-attempt inputs, and SR has the broadest such distribution. Theorem 3.1
32 shows that increasing the policy’s one-step revision value on SR-induced inputs also raises J_ϕ for
33 every revision-using TTS algorithm whose call inputs are covered by SR, up to small drift terms.

34 How, then, do we post-train an LLM to optimize the SR objective $J_{\phi_{\text{SR}}}$? Standard post-training,
35 RLHF (Ouyang et al., 2022), RLVR, and GRPO (DeepSeek-AI et al., 2025), optimizes a single-
36 shot expected reward, which is fundamentally misaligned with the multi-step nature of test-time

37 deployment. The natural fix is multi-turn RL on revision rollouts, but trajectory-level credit broadcasts
38 assign path-dependent credit. Take a wrong, wrong, correct rollout: every per-turn log-probability
39 gradient receives the same positive credit, including the wrong intermediates. Per rollout this credit is
40 biased, and only cancels in expectation across many rollouts.

41 We propose REVES, a two-stage iterative framework that targets the SR objective $J_{\phi_{\text{SR}}}$ directly
42 through its per-state structure. We show that $J_{\phi_{\text{SR}}}$ decomposes exactly into a weighted sum of
43 per-state one-step recovery probabilities along visited rollouts (Lemma 4.1), exposing single-state,
44 single-turn gradients with no horizon credit assignment. REVES is an offline realization of this
45 per-state training signal: each epoch, it runs SR rollouts under the current policy, retains those that
46 succeed within budget, converts the intermediate (“near-miss”) answers into decoupled *revision* and
47 *verification* prompts, and trains with standard single-turn RL on the augmented prompt set; the next
48 epoch refreshes the augmentation. This concentrates training on per-state recovery and avoids both
49 the bias of trajectory broadcasts and the cost of long-horizon online sampling.

50 Empirically, on LiveCodeBench with public test cases as feedback, REVES achieves +6.5 points
51 over the single-shot RL baseline and +4.0 points over standard multi-turn training. On the circle
52 packing benchmark, a Qwen3-4B base trained with REVES matches the previously reported best
53 results from much larger evolutionary search systems built on Gemini-2.0 Pro/Flash and Qwen3-8B.
54 Math benchmarks (MATH500, AIME24/25) show consistent improvements under both ground-truth
55 and self-confidence stopping, and the trained policies generalize to out-of-distribution constraint-
56 satisfaction puzzles (n_queens, mini_sudoku). Beyond standalone SR, REVES improves every
57 revision-using TTS algorithm we tested (BoN, MCTS, AB-MCTS variants, Mind Evolution), in line
58 with the conditional bound in Theorem 3.1.

59 Our contributions are threefold:

- 60 1. **Conceptual.** We formulate test-time scaling as a meta-RL problem and show that improving
61 sequential revision also improves any revision-using TTS algorithm whose call inputs are covered
62 by SR (Theorem 3.1).
- 63 2. **Methodological.** We derive an exact decomposition of $J_{\phi_{\text{SR}}}$ into per-state one-step recovery
64 probabilities and propose REVES (Lemma 4.1), a two-stage framework that converts intermediate
65 near-miss answers into decoupled revision and verification prompts and trains with single-turn
66 RL, sidestepping the path-dependent credit assignment of multi-turn RL.
- 67 3. **Empirical.** REVES delivers consistent test-time improvements across math, coding, out-of-
68 distribution puzzles, circle packing, and improves revision-using TTS algorithm we evaluated.

69 2 Related Work

70 **Test-time-aware post-training.** A growing line of work shows the importance of aligning the
71 post-training objective with the intended test-time inference strategy. For parallel-sampling families,
72 this includes optimizing pass@ k (Chen et al., 2025; Walder and Karkhanis, 2025; Tang et al., 2025)
73 when deployment uses best-of- N , majority voting (Wang et al., 2023), or learned verifiers (Chow
74 et al., 2025). Our deployment is structurally different: the test-time strategies we care about, including
75 sequential revision, MCTS variants, and Mind Evolution, are all *revision-using* algorithms whose
76 calls take as input a previous attempt together with its feedback. Aligning training with this inference
77 family therefore requires a different objective.

78 **Training sequential revision capability.** Closer to our setting, several prior methods explicitly
79 train the policy to revise, including supervised and preference-based approaches (Qu et al., 2024;
80 Xiong et al., 2025), RL-based approaches (Kumar et al., 2024; Jiang et al., 2025; Lee et al., 2025a),
81 and multi-turn RL methods that incorporate per-turn natural-language or numerical critiques (Jain
82 et al., 2025; Li et al., 2025b,a; Zhang et al., 2026). Our work shares the broader goal of strengthening
83 the model’s revision capability, but differs in framing: we primarily focus on how to align training
84 with the test-time objective of a revision-using inference strategy.

85 **Weakness-driven data synthesis and guided exploration.** SwS (Liang et al., 2025) synthesizes
86 new problems targeting the model’s weaknesses using strong external teacher models, and POPE (Qu
87 et al., 2026) uses privileged oracle hints to guide on-policy exploration on hard problems. These meth-

88 ods are orthogonal to ours: they target single-shot pass@1, whereas REVES targets a fundamentally
 89 different test-time objective.

90 3 Test-Time Scaling as a Meta-RL Problem

91 **Setup and two motivating examples.** We fix notation used throughout the paper. Denote by \mathcal{X}
 92 the distribution over input problems, \mathcal{Y} the answer space, and $r^* : \mathcal{X} \times \mathcal{Y} \rightarrow [0, 1]$ the *ground-truth*
 93 reward, where $r^*(x, y) = 1$ indicating a fully correct answer. We focus on the verifiable-reward
 94 regime where r^* is computable during both training and at test time, e.g., public test-case execution
 95 for code, and constraint satisfaction for puzzles. When r^* is unavailable at test time, a surrogate
 96 stopping rule replaces it (Section 4). We model the LLM as a policy π_θ in RL parameterized by
 97 θ . At deployment, an LLM rarely operates by a single call to π_θ ; instead a *test-time scaling (TTS)*
 98 *algorithm* ϕ orchestrates up to K policy calls and extracts a final answer $\hat{y} \in \mathcal{Y}$. Two canonical
 99 examples bracket the design space: *Best-of-N (BoN)* (Chow et al., 2025) draws N candidates in
 100 parallel from $\pi_\theta(\cdot | x)$ and selects \hat{y} via an external scorer; every call uses the same original prompt.
 101 *Sequential Revision (SR)* (Madaan et al., 2023) generates an initial response $y_1 \sim \pi_\theta(\cdot | x)$ from
 102 prompt x , and for $t = 2, \dots, K$ produces a revision $y_t \sim \pi_\theta(\cdot | x, y_{t-1}, f_{t-1})$ conditioned on the
 103 previous attempt y_{t-1} and a feedback signal f_{t-1} (e.g., compiler error messages in code generation),
 104 halting at the first correct response. Detailed pseudocode is in Protocol 1 of Appendix A.

105 **General formulation: the TTS-induced decision process.** Pairing any TTS algorithm ϕ with
 106 π_θ induces a decision process with prompt-state space \mathcal{S} and response space \mathcal{Y} . At each step t , the
 107 policy emits a response $y_t \sim \pi_\theta(\cdot | s_t)$, and ϕ constructs the next prompt via a transition kernel
 108 $P_\phi(s_{t+1} | s_{1:t}, y_{1:t})$. The initial state is $s_1 = x$. Note that the policy itself remains Markov in the
 109 current state (at every step, π_θ conditions only on the local prompt s_t); the history dependence is in
 110 the kernel, because ϕ may branch from any earlier response. We focus on the family $\Phi_{\mathcal{R}}$ of TTS
 111 algorithms that issue at least one revision call on a prior attempt; this includes Sequential Revision
 112 and adaptive tree search such as AB-MCTS (Inoue et al., 2025). For such algorithms a revision-step
 113 prompt has the form $s_{t+1} = (x, y_{i_t}, f_{i_t})$, $i_t \in \{1, \dots, t\}$, where the prior step i_t is chosen by ϕ :
 114 for SR, $i_t = t$ (always the most recent response), so s_{t+1} depends only on (s_t, y_t) and the induced
 115 process is a standard MDP; for tree-search algorithms such as AB-MCTS, i_t may be any earlier
 116 response that ϕ 's search policy decides to branch from, and the kernel is genuinely history-dependent.
 117 The procedure runs until a random stopping time $\tau := K \wedge \min\{t : r^*(x, y_t) = 1\}$ and outputs
 118 $\hat{y} = y_\tau$; that is, the procedure either *early-stops* once a verified-correct response is produced, or
 119 *exhausts the budget* K . Running π_θ inside ϕ yields the TTS-induced objective

$$J_\phi(\theta) = \mathbb{E}_{x \sim \mathcal{X}, (s_{1:\tau}, y_{1:\tau}) \sim (\pi_\theta, \phi)} [r^*(x, y_\tau)],$$

120 to be contrasted with the standard single-shot objective

$$J_{\text{OneShot}}(\theta) = \mathbb{E}_{x \sim \mathcal{X}, y \sim \pi_\theta(\cdot | x)} [r^*(x, y)],$$

121 which is what RLHF (Ouyang et al., 2022), RLVR, and GRPO (DeepSeek-AI et al., 2025) optimize.
 122 Optimizing J_ϕ rather than J_{OneShot} is a *meta-RL* problem (Figure 1): the policy must adapt across
 123 the distribution of multi-step contexts that ϕ generates, not just answer the original question. Our
 124 goal is to improve J_ϕ for every TTS algorithm $\phi \in \Phi_{\mathcal{R}}$.

125 This raises an important question: *does there exist a training target $J_{\phi_{\text{train}}}$ whose optimization*
 126 *provably improves J_ϕ for every $\phi \in \Phi_{\mathcal{R}}$?* If so, we can train against one algorithm and have the
 127 gains transfer to every other deployment-time choice in $\Phi_{\mathcal{R}}$, side-stepping the cost of running tree
 128 search or external verifiers during training. The answer is affirmative. The key fact is that every
 129 algorithm in $\Phi_{\mathcal{R}}$ at some point invokes π_θ to revise a previous attempt, so a training target whose
 130 own revision-input distribution covers the deployment one improves every member of $\Phi_{\mathcal{R}}$ when its
 131 revision-value goes up.

132 **Theorem 3.1** (Sequential-revision recovery transfers to revision-using TTS). *Let π_0 be a baseline*
 133 *policy and π_1 an updated policy. For $\phi \in \Phi_{\mathcal{R}}$, policy π , and revision-call input z , let*

$$V_\pi(z) := \mathbb{P}_{y \sim \pi(\cdot | z)}(r^*(x, y) = 1), \quad \rho_\pi^\phi(z) := \mathbb{E}_{(\phi, \pi)} \left[\sum_{t=1}^{\tau} \mathbf{1}\{z_t = z\} \right]$$

134 *denote the one-step recovery probability and the expected visit count to z . Assume:*

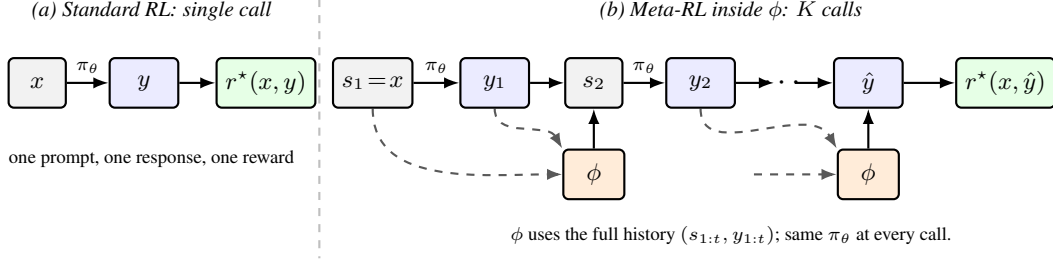


Figure 1: **(a)** Standard RL: policy π_θ maps problem x to response y with reward $r^*(x, y)$. **(b)** A TTS algorithm ϕ invokes the same π_θ over up to K steps, each prompt s_t built from the prior history (dashed arrows), and extracts a final response $\hat{y} = y_\tau$ at stopping time $\tau \leq K$; π_θ must therefore behave well across the whole context distribution induced by ϕ , not just on x .

135 (C1) Coverage by SR at π_0 : *there exists $C_\phi \geq 1$ such that $\rho_{\pi_0}^\phi(z) \geq \frac{1}{C_\phi} \rho_{\pi_0}^{\phi_{\text{SR}}}(z)$ for every z .*

136 (C2) Monotone recovery improvement: $V_{\pi_1}(z) \geq V_{\pi_0}(z)$ for every $z \in \text{supp}(\rho_{\pi_0}^\phi)$.

137 Then, we have the following guarantee:

$$J_\phi(\pi_1) - J_\phi(\pi_0) \geq \frac{1}{C_\phi} [J_{\phi_{\text{SR}}}(\pi_1) - J_{\phi_{\text{SR}}}(\pi_0)] - \|\rho_{\pi_1}^\phi - \rho_{\pi_0}^\phi\|_1 - \frac{1}{C_\phi} \|\rho_{\pi_1}^{\phi_{\text{SR}}} - \rho_{\pi_0}^{\phi_{\text{SR}}}\|_1.$$

138 The take-away is clear: improving the SR test-time objective $J_{\phi_{\text{SR}}}$ also improves J_ϕ for every
 139 revision-using $\phi \in \Phi_{\mathcal{R}}$, so we focus on SR as the training target. See Appendix C.2 for the proof.

140 **Our training target $J_{\phi_{\text{SR}}}$.** The theorem licenses any $\phi_{\text{train}} \in \Phi_{\mathcal{R}}$ whose visit measure dominates
 141 the deployment-time ϕ 's. We pick ϕ_{SR} and optimize

$$J_{\phi_{\text{SR}}}(\theta) := \mathbb{E}_{x \sim \mathcal{X}, y_{1:\tau} \sim (\pi_\theta, \phi_{\text{SR}})} [r^*(x, y_\tau)]$$

142 for two reasons. First, ϕ_{SR} has the broadest revision-input distribution in $\Phi_{\mathcal{R}}$. At every step, the
 143 next prompt is conditioned on a freshly drawn y_{t-1} , so successive revision inputs z_t are spread
 144 out. Tree-search and selection algorithms instead branch repeatedly from a small pool of cached
 145 responses, concentrating their visits on a narrow subset. (C1) therefore holds at the smallest C_ϕ
 146 over deployment-time ϕ . Second, ϕ_{SR} is itself a strong TTS algorithm. Figure 2 shows that on
 147 TravelPlanner (Xie et al., 2024), ϕ_{SR} matches or surpasses Mind Evolution (Lee et al., 2025b) and
 148 AB-MCTS (Inoue et al., 2025) (left), with the Markov-revision form $\pi_\theta(\cdot | x, y_{\text{prev}}, f)$ on the best
 149 cost-accuracy frontier (right), while being the simplest to implement.

150 **Optimizing $J_{\phi_{\text{SR}}}$ is not optimizing
 151 single-shot pass@1.** A natural next
 152 question is whether standard single-shot
 153 RLVR optimization already addresses
 154 $J_{\phi_{\text{SR}}}$. It does not: two policies can be
 155 indistinguishable under J_{OneShot} on every
 156 problem in \mathcal{X} yet differ substantially
 157 under $J_{\phi_{\text{SR}}}$.

158 **Theorem 3.2** (Objective mismatch).
 159 For any revision length $K \geq 2$, there
 160 exist policies π_1, π_2 , $\Delta(K) > 0$, and
 161 a problem distribution \mathcal{X} such that
 162 $J_{\text{OneShot}}(\pi_1) = J_{\text{OneShot}}(\pi_2)$, while $J_{\phi_{\text{SR}}}(\pi_2) - J_{\phi_{\text{SR}}}(\pi_1) \geq \Delta(K)$.

163 The implication is direct: J_{OneShot} and $J_{\phi_{\text{SR}}}$ are different objectives, and policies with identical
 164 single-shot performance can have substantially different sequential-revision performance. Our
 165 problem reduces to a single concrete question: how do we post-train π_θ so that $J_{\phi_{\text{SR}}}$ goes up? The
 166 next section presents our algorithm.

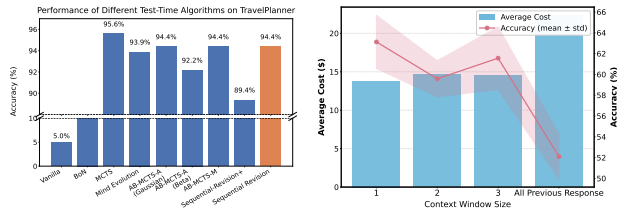


Figure 2: DeepSeek-V3 on TravelPlanner. **Left:** sequential revision vs. Mind Evolution / MCTS / AB-MCTS. **Right:** cost-accuracy as the revision context window covers the k most recent responses, $k \in \{1, 2, 3\}$, or the full history.

167 **4 REVES: From Path-Level Credit to Recovery-Level Supervision**

168 Section 3 established that improving $J_{\phi_{\text{SR}}}$ is our focus. This section turns that target into an algorithm.
 169 We start by writing $J_{\phi_{\text{SR}}}$ exactly as a sum of per-state one-step recovery probabilities along visited
 170 trajectories; this hazard decomposition makes the natural training signal explicit and reads off the
 171 algorithm directly from the figure.

172 **Lemma 4.1.** *Let $\zeta = (z_1, y_1, \dots, z_\tau, y_\tau) \sim (\pi_\theta, \phi_{\text{SR}})$ be a stopped SR trajectory under π_θ : at
 173 each step t the model samples $y_t \sim \pi_\theta(\cdot | z_t)$ and stops at the first verified correct response or
 174 at budget K , with $z_t = (x, y_{t-1}, f_{t-1})$ the revision state before generating y_t . Let $V_\pi(z_t) :=$
 175 $\mathbb{P}_{y' \sim \pi_\theta(\cdot | z_t)}(r^*(x, y') = 1)$ denote the one-step recovery probability at z_t , where y' is a fresh sample
 176 (not the trajectory’s own y_t), and let $\rho_\theta(z) := \mathbb{E}_{x \sim \mathcal{X}, \zeta \sim (\pi_\theta, \phi_{\text{SR}})} \left[\sum_{t=1}^{\tau} \mathbf{1}\{z_t = z\} \right]$ denote the
 177 expected number of times state z is visited along a SR rollout. Then*

$$\begin{aligned}
 J_{\phi_{\text{SR}}}(\theta) &= \sum_{t=1}^K \mathbb{E}_{x, \zeta} [\mathbf{1}\{\tau \geq t\} V_\pi(z_t)] = \mathbb{E}_{x \sim \mathcal{X}, \zeta \sim (\pi_\theta, \phi_{\text{SR}})} \left[\sum_{t=1}^{\tau} V_\pi(z_t) \right] \\
 &= \sum_z \rho_\theta(z) \mathbb{E}_{y' \sim \pi_\theta(\cdot | z)} [r^*(x, y')]. \tag{1}
 \end{aligned}$$

178 See Appendix C.3 for the proof. Equation (1) reads the algorithm off Figure 3 (b). The horizontal
 179 SR trajectory provides the outer samples $z \sim \rho_\theta$, and at each visited state $z = (x, y_{t-1}, f_{t-1})$ the
 180 branches drawn from $\pi_\theta(\cdot | z)$ are online rollouts that estimate and optimize the per-state factor
 181 $V_\pi(z) = \mathbb{E}_{y' \sim \pi_\theta(\cdot | z)} [r^*(x, y')]$. Since $J_{\phi_{\text{SR}}}$ is a nonnegative-weighted sum of these factors, raising
 182 V_π at any visited z directly raises $J_{\phi_{\text{SR}}}$, without horizon credit assignment.

183 **Comparison with multi-turn RL.** Multi-turn RL provides an unbiased policy-gradient estimator
 184 for the trajectory-level objective. However, its per-turn learning signal is highly indirect. Consider
 185 a wrong–wrong–correct rollout: the gradients for all three turns are weighted by the same positive
 186 terminal advantage, including the two turns that produced incorrect intermediate responses. This is
 187 valid for optimizing terminal success of the whole trajectory, but it yields a high-variance and coarse
 188 signal for improving the local revision behavior at each intermediate state.

189 Our method instead reuses such incorrect intermediate responses as revision states. Given a state
 190 $z = (x, y_{\text{prev}}, f_{\text{prev}})$, we draw fresh samples from $\pi(\cdot | z)$ and estimate the local one-step recovery
 191 probability $V_\pi(z) = \Pr_{y \sim \pi(\cdot | z)} [r^*(x, y) = 1]$. Thus, the intermediate wrong response is not treated
 192 as an action to imitate; it is treated as a state from which the model should learn to recover. This
 193 converts a successful long-horizon trajectory into several short-horizon revision problems, providing
 194 a sharper training signal for local recovery while relying on the current policy to discover useful
 195 intermediate states.

196 **From recovery target to algorithm.** An ideal on-policy implementation of Lemma 4.1 draws
 197 $z \sim \rho_\theta$ by running SR and runs single-turn RL at each z . Three design choices turn this into the
 198 algorithm we use.

199 (1) *Verification prompts alongside revision data.* At test time the verifier is typically unavailable. We
 200 therefore include both *revision prompts* and *verification prompts* (asking the policy to judge whether
 201 an intermediate response is correct) in the augmented data, so the trained policy can self-stop at
 202 deployment. The split of gains across the two prompt types is in Table 4.

203 (2) *Offline generation.* An on-policy estimator of $\nabla J_{\phi_{\text{SR}}}$ would need fresh SR rollouts every gradient
 204 step, which is intrinsically serial and prohibitively slow. We instead generate visited states once per
 205 epoch in Stage I, reuse them through Stage II, and refresh at the next epoch.

206 (3) *Successful-trajectory filtering.* For a far-from-optimal base policy, online samples at states from
 207 failed rollouts rarely produce a verifier-accepted recovery and waste tokens. We restrict Stage I
 208 to states on trajectories that succeed within budget K , where at least one recovery target exists by
 209 construction; the filter loosens as π_θ improves.

210 **Two-stage realization. Stage I (Data augmentation).** For each input x , run SR under the current
 211 π_θ : $y_1 \sim \pi_\theta(\cdot | x)$ and $y_t \sim \pi_\theta(\cdot | x, y_{t-1}, \text{Feedback}(y_{t-1}))$ for $t = 2, \dots, K$, terminating at

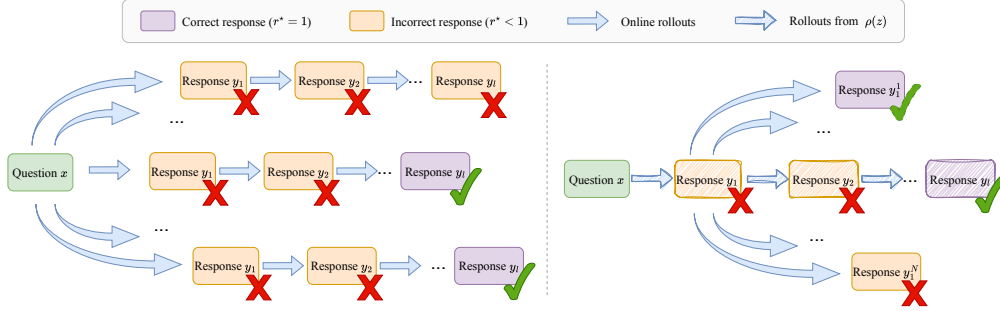


Figure 3: Mapping Eq. (1) onto the figure. **(a)** Stage I samples SR rollouts under the current π_θ and keeps those that succeed within budget K . **(b)** The horizontal hatched chain is a retained sequential-revision trajectory; it determines the visit-count weight $\rho_\theta(z)$ at each intermediate state $z = (x, y_{t-1}, f_{t-1})$. The branches at each visited z are fresh samples from $\pi_\theta(\cdot | z)$; their verifier outcomes provide the per-state estimate of $V_\pi(z) = \mathbb{E}_{y' \sim \pi_\theta(\cdot | z)}[r^*(x, y')]$.

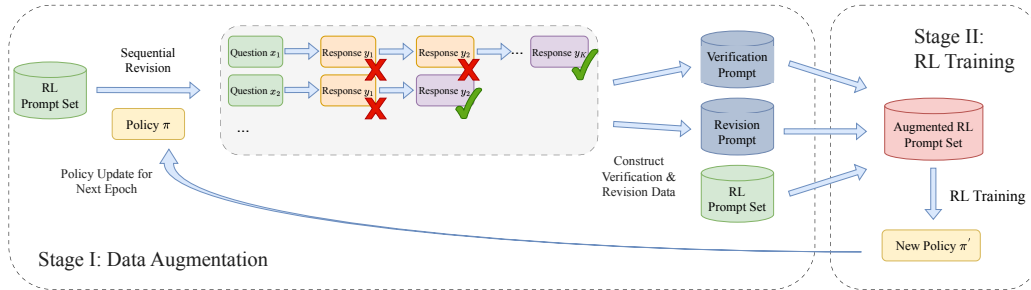


Figure 4: Overview of the proposed RL-based data augmentation framework. In Stage I, the current policy performs sequential revision; only trajectories that succeed within budget K are retained, and their intermediate states yield revision and verification prompts that are merged with the original RL prompt set. In Stage II, the policy is updated by single-turn RL on the augmented set, and the new policy regenerates the augmentation in the next epoch.

212 the first correct response or after K attempts. Discard trajectories that fail or are already correct
 213 at y_1 . From each retained trajectory (x, y_1, \dots, y_ℓ) with $r^*(x, y_\ell) = 1$, build (i) *revision data*
 214 $\{\text{Revision Prompt}(x, y_i)\}_{i=1}^\ell$, where each prompt asks the model to revise the previous response,
 215 and (ii) *verification data* $\{\text{Verification Prompt}(x, y_i)\}_{i=1}^\ell$, where each prompt asks the model
 216 to judge whether y_i is correct. Templates are in Appendix A. **Stage II (Single-turn RL)**. Train π_θ
 217 on the union of the original RL prompts and the prompts produced in Stage I. The updated policy serves
 218 as the rollout policy for the next epoch's Stage I. Figure 4 shows the loop.

219 **Efficiency.** Our algorithm is also more efficient than standard multi-turn training. First, augmented
 220 samples are generated in an *off-policy* manner. With a well-designed infrastructure, data augmentation
 221 can run in parallel with RL training, so the model continues training on the original RL prompts while
 222 new augmented data are generated asynchronously, avoiding any interruption to the main training
 223 loop. Second, training itself is single-turn. Each gradient step processes one (prompt, response) pair
 224 as in ordinary RLVR, with no per-step sequential rollout in the inner loop, which is the dominant cost
 225 of multi-turn RL with budget K . Appendix F.2 reports the wall-clock comparison.

226 **Test-time stopping rule.** At test time the oracle stopping time τ is almost always unavailable, so
 227 we adopt a surrogate stopping time $\hat{\tau}$ that replaces r^* with a domain-specific signal.

228 For mathematical reasoning tasks, we leverage the *Tail Confidence* metric of Fu et al. (2025) as
 229 the surrogate signal. At each revision step t , the token-level confidence at position i of y_t is
 230 $C_i = -\frac{1}{k} \sum_{j=1}^k \log P_i(j)$, where $P_i(j)$ is the probability of the j -th most likely token (a larger C_i
 231 indicates a more concentrated, more confident distribution). The response-level score c_t averages C_i
 232 over the tail positions $i \in T_{\text{tail}}$ of y_t . The stopping time is $\hat{\tau} := K \wedge \min\{t \geq 3 : c_t / \sum_{j=1}^t c_j > c\}$

233 for a predefined threshold c . If $\hat{\tau} = K$ (budget exhausted without trigger), we output $\hat{y} = y_{t^*}$ with
 234 $t^* = \arg \max_{1 \leq t \leq K} c_t$, breaking ties uniformly at random when several responses share the same
 235 confidence. Selecting the most confident response across all K revisions instead of just y_K improves
 236 robustness under small budgets by mitigating noise in confidence estimates at the stopping step.

237 For coding tasks, we use publicly available test cases as a surrogate reward $\tilde{r} : \mathcal{X} \times \mathcal{Y} \rightarrow \{0, 1\}$, with
 238 $\tilde{r}(x, y) = \mathbf{1}(y \text{ passes all public test cases})$. The stopping time is $\hat{\tau} := K \wedge \min\{t : \tilde{r}(x, y_t) = 1\}$,
 239 and we output $\hat{y} = y_K$ when $\hat{\tau} = K$.

240 5 Experimental Results

241 **Models and data.** We train Qwen2.5-7B, Qwen2.5-3B (Yang et al., 2024), Qwen3-4B(non-
 242 thinking) (Yang et al., 2025) and DeepSeek-R1-Distill-Qwen-7B (DeepSeek-AI et al., 2025) directly
 243 with RL. Training data is sampled from Skywork/Skywork-0R1-RL-Data (7298 math + 1015 cod-
 244 ing = 8313 examples). We evaluate on MATH500 and AIME 24/25 for math; LiveCodeBench
 245 (Jain et al., 2024) and CodeContest (Li et al., 2022) for coding; and n_queens / mini_sudoku from
 246 ReasoningGym (Stojanovski et al., 2025) for out-of-distribution puzzles.

247 **Baselines and implementation.** We compare against **RL** (single-response RL, no revision), **Multi-**
 248 **Turn** (multi-turn rollouts without explicit self-verification), and **PAG** (Jiang et al., 2025) (multi-turn
 249 with self-verification via a designed reward). Full details in Appendix E.

250 5.1 Main Results

251 **Consistent improvements on coding tasks.** Table 1 demonstrates the benefits of RL-Augmentation
 252 on coding benchmarks. Under test-time surrogate stopping criteria, where publicly available test cases
 253 are used to determine when to stop and execution results are provided as feedback for subsequent
 254 revisions (see Appendix A for an example), RL-Augmentation consistently improves performance
 255 across different test-time budgets, demonstrating robustness beyond symbolic rewards. Additional
 256 results on DeepSeek-R1-Distill-7B (DeepSeek-AI et al., 2025) are provided in Appendix G.

Table 1: Sequential-revision test-time scaling performance on coding benchmarks (LiveCodeBench, CodeContest). The best result in each column is shown in bold, and the second-best is underlined. TestCases- B is sequential-revision (SR) test-time scaling with budget B , using execution-based test cases as the verifier and feeding the resulting feedback into the next revision step. Best results in each column are in bold; second-best are underlined.

Model	LiveCodeBench (08/01/24–01/01/25)			LiveCodeBench (01/01/25–05/01/25)		CodeContest		
	OneShot	TestCases-32	TestCases-4	OneShot	TestCases-32	OneShot	TestCases-32	TestCases-4
Qwen-2.5-7B	4.78	6.13	6.13	6.04	18.86	0.41	5.73	5.86
Qwen-2.5-7B-RL	18.02	23.04	20.83	19.96	24.18	7.07	11.09	7.47
Qwen-2.5-7B-Multi-turn	<u>19.36</u>	25.49	21.33	20.70	<u>27.47</u>	<u>6.87</u>	14.34	8.08
Qwen-2.5-7B-PAG	19.49	<u>25.73</u>	<u>23.16</u>	19.78	26.37	6.46	<u>14.95</u>	<u>9.09</u>
Qwen-2.5-7B-RL-Aug	18.87	29.54	23.77	<u>20.51</u>	30.03	7.07	16.57	11.72
Qwen-2.5-3B	5.88	15.81	7.23	7.87	18.31	1.62	6.26	2.43
Qwen-2.5-3B-RL	9.19	12.01	<u>11.52</u>	14.11	17.76	3.03	4.24	3.64
Qwen-2.5-3B-Multi-turn	9.93	<u>15.69</u>	12.13	13.19	<u>19.23</u>	<u>3.23</u>	<u>7.88</u>	4.44
Qwen-2.5-3B-PAG	8.09	14.71	11.15	<u>13.37</u>	18.50	3.64	5.66	<u>4.85</u>
Qwen-2.5-3B-RL-Aug	<u>9.56</u>	17.65	11.03	14.29	21.61	2.22	8.48	5.45
Qwen3-4B (non-thinking)	17.60	22.10	20.20	20.90	22.00	8.30	12.80	7.70
Qwen3-4B-RL	<u>22.40</u>	30.90	<u>29.00</u>	26.90	29.70	21.80	28.20	26.30
Qwen3-4B-Multi-turn	20.20	30.10	28.30	24.70	<u>31.90</u>	11.50	39.10	<u>28.80</u>
Qwen3-4B-PAG	<u>22.40</u>	27.90	24.30	23.10	28.60	10.90	32.10	28.20
Qwen3-4B-RL-Aug	23.90	32.40	30.90	<u>25.80</u>	33.50	<u>14.10</u>	<u>37.20</u>	32.70

257 **Improved sequential revision capability for mathematical reasoning.** Table 2 shows that RL-
 258 Augmentation substantially improves sequential revision performance across all evaluated math
 259 benchmarks. Under oracle stopping (the Oracle-32 and Oracle-4 columns, where the rollout halts
 260 as soon as a generated answer matches the ground truth, isolating the model’s ability to revise
 261 an incorrect solution into a correct one), our method achieves consistently large gains. These
 262 improvements persist under the SelfConf-4 column, the practical test-time stopping rule based on the
 263 model’s own confidence, indicating that the learned capability transfers beyond the oracle setting.

Method	Model	$n = 26$ (\uparrow)
AlphaEvolve (Novikov et al., 2025)	Gemini-2.0 Pro + Flash	2.635862
AlphaEvolve V2 (Georgiev et al., 2025)	Gemini-2.0 Pro + Flash	2.635983
ShinkaEvolve (Lange et al., 2025)	Ensemble	2.635982
ThetaEvolve (Wang et al., 2025)	R1-Qwen3-8B	2.635983
TTT-Discover (Yuksekgonul et al., 2026)	Qwen3-8B	2.635983
Ours (REVES)	Qwen3-4B	2.635983

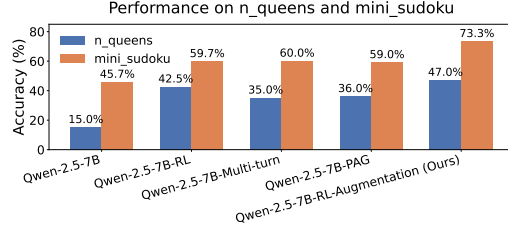


Figure 5: REVES performs well beyond standard reasoning benchmarks. **Left:** circle packing benchmark (sum of radii at $n = 26$, \uparrow); REVES on Qwen3-4B matches the best previously reported number despite using a smaller base model and less rollouts. **Right:** out-of-distribution puzzle benchmarks (n_queens and mini_sudoku) under verifier-based stopping.

264 **Matching state-of-the-art on circle packing.** We further evaluate REVES on the circle packing
 265 benchmark, where the goal is to pack n unit-area circles into a unit square to maximize the sum of
 266 their radii. We directly optimize Eq. (1) on this task. With a Qwen3-4B, our trained policy matches
 267 the previously reported best results from much larger evolutionary search systems built on Gemini-2.0
 268 Pro/Flash and Qwen3-8B (Figure 5). Full implementation details are in Appendix H.

Table 2: Sequential-revision test-time scaling performance on math benchmarks (AIME24/25, MATH500). SelfConf- B and Oracle- B are sequential-revision (SR) test-time scaling with budget B , using the model’s own confidence as a surrogate stopping rule and the ground-truth answer as the oracle stopping rule, respectively. \dagger denotes evaluation using the authors’ original codebase and protocol. Best and second-best results are shown in bold and underline.

Model	AIME 24				AIME 25				MATH 500			
	OneShot	Oracle-32	Oracle-4	SelfConf-4	OneShot	Oracle-32	Oracle-4	SelfConf-4	OneShot	Oracle-32	Oracle-4	SelfConf-4
Qwen-2.5-7B	2.03	13.96	4.79	2.34	0.26	6.25	1.25	0.57	23.83	50.05	28.98	26.10
Qwen-2.5-7B-RL	17.86	<u>33.49</u>	<u>21.41</u>	<u>18.65</u>	10.35	22.86	<u>13.44</u>	<u>12.45</u>	<u>76.13</u>	85.90	79.43	78.60
Qwen-2.5-7B-Multi-turn	14.58	30.26	20.00	13.80 \dagger	7.19	25.31	11.77	9.84 \dagger	75.93	<u>87.13</u>	80.13	<u>77.95</u> \dagger
Qwen-2.5-7B-PAG	<u>15.31</u>	31.09	19.58	12.97 \dagger	7.45	<u>26.82</u>	12.81	9.43 \dagger	<u>76.13</u>	86.68	<u>80.75</u>	77.20 \dagger
Qwen-2.5-7B-RL-Aug	14.90	45.73	25.94	19.58	<u>9.27</u>	40.52	22.34	15.10	76.75	94.68	85.60	76.20
Qwen-2.5-3B	1.15	12.92	3.44	1.93	0.47	7.92	1.30	0.57	30.53	66.98	39.33	38.40
Qwen-2.5-3B-RL	<u>6.41</u>	16.93	8.13	7.40	1.82	9.01	3.91	2.40	<u>64.93</u>	80.53	69.57	65.63
Qwen-2.5-3B-Multi-turn	7.60	<u>20.52</u>	<u>11.98</u>	<u>7.97</u> \dagger	2.60	<u>12.71</u>	<u>5.16</u>	2.29 \dagger	65.28	<u>81.23</u>	<u>72.65</u>	67.00 \dagger
Qwen-2.5-3B-PAG	6.30	14.79	7.92	5.31 \dagger	3.44	11.41	4.69	<u>4.11</u> \dagger	64.00	78.75	69.90	65.63 \dagger
Qwen-2.5-3B-RL-Aug	6.20	33.70	14.74	8.23	<u>2.76</u>	24.36	7.19	4.17	64.47	89.90	77.78	<u>66.03</u>

269 **Generalization to out-of-distribution puzzle tasks.** As illustrated in Figure 5, we evaluate puzzle
 270 tasks using checkpoints trained only on math and code data, without any puzzle-specific training.
 271 Despite this distribution shift, RL-Augmentation substantially improves performance on out-of-
 272 distribution puzzle benchmarks, indicating strong cross-task generalization.

273 **$J_{\phi_{SR}}$ -trained policies also improve other TTS algorithms.** Table 3 shows that models trained
 274 with RL-Augmentation consistently improve the performance of a wide range of test-time algorithms
 275 that rely on sequential revision, including Best-of- N , MCTS (Inoue et al., 2025), AB-MCTS variants
 276 (Inoue et al., 2025) and Mind Evolution (Lee et al., 2025b). This suggests that sequential revision
 277 capability serves as a fundamental backbone that raises the performance ceiling of diverse test-time
 278 inference strategies.

Table 3: Comparison of model performance on different test-time strategies. Best results in each column are in bold; second-best are underlined.

Model	OneShot	SeqRev	BoN	MCTS	Mind Evolution	AB-MCTS-A (Beta)	AB-MCTS-A (Gaussian)	AB-MCTS-M
Qwen-2.5-7B-RL	18.02	23.04	27.82	22.67	25.74	25.98	27.33	26.23
Qwen-2.5-7B-Multi-turn	<u>19.36</u>	25.49	28.68	26.10	26.96	28.55	<u>29.04</u>	27.82
Qwen-2.5-7B-PAG	19.49	<u>25.73</u>	30.39	<u>26.35</u>	<u>27.70</u>	<u>29.66</u>	28.06	<u>28.68</u>
Qwen-2.5-7B-RL-Augmentation	18.87	29.54	<u>29.29</u>	28.92	29.04	31.00	30.02	30.39

Table 4: Effect of revision and verification prompts on AIME25. **Left:** OneShot, Oracle-4, and SelfConf-4 accuracy across the four configurations: *RL* (no augmentation), *VerificationOnly* (verification prompts only), *RevisionOnly* (revision prompts only), and *RL-Augmentation* (both, our method). RL-Augmentation dominates under Oracle-4 and SelfConf-4. **Right:** TailConfidence distributions for RevisionOnly and RL-Aug; RL-Augmentation produces sharper, better-separated correct/incorrect distributions, indicating improved confidence calibration.

Model	OneShot	Oracle-4	SelfConf-4
Qwen-2.5-7B-RL	10.35	13.44	12.45
Qwen-2.5-7B-VerificationOnly	6.25	10.26	7.08
Qwen-2.5-7B-RevisionOnly	7.86	16.46	8.33
Qwen-2.5-7B-RL-Augmentation	9.48	21.16	15.83

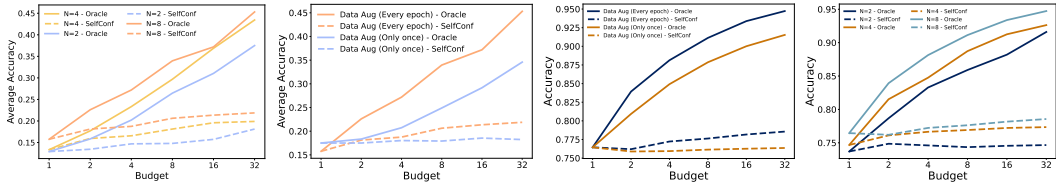


Figure 6: Ablations on data augmentation. From left to right: **(a)** effect of augmentation budget K (the maximum trajectory length allowed during augmentation) on AIME24; **(b)** continual versus one-shot augmentation on AIME24; **(c)** effect of augmentation frequency on MATH500; **(d)** effect of allowed revision length on MATH500. The MATH500 panels mirror the AIME24 trends, showing that the design choices generalize across datasets.

279 5.2 Ablation Study

280 **The role of revision and verification prompts.** As shown in Table 4, we investigate the roles
 281 of revision and verification prompts in sequential revision training. Qwen-2.5-7B-RevisionOnly
 282 augments RL-generated data using only revision prompts, while Qwen-2.5-7B-VerificationOnly
 283 augments data using only verification prompts. Qwen-2.5-7B-RL-Augmentation combines both
 284 prompts during data augmentation. Results under Oracle-4 indicate that revision prompts primarily
 285 improve *sequential revision capability*, whereas verification prompts do not directly enhance revision
 286 but instead train the model to assess solution validity.

287 To isolate the effect of verification prompts, we evaluate confidence calibration using Area Under
 288 the Receiver Operating Characteristic curve (AUROC) on AIME25. For each problem, we sample
 289 multiple model responses, label them as correct or incorrect based on the final answer, and compute
 290 the probability that a randomly sampled correct response receives a higher TailConfidence score
 291 than a randomly sampled incorrect one. The final AUROC is obtained by averaging over problems
 292 with sufficient correct samples. Under this evaluation, Qwen-2.5-7B-RL-Augmentation achieves a
 293 higher AUROC (74.1%) than RevisionOnly (72.1%), demonstrating that verification prompts improve
 294 confidence calibration. Thus, combining revision and verification prompts yields consistent gains by
 295 jointly enhancing revision capability and confidence estimation.

296 **Necessity of continual data augmentation.** As shown in Figure 6(b), augmenting data only at
 297 the initial epoch leads to substantially worse performance than continuous augmentation. This is
 298 because once the model stops making earlier mistakes, stale augmented data no longer provides
 299 informative learning signals, whereas regenerating data each epoch keeps the supervision aligned
 300 with the model’s current failure modes.

301 **Effect of augmentation budget.** We study the effect of the augmentation budget on training
 302 performance. As shown in Figure 6(a), increasing the augmentation budget consistently improves
 303 performance, indicating that our method exhibits strong scaling behavior with respect to the augmen-
 304 tation budget. At the same time, allowing a larger augmentation length N naturally incurs higher
 305 computational cost, reflecting a trade-off between performance and compute in data augmentation.

306 References

- 307 Chen, Z., Qin, X., Wu, Y., Ling, Y., Ye, Q., Zhao, W. X. and Shi, G. (2025). Pass@k training for
308 adaptively balancing exploration and exploitation of large reasoning models.
309 <https://arxiv.org/abs/2508.10751>
- 310 Chow, Y., Tennenholtz, G., Gur, I., Zhuang, V., Dai, B., Thiagarajan, S., Boutilier, C., Agarwal, R.,
311 Kumar, A. and Faust, A. (2025). Inference-aware fine-tuning for best-of-n sampling in large
312 language models.
313 <https://arxiv.org/abs/2412.15287>
- 314 DeepSeek-AI, Guo, D., Yang, D., Zhang, H., Song, J., Zhang, R., Xu, R., Zhu, Q., Ma, S.,
315 Wang, P., Bi, X., Zhang, X., Yu, X., Wu, Y., Wu, Z. F., Gou, Z., Shao, Z., Li, Z., Gao, Z.,
316 Liu, A., Xue, B., Wang, B., Wu, B., Feng, B., Lu, C., Zhao, C., Deng, C., Zhang, C., Ruan, C.,
317 Dai, D., Chen, D., Ji, D., Li, E., Lin, F., Dai, F., Luo, F., Hao, G., Chen, G., Li, G., Zhang, H.,
318 Bao, H., Xu, H., Wang, H., Ding, H., Xin, H., Gao, H., Qu, H., Li, H., Guo, J., Li, J., Wang, J.,
319 Chen, J., Yuan, J., Qiu, J., Li, J., Cai, J. L., Ni, J., Liang, J., Chen, J., Dong, K., Hu, K., Gao, K.,
320 Guan, K., Huang, K., Yu, K., Wang, L., Zhang, L., Zhao, L., Wang, L., Zhang, L., Xu, L., Xia, L.,
321 Zhang, M., Zhang, M., Tang, M., Li, M., Wang, M., Li, M., Tian, N., Huang, P., Zhang, P.,
322 Wang, Q., Chen, Q., Du, Q., Ge, R., Zhang, R., Pan, R., Wang, R., Chen, R. J., Jin, R. L.,
323 Chen, R., Lu, S., Zhou, S., Chen, S., Ye, S., Wang, S., Yu, S., Zhou, S., Pan, S., Li, S. S.,
324 Zhou, S., Wu, S., Ye, S., Yun, T., Pei, T., Sun, T., Wang, T., Zeng, W., Zhao, W., Liu, W.,
325 Liang, W., Gao, W., Yu, W., Zhang, W., Xiao, W. L., An, W., Liu, X., Wang, X., Chen, X.,
326 Nie, X., Cheng, X., Liu, X., Xie, X., Liu, X., Yang, X., Li, X., Su, X., Lin, X., Li, X. Q., Jin, X.,
327 Shen, X., Chen, X., Sun, X., Wang, X., Song, X., Zhou, X., Wang, X., Shan, X., Li, Y. K.,
328 Wang, Y. Q., Wei, Y. X., Zhang, Y., Xu, Y., Li, Y., Zhao, Y., Sun, Y., Wang, Y., Yu, Y., Zhang, Y.,
329 Shi, Y., Xiong, Y., He, Y., Piao, Y., Wang, Y., Tan, Y., Ma, Y., Liu, Y., Guo, Y., Ou, Y., Wang, Y.,
330 Gong, Y., Zou, Y., He, Y., Xiong, Y., Luo, Y., You, Y., Liu, Y., Zhou, Y., Zhu, Y. X., Xu, Y.,
331 Huang, Y., Li, Y., Zheng, Y., Zhu, Y., Ma, Y., Tang, Y., Zha, Y., Yan, Y., Ren, Z. Z., Ren, Z.,
332 Sha, Z., Fu, Z., Xu, Z., Xie, Z., Zhang, Z., Hao, Z., Ma, Z., Yan, Z., Wu, Z., Gu, Z., Zhu, Z.,
333 Liu, Z., Li, Z., Xie, Z., Song, Z., Pan, Z., Huang, Z., Xu, Z., Zhang, Z. and Zhang, Z. (2025).
334 Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning.
335 <https://arxiv.org/abs/2501.12948>
- 336 Fu, Y., Wang, X., Tian, Y. and Zhao, J. (2025). Deep think with confidence.
337 <https://arxiv.org/abs/2508.15260>
- 338 Georgiev, B., Gómez-Serrano, J., Tao, T. and Wagner, A. Z. (2025). Mathematical exploration and
339 discovery at scale. *arXiv preprint arXiv:2511.02864*.
- 340 Inoue, Y., Misaki, K., Imajuku, Y., Kuroki, S., Nakamura, T. and Akiba, T. (2025). Wider or deeper?
341 scaling llm inference-time compute with adaptive branching tree search.
342 <https://arxiv.org/abs/2503.04412>
- 343 Jain, A. K., Gonzalez-Pumariiega, G., Chen, W., Rush, A. M., Zhao, W. and Choudhury, S. (2025).
344 Multi-turn code generation through single-step rewards.
345 <https://arxiv.org/abs/2502.20380>
- 346 Jain, N., Han, K., Gu, A., Li, W.-D., Yan, F., Zhang, T., Wang, S., Solar-Lezama, A., Sen, K. and
347 Stoica, I. (2024). Livecodebench: Holistic and contamination free evaluation of large language
348 models for code.
349 <https://arxiv.org/abs/2403.07974>
- 350 Jiang, Y., Xiong, Y., Yuan, Y., Xin, C., Xu, W., Yue, Y., Zhao, Q. and Yan, L. (2025). Pag: Multi-
351 turn reinforced llm self-correction with policy as generative verifier.
352 <https://arxiv.org/abs/2506.10406>
- 353 Kumar, A., Zhuang, V., Agarwal, R., Su, Y., Co-Reyes, J. D., Singh, A., Baumli, K., Iqbal, S.,
354 Bishop, C., Roelofs, R., Zhang, L. M., McKinney, K., Shrivastava, D., Paduraru, C., Tucker, G.,
355 Precup, D., Behbahani, F. and Faust, A. (2024). Training language models to self-correct via
356 reinforcement learning.
357 <https://arxiv.org/abs/2409.12917>

- 358 Lange, R. T., Imajuku, Y. and Cetin, E. (2025). Shinkaevolve: Towards open-ended and sample-
359 efficient program evolution. *arXiv preprint arXiv:2509.19349*.
- 360 Lee, H., Oh, S., Kim, J., Shin, J. and Tack, J. (2025a). Revise: Learning to refine at test-time via
361 intrinsic self-verification.
362 <https://arxiv.org/abs/2502.14565>
- 363 Lee, K.-H., Fischer, I., Wu, Y.-H., Marwood, D., Baluja, S., Schuurmans, D. and Chen, X. (2025b).
364 Evolving deeper llm thinking.
365 <https://arxiv.org/abs/2501.09891>
- 366 Li, A., Wang, Y., Yuan, Z., Jegelka, S. and Wang, Y. (2025a). Lanpo: Bootstrapping language and
367 numerical feedback for reinforcement learning in llms.
368 <https://arxiv.org/abs/2510.16552>
- 369 Li, Y., Choi, D., Chung, J., Kushman, N., Schrittwieser, J., Leblond, R., Eccles, T., Keeling, J.,
370 Gimeno, F., Lago, A. D., Hubert, T., Choy, P., de Masson d’Autume, C., Babuschkin, I., Chen, X.,
371 Huang, P.-S., Welbl, J., Gowal, S., Cherepanov, A., Molloy, J., Mankowitz, D. J., Robson, E. S.,
372 Kohli, P., de Freitas, N., Kavukcuoglu, K. and Vinyals, O. (2022). Competition-level code genera-
373 tion with alphacode. *Science*, **378** 1092–1097.
374 <https://www.science.org/doi/abs/10.1126/science.abq1158>
- 375 Li, Y., Lyu, M. and Wang, L. (2025b). Learning to reason from feedback at test-time.
376 <https://arxiv.org/abs/2502.15771>
- 377 Liang, X., Li, Z.-Z., Gong, Y., Wang, Y., Zhang, H., Shen, Y., Wu, Y. N. and Chen, W. (2025). Sws:
378 Self-aware weakness-driven problem synthesis in reinforcement learning for llm reasoning.
379 <https://arxiv.org/abs/2506.08989>
- 380 Madaan, A., Tandon, N., Gupta, P., Hallinan, S., Gao, L., Wiegrefe, S., Alon, U., Dziri, N.,
381 Prabhumoye, S., Yang, Y., Gupta, S., Majumder, B. P., Hermann, K., Welleck, S.,
382 Yazdanbakhsh, A. and Clark, P. (2023). Self-refine: Iterative refinement with self-feedback.
383 <https://arxiv.org/abs/2303.17651>
- 384 Novikov, A., Vū, N., Eisenberger, M., Dupont, E., Huang, P.-S., Wagner, A. Z., Shirobokov, S.,
385 Kozlovskii, B., Ruiz, F. J., Mehrabian, A. et al. (2025). Alphaevolve: A coding agent for sci-
386 entific and algorithmic discovery. *arXiv preprint arXiv:2506.13131*.
- 387 Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C. L., Mishkin, P., Zhang, C., Agarwal, S.,
388 Slama, K., Ray, A., Schulman, J., Hilton, J., Kelton, F., Miller, L., Simens, M., Askell, A.,
389 Welinder, P., Christiano, P., Leike, J. and Lowe, R. (2022). Training language models to follow
390 instructions with human feedback.
391 <https://arxiv.org/abs/2203.02155>
- 392 Qu, Y., Setlur, A., Smith, V., Salakhutdinov, R. and Kumar, A. (2026). Pope: Learning to reason on
393 hard problems via privileged on-policy exploration.
394 <https://arxiv.org/abs/2601.18779>
- 395 Qu, Y., Zhang, T., Garg, N. and Kumar, A. (2024). Recursive introspection: Teaching language
396 model agents how to self-improve.
397 <https://arxiv.org/abs/2407.18219>
- 398 Shinn, N., Cassano, F., Berman, E., Gopinath, A., Narasimhan, K. and Yao, S. (2023). Reflexion:
399 Language agents with verbal reinforcement learning.
400 <https://arxiv.org/abs/2303.11366>
- 401 Stojanovski, Z., Stanley, O., Sharratt, J., Jones, R., Adefioye, A., Kaddour, J. and Köpf, A. (2025).
402 Reasoning gym: Reasoning environments for reinforcement learning with verifiable rewards.
403 <https://arxiv.org/abs/2505.24760>
- 404 Tang, Y., Zheng, K., Synnaeve, G. and Munos, R. (2025). Optimizing language models for inference
405 time objectives using reinforcement learning.
406 <https://arxiv.org/abs/2503.19595>

- 407 Walder, C. and Karkhanis, D. (2025). Pass@k policy optimization: Solving harder reinforcement
408 learning problems.
409 <https://arxiv.org/abs/2505.15201>
- 410 Wang, X., Wei, J., Schuurmans, D., Le, Q., Chi, E., Narang, S., Chowdhery, A. and Zhou, D. (2023).
411 Self-consistency improves chain of thought reasoning in language models.
412 <https://arxiv.org/abs/2203.11171>
- 413 Wang, Y., Su, S.-R., Zeng, Z., Xu, E., Ren, L., Yang, X., Huang, Z., He, X., Ma, L., Peng, B. et al.
414 (2025). Thetaevolve: Test-time learning on open problems. *arXiv preprint arXiv:2511.23473*.
- 415 Xie, J., Zhang, K., Chen, J., Zhu, T., Lou, R., Tian, Y., Xiao, Y. and Su, Y. (2024). Travelplanner: A
416 benchmark for real-world planning with language agents.
417 <https://arxiv.org/abs/2402.01622>
- 418 Xiong, W., Shi, C., Shen, J., Rosenberg, A., Qin, Z., Calandriello, D., Khalman, M., Joshi, R.,
419 Piot, B., Saleh, M., Jin, C., Zhang, T. and Liu, T. (2025). Building math agents with multi-turn
420 iterative preference learning.
421 <https://arxiv.org/abs/2409.02392>
- 422 Yang, A., Li, A., Yang, B., Zhang, B., Hui, B., Zheng, B., Yu, B., Gao, C., Huang, C., Lv, C.,
423 Zheng, C., Liu, D., Zhou, F., Huang, F., Hu, F., Ge, H., Wei, H., Lin, H., Tang, J., Yang, J.,
424 Tu, J., Zhang, J., Yang, J., Yang, J., Zhou, J., Zhou, J., Lin, J., Dang, K., Bao, K., Yang, K.,
425 Yu, L., Deng, L., Li, M., Xue, M., Li, M., Zhang, P., Wang, P., Zhu, Q., Men, R., Gao, R.,
426 Liu, S., Luo, S., Li, T., Tang, T., Yin, W., Ren, X., Wang, X., Zhang, X., Ren, X., Fan, Y., Su, Y.,
427 Zhang, Y., Zhang, Y., Wan, Y., Liu, Y., Wang, Z., Cui, Z., Zhang, Z., Zhou, Z. and Qiu, Z. (2025).
428 Qwen3 technical report.
429 <https://arxiv.org/abs/2505.09388>
- 430 Yang, A., Yang, B., Hui, B., Zheng, B., Yu, B., Zhou, C., Li, C., Li, C., Liu, D., Huang, F., Dong, G.,
431 Wei, H., Lin, H., Tang, J., Wang, J., Yang, J., Tu, J., Zhang, J., Ma, J., Yang, J., Xu, J., Zhou, J.,
432 Bai, J., He, J., Lin, J., Dang, K., Lu, K., Chen, K., Yang, K., Li, M., Xue, M., Ni, N., Zhang, P.,
433 Wang, P., Peng, R., Men, R., Gao, R., Lin, R., Wang, S., Bai, S., Tan, S., Zhu, T., Li, T., Liu, T.,
434 Ge, W., Deng, X., Zhou, X., Ren, X., Zhang, X., Wei, X., Ren, X., Liu, X., Fan, Y., Yao, Y.,
435 Zhang, Y., Wan, Y., Chu, Y., Liu, Y., Cui, Z., Zhang, Z., Guo, Z. and Fan, Z. (2024). Qwen2
436 technical report.
437 <https://arxiv.org/abs/2407.10671>
- 438 Yuksekogonul, M., Koceja, D., Li, X., Bianchi, F., McCaleb, J., Wang, X., Kautz, J., Choi, Y., Zou, J.,
439 Guestrin, C. and Sun, Y. (2026). Learning to discover at test time.
440 <https://arxiv.org/abs/2601.16175>
- 441 Zhang, X., Zhang, Y., Sun, H., Feng, K., Lu, C., Yang, C. and Meng, H. (2026). Critique-grpo:
442 Advancing llm reasoning with natural language and numerical feedback.
443 <https://arxiv.org/abs/2506.03106>

444 **A Sequential Revision Protocol**

445 In this section, we elaborate on the test-time protocol described in Protocol 1 and illustrate its practical
 446 usage with several concrete examples, including a mathematical reasoning task in Appendix A.1 and
 447 a coding task in Appendix A.2.

Protocol 1 Sequential Revision Protocol

```

1: Input: revision budget  $K$ , question  $x$ , model  $\pi$ 
2: Sample initial response  $y_1 \sim \pi(\cdot | x)$ 
3: for  $i = 2$  to  $K$  do
4:   Generate revised response  $y_i \sim \pi(\cdot | x, y_{i-1}, \text{Feedback}(y_{i-1}))$ 
5:   if  $r^*(x, y_i) = 1$  then
6:     # Early stopping
7:     Set  $\hat{y} = y_i$ 
8:     break for loop
9:   end if
10: end for
11: if No early stopping then
12:    $\hat{y} = y_K$ 
13: end if
  
```

448 **A.1 Example Math Task**

449 In mathematical tasks, the feedback consists solely of the sentence “The response is incorrect.”
 450 (marked in red in the figure), without offering further diagnostic information.

Input (Turn 1)

451 Solve the following math problem. Show brief reasoning and provide the final answer.
 Problem: A 90° rotation around $-4 - 5i$ in the clockwise direction is applied to 3 (as a complex number). What is the resulting complex number?

LLM Output (Turn 1)

452 To solve this problem, we need to perform a 90° clockwise rotation around the complex number $-4 - 5i$.
 First, it’s useful to understand that a 90° clockwise rotation of a complex number z around another complex number a is equivalent to the following transformation:

$$z' = a + (a - z) \cdot i$$

However, since we are rotating around $-4 - 5i$, let’s denote this center of rotation as $a = -4 - 5i$. Our complex number $z = 3$ needs to be rotated.

1. Find the difference between z and a :

$$z - a = 3 - (-4 - 5i) = 3 + 4 + 5i = 7 + 5i$$

2. Perform the 90° clockwise rotation by multiplying this difference by i and adding it back to a :

$$z' = a + (z - a) \cdot i = (-4 - 5i) + (7 + 5i) \cdot i$$

Next, simplify the expression:

$$(7 + 5i) \cdot i = 7i + 5i^2$$

Since $i^2 = -1$:

$$7i + 5(-1) = 7i - 5$$

Therefore, the final answer is $-9 + 2i$.

Input (Turn 2)

Problem: A 90° rotation around $-4 - 5i$ in the clockwise direction is applied to 3 (as a complex number). What is the resulting complex number?

Previous Attempt: To solve this problem, we need to perform a 90° clockwise rotation around the complex number $-4 - 5i$.

First, it's useful to understand that a 90° clockwise rotation of a complex number z around another complex number a is equivalent to the following transformation:

$$z' = a + (a - z) \cdot i$$

However, since we are rotating around $-4 - 5i$, let's denote this center of rotation as $a = -4 - 5i$. Our complex number $z = 3$ needs to be rotated.

1. Find the difference between z and a :

$$z - a = 3 - (-4 - 5i) = 3 + 4 + 5i = 7 + 5i$$

2. Perform the 90° clockwise rotation by multiplying this difference by i and adding it back to a :

$$z' = a + (z - a) \cdot i = (-4 - 5i) + (7 + 5i) \cdot i$$

Next, simplify the expression:

$$(7 + 5i) \cdot i = 7i + 5i^2$$

Since $i^2 = -1$:

$$7i + 5(-1) = 7i - 5$$

...

So, the resulting complex number after the rotation is $-9 + 2i$.

Therefore, the final answer is $\boxed{-9 + 2i}$.

Verifiable Feedback: **The response is incorrect.**

Now produce a corrected solution. Keep reasoning concise.

453

LLM Output (Turn 2)

To correctly rotate a complex number 3 around $-4 - 5i$ by 90° clockwise, we need to use the appropriate transformation for rotations in the complex plane.

Given a complex number $z = 3$, and the center of rotation $a = -4 - 5i$, a 90° clockwise rotation of z around a can be represented as:

$$z' = a + (a - z) \cdot (-i)$$

Let's go through the steps in detail:

1. **Calculate $a - z$:**

$$a - z = (-4 - 5i) - 3 = -4 - 5i - 3 = -7 - 5i$$

2. **Multiply by $-i$:**

$$(a - z) \cdot (-i) = (-7 - 5i) \cdot (-i) = (-7)(-i) + (-5i)(-i) = 7i + 5i^2$$

Since $i^2 = -1$:

$$7i + 5(-1) = 7i - 5$$

3. **Add this result back to a :**

$$z' = a + (7i - 5) = (-4 - 5i) + (-5 + 7i)$$

$$z' = -4 - 5i - 5 + 7i$$

Combine the real and imaginary parts:

$$z' = (-4 - 5) + (-5i + 7i) = -9 + 2i$$

Thus, after performing the 90° clockwise rotation, the resulting complex number is

$$\boxed{-9 + 2i}$$

454

Input (Turn 3)

Problem: A 90° rotation around $-4 - 5i$ in the clockwise direction is applied to 3 (as a complex number). What is the resulting complex number?

Previous Attempt: To correctly rotate a complex number 3 around $-4 - 5i$ by 90° clockwise, we need to use the appropriate transformation for rotations in the complex plane.

Given a complex number $z = 3$, and the center of rotation $a = -4 - 5i$, a 90° clockwise rotation of z around a can be represented as:

$$z' = a + (a - z) \cdot (-i)$$

Let's go through the steps in detail:

1. **Calculate $a - z$:**

$$a - z = (-4 - 5i) - 3 = -4 - 5i - 3 = -7 - 5i$$

2. **Multiply by $-i$:**

$$(a - z) \cdot (-i) = (-7 - 5i) \cdot (-i) = (-7)(-i) + (-5i)(-i) = 7i + 5i^2$$

Since $i^2 = -1$:

$$7i + 5(-1) = 7i - 5$$

3. **Add this result back to a :**

$$z' = a + (7i - 5) = (-4 - 5i) + (-5 + 7i)$$

$$z' = -4 - 5i - 5 + 7i$$

Combine the real and imaginary parts:

$$z' = (-4 - 5) + (-5i + 7i) = -9 + 2i$$

Thus, after performing the 90° clockwise rotation, the resulting complex number is

$$\boxed{-9 + 2i}.$$

Verifiable Feedback: **The response is incorrect.**

Now produce a corrected solution. Keep reasoning concise.

LLM Output (Turn 3)

To find the result after a 90° clockwise rotation around $-4 - 5i$, we can use the following transformation steps. We start by considering a complex number $z = 3$ and the center of rotation $a = -4 - 5i$:

1. **Move z to the origin relative to a :**

$$z - a = 3 - (-4 - 5i) = 3 + 4 + 5i = 7 + 5i$$

2. **Rotate by 90° clockwise:** A 90° clockwise rotation of a complex number w is equivalent to multiplying w by $-i$. So, we need to multiply $7 + 5i$ by $-i$:

$$(7 + 5i) \times (-i) = 7 \times (-i) + 5i \times (-i) = -7i - 5i^2$$

Since $i^2 = -1$:

$$-7i - 5(-1) = -7i + 5 = 5 - 7i$$

3. **Translate back by adding a :** Adding a to the resultant, we get:

$$z' = (5 - 7i) + (-4 - 5i) = 5 - 4 - 7i - 5i = 1 - 12i$$

Thus, the resulting complex number after the 90° clockwise rotation is $\boxed{1 - 12i}$.

455

456

457 A.2 Example Coding Task

458 At test time for coding tasks, we follow an iterative execution-based evaluation protocol. The model
459 first generates a candidate program, which is then executed against the publicly available test cases
460 provided by the problem. If the program fails any test, the resulting error messages or failing test
461 feedback are collected and appended to the context for the next generation step, allowing the model
462 to revise its solution based on execution feedback. Note that passing the public test cases does not
463 guarantee correctness: a solution is considered correct only if it passes all test cases, including both
464 public and private ones.

Input (Turn 1)

You are an expert Python programmer. You will be given a problem specification and must produce a correct Python program that passes the tests.

Question: You are given a permutation $P = (P_1, P_2, \dots, P_N)$ of $(1, 2, \dots, N)$. You want to satisfy $P_i = i$ for all $i = 1, 2, \dots, N$ by performing the following operation zero or more times:

- Choose an integer k such that $1 \leq k \leq N$. If $k \geq 2$, sort the 1-st through $(k - 1)$ -th terms of P in ascending order. Then, if $k \leq N - 1$, sort the $(k + 1)$ -th through N -th terms of P in ascending order.

It can be proved that under the constraints of this problem, it is possible to satisfy $P_i = i$ for all $i = 1, 2, \dots, N$ with a finite number of operations for any P . Find the minimum number of operations required. You have T test cases to solve.

Input ...

Output

Print T lines. The i -th line should contain the answer for the i -th test case.

Constraints ...

Sample Input 1

- 3
- 5
- ...

Sample Output 1

- 1
- ...

Format: Read inputs from stdin and write the answer to stdout. Enclose your final code in the following box.

```
# YOUR CODE HERE
```

Answer: Provide *only* the final Python code.

465

LLM Output (Turn 1)

To solve this problem, we aim to determine the minimum number of operations required to transform a given permutation P into the identity permutation $(1, 2, \dots, N)$. The key insight is that each operation can be viewed as a global reordering step that partially resolves disorder in the permutation. Intuitively, such operations progressively eliminate inconsistencies between positions and values.

```
...  
def min_operations(N, P):  
    count = 0  
    for i in range(N):  
        if P[i] != i + 1:  
            count += 1  
    return count  
...  
for _ in range(T):  
    N = int(input())  
    P = list(map(int, input().split()))  
    print(min_operations(N, P))
```

466

Input (Turn 2)

You are an expert Python programmer. Given a problem specification, produce a correct Python program that passes all tests.

Question: Given a permutation $P = (P_1, \dots, P_N)$ of $(1, \dots, N)$, the goal is to transform P into the identity permutation ...

Input / Output. ...

Previous Codes

code_v1

```
# Your task is to find the minimum number of operations required
to sort the permutation P.

def min_operations(N, P):
    operations = 0
    for i in range(N):
        while P[i] != i + 1:
            ...
    N = case[0]
    P = case[1:]
    print(min_operations(N, P))
```

467

Issues with code_v1

For the input

```
3
5
2 1 3 5 4
3
1 2 3
7
3 2 1 7 5 6 4
```

the expected output is

```
1
0
2
```

but code_v1 results in a runtime error. The failure is primarily due to incorrect handling of index-based updates and repeated linear searches ($P.index$), which lead to inconsistent state updates and excessive runtime overhead.

LLM Output (Turn 2)

To solve this problem, we aim to determine the minimum number of operations required to transform a given permutation P into the identity permutation $(1, 2, \dots, N)$. The key insight is that each operation can be viewed as a global reordering step that partially resolves disorder in the permutation. Intuitively, such operations progressively eliminate inconsistencies between positions and values.

...

```
def min_operations(N, P):
    count = 0
    for i in range(N):
        if P[i] != i + 1:
            count += 1
    return count

...

for _ in range(T):
    N = int(input())
    P = list(map(int, input().split()))
    print(min_operations(N, P))
```

468

469 **B Limitation and Future Direction**

470 **Limitation.** REVES requires a ground-truth verifier during training. Math, coding, puzzle, and
 471 circle-packing benchmarks all admit deterministic verifiers, so REVES applies directly. For tasks
 472 without such verifiers (free-form QA, creative writing, RLHF-style preference learning), REVES
 473 does not directly extend; the verification prompts trained alongside revision data can be repurposed
 474 for self-verification at test time, but their reliability depends on the policy’s own calibration.

475 **Future direction.** The verifier requirement points to a natural class of problems where REVES is
 476 especially well suited: scientific tasks whose objectives admit a clear, automated reward. Examples
 477 include solving partial differential equations (PDEs), where residual norms or analytic checks serve
 478 as verifiers; symbolic regression, where fit and complexity supply a numeric score; theorem proving,
 479 where a proof checker provides binary feedback; and discovery-style optimization tasks like circle
 480 packing (Appendix H), where the verifier is the geometric solver itself. In all of these settings, the
 481 cost of producing the verifier signal is offset by an essentially unlimited supply of training problems
 482 and a precise notion of correctness, both of which line up well with REVES’s reliance on per-state
 483 one-step recovery as the gradient signal. Scaling REVES to such scientific-discovery domains is a
 484 natural next step.

485 **C Proofs**

486 **C.1 Proof of Theorem 3.2**

487 **Theorem C.1** (Objective mismatch). *For any revision length $K \geq 2$, there exist policies π_1, π_2 and*
 488 *a problem distribution \mathcal{X} such that*

$$J_{\text{OneShot}}(\pi_1) = J_{\text{OneShot}}(\pi_2), \quad J_{\phi_{\text{SR}}}(\pi_2) - J_{\phi_{\text{SR}}}(\pi_1) \geq \Delta(K),$$

489 where $\Delta(K) > 0$.

490 *Proof. Construction.* Take \mathcal{X} supported on a single problem x with $N \geq 2$ candidate answers
 491 $\{y^{(1)}, \dots, y^{(N)}\}$ and ground-truth y^* drawn uniformly from this set; feedback is binary correctness.
 492 Define two history-independent policies, each sampling from a fixed distribution at every step:

$$\pi_1(\cdot | \cdot) = \delta_{y^{(1)}}, \quad \pi_2(\cdot | \cdot) = \text{Uniform}\{y^{(1)}, \dots, y^{(N)}\},$$

493 with π_2 ’s draws independent across steps.

494 **One-shot.** Since y^* is uniform,

$$J_{\text{OneShot}}(\pi_1) = \mathbb{P}(y^* = y^{(1)}) = \frac{1}{N}, \quad J_{\text{OneShot}}(\pi_2) = \sum_{i=1}^N \frac{1}{N} \cdot \frac{1}{N} = \frac{1}{N},$$

495 so $J_{\text{OneShot}}(\pi_1) = J_{\text{OneShot}}(\pi_2)$.

496 **Sequential revision (budget K).** π_1 outputs $y^{(1)}$ at every step and therefore succeeds iff $y^* = y^{(1)}$,
 497 giving $J_{\phi_{\text{SR}}}(\pi_1) = 1/N$. For π_2 , the K i.i.d. uniform draws are independent of y^* , so the probability
 498 that none equals y^* is $(1 - 1/N)^K$, giving $J_{\phi_{\text{SR}}}(\pi_2) = 1 - (1 - 1/N)^K$.

499 **Gap.** For $K \geq 2$, monotonicity of $K \mapsto (1 - 1/N)^K$ on $K \geq 1$ gives $(1 - 1/N)^K \leq (1 - 1/N)^2 =$
 500 $1 - 2/N + 1/N^2$, hence

$$\Delta(K) := J_{\phi_{\text{SR}}}(\pi_2) - J_{\phi_{\text{SR}}}(\pi_1) = 1 - (1 - \frac{1}{N})^K - \frac{1}{N} \geq \frac{N-1}{N^2} > 0. \quad \square$$

501 **C.2 Proof of Theorem 3.1**

502 The argument uses the recovery decomposition $J_\phi(\pi) = \sum_z \rho_\pi^\phi(z) V_\pi(z)$, which holds for any
 503 first-success transition strategy ϕ paired with any policy π . The proof of Lemma 4.1 (Appendix C.3)
 504 establishes this for $\phi = \phi_{\text{SR}}$; the same three steps – (i) writing the success event at step t as
 505 $\{\tau \geq t, y_t \in G_x\}$, (ii) Fubini against $\sum_{t=1}^K \mathbf{1}\{\tau \geq t\} = \tau$, (iii) commuting the visit-count sum with
 506 the outer expectation – depend only on the first-success stopping rule and the Markov property of π ,
 507 so they generalize verbatim from $\rho_\pi^{\phi_{\text{SR}}}$ to ρ_π^ϕ for any $\phi \in \Phi_{\mathcal{R}}$ that early-stops.

508 *Proof of Theorem 3.1.* By the recovery decomposition applied to π_1 and π_0 ,

$$J_\phi(\pi_1) - J_\phi(\pi_0) = \sum_z \rho_{\pi_1}^\phi(z) V_{\pi_1}(z) - \sum_z \rho_{\pi_0}^\phi(z) V_{\pi_0}(z).$$

509 Add and subtract $\sum_z \rho_{\pi_0}^\phi(z) V_{\pi_1}(z)$ and group:

$$J_\phi(\pi_1) - J_\phi(\pi_0) = \underbrace{\sum_z \rho_{\pi_0}^\phi(z) (V_{\pi_1}(z) - V_{\pi_0}(z))}_{(\star)} + \underbrace{\sum_z (\rho_{\pi_1}^\phi(z) - \rho_{\pi_0}^\phi(z)) V_{\pi_1}(z)}_{(\dagger)}. \quad (2)$$

510 **Bounding** (\dagger) . Since $V_{\pi_1}(z) \in [0, 1]$,

$$(\dagger) \geq - \sum_z |\rho_{\pi_1}^\phi(z) - \rho_{\pi_0}^\phi(z)| = -\|\rho_{\pi_1}^\phi - \rho_{\pi_0}^\phi\|_1. \quad (3)$$

511 **Bounding** (\star) . We claim the pointwise inequality

$$\rho_{\pi_0}^\phi(z) (V_{\pi_1}(z) - V_{\pi_0}(z)) \geq \frac{1}{C_\phi} \rho_{\pi_0}^{\phi_{\text{SR}}}(z) (V_{\pi_1}(z) - V_{\pi_0}(z)) \quad \text{for every } z. \quad (4)$$

512 *Case 1:* $z \in \text{supp}(\rho_{\pi_0}^\phi)$. By (C2), $V_{\pi_1}(z) - V_{\pi_0}(z) \geq 0$. By (C1), $\rho_{\pi_0}^\phi(z) \geq C_\phi^{-1} \rho_{\pi_0}^{\phi_{\text{SR}}}(z) \geq 0$.
 513 Multiplying both sides of (C1) by the nonnegative quantity $V_{\pi_1}(z) - V_{\pi_0}(z)$ preserves the inequality
 514 and gives (4).

515 *Case 2:* $z \notin \text{supp}(\rho_{\pi_0}^\phi)$. Then $\rho_{\pi_0}^\phi(z) = 0$, and (C1) forces $0 \geq C_\phi^{-1} \rho_{\pi_0}^{\phi_{\text{SR}}}(z)$. Combined with
 516 $\rho_{\pi_0}^{\phi_{\text{SR}}}(z) \geq 0$, this gives $\rho_{\pi_0}^{\phi_{\text{SR}}}(z) = 0$, so both sides of (4) vanish.

517 Summing (4) over z ,

$$(\star) \geq \frac{1}{C_\phi} \sum_z \rho_{\pi_0}^{\phi_{\text{SR}}}(z) (V_{\pi_1}(z) - V_{\pi_0}(z)). \quad (5)$$

518 **Converting** (5) to $J_{\phi_{\text{SR}}}(\pi_1) - J_{\phi_{\text{SR}}}(\pi_0)$. Apply the recovery decomposition to $\phi = \phi_{\text{SR}}$ at π_1 and
 519 π_0 , then add and subtract $\sum_z \rho_{\pi_0}^{\phi_{\text{SR}}}(z) V_{\pi_1}(z)$:

$$J_{\phi_{\text{SR}}}(\pi_1) - J_{\phi_{\text{SR}}}(\pi_0) = \sum_z \rho_{\pi_0}^{\phi_{\text{SR}}}(z) (V_{\pi_1}(z) - V_{\pi_0}(z)) + \sum_z (\rho_{\pi_1}^{\phi_{\text{SR}}}(z) - \rho_{\pi_0}^{\phi_{\text{SR}}}(z)) V_{\pi_1}(z).$$

520 Rearranging and using $V_{\pi_1}(z) \in [0, 1]$,

$$\sum_z \rho_{\pi_0}^{\phi_{\text{SR}}}(z) (V_{\pi_1}(z) - V_{\pi_0}(z)) \geq [J_{\phi_{\text{SR}}}(\pi_1) - J_{\phi_{\text{SR}}}(\pi_0)] - \|\rho_{\pi_1}^{\phi_{\text{SR}}} - \rho_{\pi_0}^{\phi_{\text{SR}}}\|_1. \quad (6)$$

521 **Combine.** Chaining (6) into (5) and adding (3),

$$J_\phi(\pi_1) - J_\phi(\pi_0) \geq \frac{1}{C_\phi} [J_{\phi_{\text{SR}}}(\pi_1) - J_{\phi_{\text{SR}}}(\pi_0)] - \|\rho_{\pi_1}^\phi - \rho_{\pi_0}^\phi\|_1 - \frac{1}{C_\phi} \|\rho_{\pi_1}^{\phi_{\text{SR}}} - \rho_{\pi_0}^{\phi_{\text{SR}}}\|_1. \quad \square$$

522 **Connection to our method.** The third term $C_\phi^{-1} \|\rho_{\pi_1}^{\phi_{\text{SR}}} - \rho_{\pi_0}^{\phi_{\text{SR}}}\|_1$ in the bound is the price for
 523 placing the on-policy SR test-time objective $J_{\phi_{\text{SR}}}(\pi) = \sum_z \rho_\pi^{\phi_{\text{SR}}}(z) V_\pi(z)$ on the right-hand side: as
 524 the policy improves, SR's own visit measure $\rho_\pi^{\phi_{\text{SR}}}$ also shifts, and that shift contaminates $J_{\phi_{\text{SR}}}(\pi_1) -$
 525 $J_{\phi_{\text{SR}}}(\pi_0)$. The off-policy fixed-reference quantity

$$J_{\phi_{\text{SR}}}^{\mathcal{R}}(\pi) := \sum_z \rho_{\pi_0}^{\phi_{\text{SR}}}(z) V_\pi(z)$$

526 absorbs this shift by holding the visit measure at the baseline; $J_{\phi_{\text{SR}}}^{\mathcal{R}}(\pi_1) - J_{\phi_{\text{SR}}}^{\mathcal{R}}(\pi_0) =$
 527 $\sum_z \rho_{\pi_0}^{\phi_{\text{SR}}}(\cdot) (V_{\pi_1} - V_{\pi_0})$ exactly, and (5) terminates the proof in one step:

$$J_\phi(\pi_1) - J_\phi(\pi_0) \geq \frac{1}{C_\phi} [J_{\phi_{\text{SR}}}^{\mathcal{R}}(\pi_1) - J_{\phi_{\text{SR}}}^{\mathcal{R}}(\pi_0)] - \|\rho_{\pi_1}^\phi - \rho_{\pi_0}^\phi\|_1.$$

528 This off-policy form is exactly the training signal REVES optimizes (Section 4): Stage I freezes the
 529 rollout distribution at the current π_0 , generating SR trajectories whose visited states constitute samples
 530 from $\rho_{\pi_0}^{\phi_{\text{SR}}}$; Stage II then trains V_π on those fixed states via single-turn RL. Each epoch's rollout
 531 refresh resets the baseline so the ℓ_1 shift on $\rho_{\pi_0}^{\phi_{\text{SR}}}$ stays small. The two forms are mathematically
 532 equivalent; the off-policy form makes the algorithmic interpretation transparent.

533 **Interpretation.** The bound is conditional, not unconditional. It bites when (i) SR’s revision-input
534 distribution covers ϕ ’s (small C_ϕ in (C1)), (ii) the policy update raises one-step recovery on those
535 states ((C2)), and (iii) the visit measures ρ^ϕ and $\rho^{\phi_{\text{SR}}}$ are stable under the update (small ℓ_1 shifts).
536 The empirical confirmation that $J_{\phi_{\text{SR}}}$ -trained policies improve other revision-using TTS algorithms
537 (tree search, BoN-with-revision, evolutionary refinement) is in Table 3.

538 C.3 Proof of Lemma 4.1

539 *Proof.* We prove the three identities in (1) in order.

540 *First identity.* SR stops at the first success, so $r^*(x, y_\tau) = 1 \iff \exists t \leq K : r^*(x, y_t) = 1$, and the
541 success event at step t equals $\{\tau \geq t, y_t \in G_x\}$, where $G_x := \{y : r^*(x, y) = 1\}$. By the tower
542 property and the fact that π_θ is Markov in z_t ,

$$\mathbb{P}(\tau \geq t, y_t \in G_x) = \mathbb{E}[\mathbf{1}\{\tau \geq t\} \mathbb{P}(y_t \in G_x \mid z_t)] = \mathbb{E}[\mathbf{1}\{\tau \geq t\} V_\pi(z_t)].$$

543 Summing over $t \leq K$,

$$J_{\phi_{\text{SR}}}(\theta) = \mathbb{P}(\exists t \leq K : r^*(x, y_t) = 1) = \sum_{t=1}^K \mathbb{E}[\mathbf{1}\{\tau \geq t\} V_\pi(z_t)].$$

544 *Second identity.* Using $\sum_{t=1}^K \mathbf{1}\{\tau \geq t\} = \tau$ and Fubini,

$$\sum_{t=1}^K \mathbb{E}[\mathbf{1}\{\tau \geq t\} V_\pi(z_t)] = \mathbb{E}\left[\sum_{t=1}^{\tau} V_\pi(z_t)\right].$$

545 *Third identity.* Substituting the definition $V_\pi(z) = \mathbb{E}_{y' \sim \pi_\theta(\cdot|z)}[r^*(x, y')]$ and exchanging the visit-
546 count sum with the outer expectation gives

$$\mathbb{E}\left[\sum_{t=1}^{\tau} V_\pi(z_t)\right] = \sum_z V_\pi(z) \mathbb{E}\left[\sum_{t=1}^{\tau} \mathbf{1}\{z_t = z\}\right] = \sum_z \rho_\theta(z) V_\pi(z) = \sum_z \rho_\theta(z) \mathbb{E}_{y' \sim \pi_\theta(\cdot|z)}[r^*(x, y')].$$

□

547 D Training Prompt Template

548 In this section, we present the training prompt templates used for data augmentation. An important
549 distinction from the test-time setting is that we do not append execution error information to the
550 context during training. Instead, as illustrated in Figure 10 and Figure 11, the training prompts rely
551 on structured supervision rather than explicit error feedback.

552 E Implementation Details

553 In this section, we provide details of the training and evaluation implementations used in our
554 experiments. All experiments are conducted on machines equipped with $8 \times$ NVIDIA H200 GPUs.

555 **RL Baseline.** Our RL baseline is implemented using the official `Skywork-0R1` codebase. We
556 directly adopt their provided training scripts and configurations, which optimize a single-shot response
557 using reinforcement learning. No revision or multi-turn generation is involved during training for this
558 baseline.

559 **Multi-Turn and PAG Baselines.** The Multi-Turn and PAG baselines are implemented based on
560 the `Policy-As-GenVerifier` codebase. To ensure comparability with our setting, we make two
561 modifications to the original implementation. First, we increase the maximum response length to
562 8092 tokens to match our training setup. Second, we adapt the reward computation function so that it
563 correctly evaluates responses under our datasets, while keeping the reward definition itself identical to
564 that used in our methods. Apart from these changes, the original training pipelines and optimization
565 settings are preserved.

Prompt Template for Revision in Math Tasks

Problem
 A frequency distribution of the scores for Mr. Sampson’s algebra class is shown. What percent of the class received a score in the 60% – 69% range?

Test Scores	Frequencies
90% – 100%	IIII
80% – 89%	IIII IIII
70% – 79%	IIII II
60% – 69%	IIII I
Below 60%	II

Previous Answer
 To solve this problem, we first count the number of students in each score range. The total number of students is computed as $4 + 8 + 6 + 5 + 2 = 25$. The number of students scoring in the 60% – 69% range is 5. ...

Verifiable Feedback:
 The previous answer is incorrect.

Instruction:
 Now produce a corrected solution. Identify and fix potential mistakes in the previous attempt. Keep the reasoning concise and provide a corrected final answer. Do not assume access to the correct answer.

Figure 7: Prompt template for the revision stage with binary correctness feedback. Concrete problem content and partial incorrect answers are shown for illustration, with remaining content omitted for brevity.

Prompt Template for Verification in Math Tasks

Problem
 A frequency distribution of the scores for Mr. Sampson’s algebra class is shown. What percent of the class received a score in the 60% – 69% range?

Test Scores	Frequencies
90% – 100%	IIII
80% – 89%	IIII IIII
70% – 79%	IIII II
60% – 69%	IIII I
Below 60%	II

Proposed Answer
 To solve this problem, we first identify the number of students in each score range from the frequency distribution. The total number of students is computed as $4 + 8 + 6 + 5 + 2 = 25$. Among them, 5 students fall into ...

Instruction:
 Judge whether the proposed answer is correct. Respond concisely and end the final line with VERDICT: True or VERDICT: False.

Figure 8: Prompt template used for the verification stage. The verifier observes the problem and a proposed solution and outputs a binary correctness judgment.

566 **RL-Augmentation (Ours).** Our method builds on top of the same training framework as the
 567 baselines. The key difference lies in the data augmentation procedure. Before training, we first
 568 perform an offline augmentation pass to initialize the augmented dataset. During training, at the
 569 beginning of each epoch, we further augment the dataset by generating new revision trajectories
 570 using the current policy, and the newly generated samples are added to the training data. When
 571 generating augmented data, we allow a maximum revision budget of $N = 8$. The training prompts
 572 used for augmentation follow the templates described in Appendix D. The model is then trained on
 573 the mixture of original and augmented data using the same RL objective as the baseline.

574 **Evaluation Protocol.** Unless otherwise specified, all experiments are evaluated using our unified
 575 evaluation framework, which will be publicly released upon paper acceptance. Given an input

576 problem, the model follows the corresponding test-time protocol (e.g., sequential revision when
577 applicable) to generate a final response. The generated output is serialized into a JSON file and
578 evaluated by an external, task-specific verifier. This pipeline is applied uniformly across all tasks and
579 methods to ensure fair comparison: generation and evaluation are fully decoupled, and correctness is
580 determined solely by the external verifier rather than any internal model signals. During evaluation,
581 the maximum response length is set to 16,384 tokens for coding tasks and AIME24/25, and 8,096
582 tokens for MATH500. The decoding temperature is fixed to 0.6 for coding tasks and 1.0 for all other
583 tasks, and the confidence threshold in the SelfConf setting is set to $c = 0.5$.

584 For math tasks, when evaluating the Multi-Turn and PAG baselines, we follow the original evaluation
585 implementation provided in the Policy-As-GenVerifier repository to remain faithful to their
586 protocol. Under this evaluation pipeline, we use the same maximum response length settings for 7B
587 models as described above. For 3B models with revision budget $K = 4$, where the effective sequence
588 length exceeds the model context limit, we reduce the maximum response length to 4,096 tokens.
589 All math evaluations are conducted using `math_eval`, ensuring that results are directly comparable
590 across evaluation settings.

591 **F Method Details: Multi-Turn Comparison, Efficiency, and Test-Time** 592 **Sensitivity**

593 **F.1 Comparison with Naive Multi-Turn RL**

594 A natural alternative to our offline scheme is online sequential revision with standard multi-turn credit
595 assignment. Two such schemes appear in the recent literature. The first, *trajectory-level broadcast*,
596 samples G trajectories of M turns per prompt, scores each trajectory by its final reward, computes
597 a trajectory-level advantage, and broadcasts that single advantage uniformly to every turn in the
598 trajectory. The second, *pooled per-response*, pools all $G \times M$ responses across the G trajectories,
599 treats each response as an independent sample, and computes advantages within the pool.

600 Neither scheme suits our setting. Sequential revision is structurally different from a vanilla multi-turn
601 task: each turn is a fresh attempt at the *same* problem x , not a step in a longer-horizon decision
602 process with intermediate sub-rewards. Trajectory-level broadcast assigns identical credit to every
603 turn and therefore ignores that some intermediate responses are much closer to correctness than others.
604 Pooled per-response also fails to capture this structure: the wrong intermediate responses are treated
605 identically, even though the near-miss failures are exactly the ones we most want the policy to learn
606 from. Our offline selection-then-replay scheme bypasses both pathologies. We identify near-miss
607 failures explicitly through the success-trajectory filter and replay them as single-turn revision prompts,
608 where standard single-turn credit assignment applies cleanly.

609 **F.2 Wall-Clock Efficiency**

610 Beyond avoiding credit-assignment pitfalls, our offline scheme is substantially faster than online
611 multi-turn training. Multi-turn RL pays a length- M serial sequential-revision cost per prompt at
612 every gradient step, multiplied by the G rollouts needed to estimate the advantage. Our method pays
613 the length- K serial cost *only once* per prompt at the offline stage, and during the online RL stage
614 performs only single-turn sampling on the extracted prompts, which can be parallelized across N
615 Best-of- N samples. Single-turn parallel sampling is much cheaper per gradient step than M -turn
616 serial generation across G rollouts, since serial generation does not benefit from the parallelism that
617 modern inference engines provide for single-turn sampling.

618 Empirically, on $8 \times$ NVIDIA H200 GPUs, three full iterations of our method completed in roughly 30
619 hours, while the equivalent online multi-turn baseline required over 48 hours on identical hardware
620 (Table 5). With a more carefully optimized infrastructure, Stage I data generation and Stage II RL
621 training could run concurrently (Stage II first trains on the original RL data while Stage I generates
622 augmented data asynchronously), yielding further gains.

623 **F.3 Sensitivity to the Confidence Threshold**

624 We study the sensitivity of our test-time stopping rule to the confidence threshold c used in the
625 TailConfidence stopping criterion (Section 4). Table 6 reports accuracy and the average number

Table 5: Per-iteration wall-clock breakdown of our method on $8 \times$ NVIDIA H200 GPUs.

Iteration	Data Generation (Stage I)	RL Training (Stage II)	Total
Iter 1	3h 4m	7h 15m	10h 19m
Iter 2	2h 41m	7h 26m	10h 7m
Iter 3	3h 13m	6h 23m	9h 36m
Total	8h 58m	21h 4m	30h 2m

626 of revision rounds used at budget $K = 16$ across three math benchmarks. Higher thresholds enforce
 627 stricter stopping (more rounds before accepting a response), and accuracy increases monotonically
 628 with the threshold over the range we tested. The number of rounds also grows with the threshold,
 629 reflecting the cost-quality trade-off.

Table 6: Sensitivity of test-time accuracy to the TailConfidence threshold c on Qwen-2.5-7B-RL-Augmentation at budget $K = 16$. ‘‘Rounds’’ is the average number of revision rounds before stopping.

Threshold c	AIME25		AIME24		MATH500	
	Acc. (%)	Rounds	Acc. (%)	Rounds	Acc. (%)	Rounds
0.4	16.56	10.22	21.15	9.04	77.80	3.99
0.5	16.88	11.91	21.56	10.06	79.17	4.46
0.6	17.50	13.31	22.08	12.19	79.10	5.64
0.7	18.23	15.45	23.65	14.97	81.58	8.76

630 G Additional Experiments and Ablations

631 We additionally evaluate our method using a strong reasoning-oriented base model, DeepSeek-R1-
 632 Distill-7B (DeepSeek-AI et al., 2025), and compare against both the original base model and an RL
 633 baseline. Due to the substantially longer responses produced by this model (approximately 9k tokens
 634 on average, compared to $\sim 2k$ tokens for Qwen-7B), multi-turn training becomes significantly more
 635 expensive, which limits the scale of these experiments.

636 As shown in Table 7, our RL-Augmented model consistently improves test-case-based performance
 637 over both the base model and the RL baseline. We expect larger gains with increased data scale and
 638 longer training.

639 H Circle Packing Implementation Details

640 This section gives the implementation details behind the circle packing results in Figure 5. Circle
 641 packing maximizes the sum of radii of n non-overlapping circles inscribed in the unit square, with
 642 verifier feedback given by a non-overlap and inscription check together with the realized objective
 643 value. We instantiate REVES on this benchmark by directly optimizing $J_{\phi_{SR}}$ in the recovery form of
 644 Lemma 4.1: each visited state z_t corresponds to a partially-improved candidate configuration, the
 645 policy proposes a refined configuration $y_t \sim \pi_{\theta}(\cdot | z_t)$, and the verifier returns $r^*(x, y_t) = 1$ iff y_t
 646 improves on the running best while remaining feasible.

Table 7: Performance on LiveCodeBench under different evaluation periods.

Method	LiveCodeBench (01/01/25–05/01/25)		LiveCodeBench (08/01/24–01/01/25)	
	OneShot	TestCases-32	OneShot	TestCases-32
Base (Distill-7B)	34.25	44.51	38.24	52.21
Base + RL	36.63	43.96	38.79	53.19
Base + RL-Aug	35.53	46.33	40.63	54.04

647 **Reward shaping.** The raw score of a candidate is the sum of radii itself (a continuous quantity),
 648 which clusters in a narrow band near the optimum (≈ 2.63 for the $n = 26$ instance we evaluate)
 649 and differs across candidates only at the third or fourth decimal. Used directly, this signal is too
 650 compressed to drive learning. We therefore shape it in two steps. First, define the raw score $s(y)$ as

$$s(y) = \begin{cases} (\text{sum of radii of } y), & y \text{ is feasible,} \\ (\text{sum of radii}) - (\text{constraint-violation amount}), & y \text{ is infeasible and } (\text{sum of radii}) \leq 2.64, \\ -1, & y \text{ is infeasible and } (\text{sum of radii}) > 2.64, \end{cases}$$

651 i.e., minor violations are absorbed as a soft penalty, while implausibly high sums coupled with
 652 infeasibility are flagged as cheating and assigned $s = -1$. Second, we amplify by squaring the offset
 653 from 2.5: $R(s) := (s - 2.5)^2$. Since feasible raw scores live above 2.5, this magnifies small absolute
 654 differences, e.g., a $\Delta s = 10^{-3}$ gain near $s = 2.6$ produces a $\Delta R \approx 2 \cdot 10^{-4}$, which is no longer
 655 dwarfed by gradient noise. The reward used for a revision y from parent state x is the resulting
 656 improvement,

$$U(y | x) := R(s(y)) - R(s(x)),$$

657 so the policy is rewarded for any move that improves on the parent and penalized for any move that
 658 worsens it.

659 **Stage I: search rollouts and prompt-set construction.** We run 4 sequential-revision rollouts under
 660 the current π_θ with budget $K = 100$, where each step proposes a refinement of the previous candidate
 661 configuration and the verifier returns the shaped reward $U(y | x)$ defined above. From the visited
 662 intermediate states across these rollouts, we filter to retain those whose subsequent revisions produce
 663 verified improvements, and assemble the resulting revision prompts into the augmented prompt set.

664 **Stage II: RL training.** We train Qwen3-4B with GRPO (DeepSeek-AI et al., 2025) on the prompt
 665 set from Stage I. Each gradient step samples 8 prompts; for each prompt we draw 64 rollouts to form
 666 the group used for the GRPO advantage estimate.

667 **A representative solution.** Figure 9 visualizes one of the configurations REVES discovers, and
 668 Listing 1 reproduces the Python program emitted by the policy that generates it: an initial grid
 669 placement followed by an SLSQP refinement against the non-overlap and inscription constraints.

Listing 1: Python program output by the REVES-trained Qwen3-4B policy that produces the configuration in Figure 9.

```
670 import numpy as np
671 from scipy.optimize import minimize
672
673
674 def run_packing():
675     num_circles = 26
676     r_initial = 0.08
677     x_positions = np.linspace(r_initial, 1 - r_initial, 6)
678     y_positions = np.linspace(r_initial, 1 - r_initial, 5)
679     positions = [(x, y) for x in x_positions for y in y_positions
680                 ][:25]
681     positions.append((0.3, 0.3))
682     positions = np.unique(positions, axis=0)
683     if len(positions) < 26:
684         positions = np.append(positions, [(0.5, 0.5)] * (26 - len(
685             positions)), axis=0)
686
687     radii = np.full(num_circles, r_initial)
688
689     variables = np.zeros(num_circles * 3)
690     for i in range(num_circles):
691         variables[2 * i] = positions[i][0]
692         variables[2 * i + 1] = positions[i][1]
693     for i in range(num_circles):
694         variables[2 * num_circles + i] = radii[i]
695
```

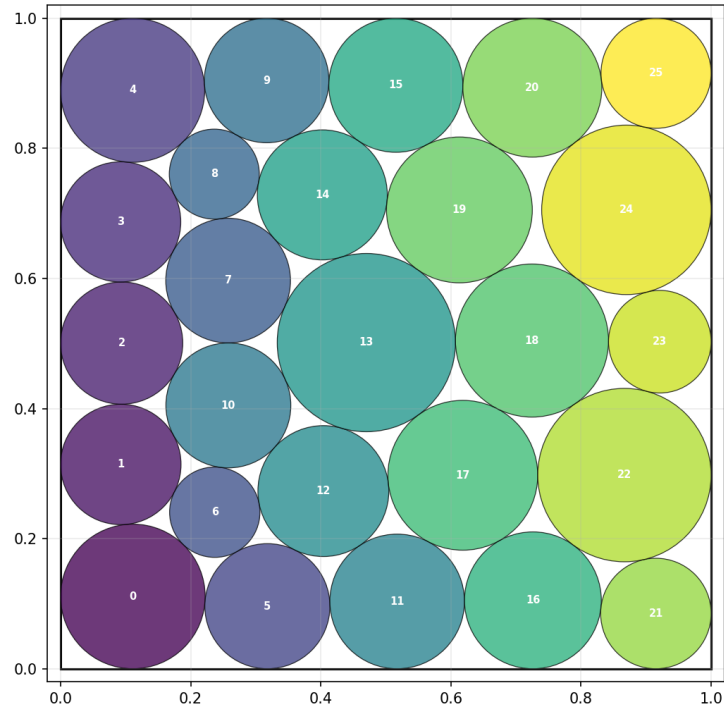


Figure 9: A circle packing configuration found by REVES on the $n = 26$ instance, achieving sum-of-radii = 2.635983.

```

696 def objective(vars):
697     return -np.sum(vars[num_circles * 2:])
698
699 constraints = []
700 for i in range(num_circles):
701     def con1(vars, i=i):
702         centers = vars[:num_circles * 2].reshape(num_circles, 2)
703         r = vars[2 * num_circles + i]
704         return centers[i, 0] - r
705     constraints.append({'type': 'ineq', 'fun': con1})
706
707     def con2(vars, i=i):
708         centers = vars[:num_circles * 2].reshape(num_circles, 2)
709         r = vars[2 * num_circles + i]
710         return 1 - centers[i, 0] - r
711     constraints.append({'type': 'ineq', 'fun': con2})
712
713     def con3(vars, i=i):
714         centers = vars[:num_circles * 2].reshape(num_circles, 2)
715         r = vars[2 * num_circles + i]
716         return centers[i, 1] - r
717     constraints.append({'type': 'ineq', 'fun': con3})
718
719     def con4(vars, i=i):
720         centers = vars[:num_circles * 2].reshape(num_circles, 2)
721         r = vars[2 * num_circles + i]
722         return 1 - centers[i, 1] - r
723     constraints.append({'type': 'ineq', 'fun': con4})
724
725 for i in range(num_circles):
726     def con5(vars, i=i):
727         return vars[2 * num_circles + i]

```

```

728     constraints.append({'type': 'ineq', 'fun': con5})
729
730     for i in range(num_circles):
731         for j in range(i + 1, num_circles):
732             def con6(vars, i=i, j=j):
733                 centers = vars[:num_circles * 2].reshape(num_circles,
734                                                            2)
735                 r_i = vars[2 * num_circles + i]
736                 r_j = vars[2 * num_circles + j]
737                 dx = centers[i, 0] - centers[j, 0]
738                 dy = centers[i, 1] - centers[j, 1]
739                 return dx**2 + dy**2 - (r_i + r_j)**2
740             constraints.append({'type': 'ineq', 'fun': con6})
741
742     result = minimize(
743         objective,
744         variables,
745         constraints=constraints,
746         method='SLSQP',
747         options={'maxiter': 2000000, 'ftol': 1e-10, 'disp': False}
748     )
749
750     centers = result.x[:num_circles * 2].reshape(num_circles, 2)
751     radii = result.x[num_circles * 2:]
752     return centers, radii, float(np.sum(radii))

```

Prompt Template for Code Revision

Problem

You will be given a problem specification and are asked to generate a correct Python program that satisfies the requirements and passes all tests. The task involves extracting two integers and an operator (“gains” or “loses”) from a natural language string and returning the result of the corresponding calculation. ...

Previous Answer

To address the issue in the provided code, there are a few key points ...

```

class Solution:
    def calculate(self, string: str) -> int:
        words = string.split()
        ...
        joined_string = "".join(words)
        result = eval(joined_string)
        return result

```

...

Verifiable Feedback:

The previous answer is incorrect.

Instruction:

Please repair the code accordingly. Enclose your fixed implementation within the specified Python code delimiters. Do not assume access to the correct implementation.

Figure 10: Prompt template for the code revision stage with binary correctness feedback. Partial incorrect code is shown for illustration, with remaining content omitted for brevity.

Prompt Template for Code Verification

Problem

You will be given a problem specification and are asked to generate a correct Python program that satisfies the requirements and passes all tests. The task involves extracting two integers and an operator (“gains” or “loses”) from a natural language string and returning the result of the corresponding calculation. ...

Candidate Answer:

To address the issue in the provided code, there are a few key points ...

```
class Solution:
    def calculate(self, string: str) -> int:
        words = string.split()
        ...
        joined_string = "".join(words)
        result = eval(joined_string)
        return result
```

...

Instruction:

Please verify whether the candidate answer is correct. At the end, output the final line exactly as:
VERDICT: \boxed{True/False}

Figure 11: Prompt template for the code verification stage. The verifier evaluates a candidate solution and outputs a binary correctness verdict.

753 **NeurIPS Paper Checklist**

754 The checklist is designed to encourage best practices for responsible machine learning research,
755 addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove
756 the checklist: **The papers not including the checklist will be desk rejected.** The checklist should
757 follow the references and follow the (optional) supplemental material. The checklist does NOT count
758 towards the page limit.

759 Please read the checklist guidelines carefully for information on how to answer these questions. For
760 each question in the checklist:

- 761 • You should answer [Yes], [No], or [N/A].
- 762 • [N/A] means either that the question is Not Applicable for that particular paper or the
763 relevant information is Not Available.
- 764 • Please provide a short (1–2 sentence) justification right after your answer (even for [N/A]).

765 **The checklist answers are an integral part of your paper submission.** They are visible to the
766 reviewers, area chairs, senior area chairs, and ethics reviewers. You will also be asked to include it
767 (after eventual revisions) with the final version of your paper, and its final version will be published
768 with the paper.

769 The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation.
770 While [Yes] is generally preferable to [No], it is perfectly acceptable to answer [No] provided a
771 proper justification is given (e.g., error bars are not reported because it would be too computationally
772 expensive” or “we were unable to find the license for the dataset we used”). In general, answering
773 [No] or [N/A] is not grounds for rejection. While the questions are phrased in a binary way, we
774 acknowledge that the true answer is often more nuanced, so please just use your best judgment and
775 write a justification to elaborate. All supporting evidence can appear either in the main paper or the
776 supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification
777 please point to the section(s) where related material for the question can be found.

778 **IMPORTANT, please:**

- 779 • **Delete this instruction block, but keep the section heading “NeurIPS Paper Checklist”,**
- 780 • **Keep the checklist subsection headings, questions/answers and guidelines below.**
- 781 • **Do not modify the questions and only use the provided macros for your answers.**

782 **1. Claims**

783 Question: Do the main claims made in the abstract and introduction accurately reflect the
784 paper’s contributions and scope?

785 Answer: [Yes]

786 Justification: The main claims made in the abstract and introduction reflect the paper’s
787 contributions and scope.

788 Guidelines:

- 789 • The answer [N/A] means that the abstract and introduction do not include the claims
790 made in the paper.
- 791 • The abstract and/or introduction should clearly state the claims made, including the
792 contributions made in the paper and important assumptions and limitations. A [No] or
793 [N/A] answer to this question will not be perceived well by the reviewers.
- 794 • The claims made should match theoretical and experimental results, and reflect how
795 much the results can be expected to generalize to other settings.
- 796 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
797 are not attained by the paper.

798 **2. Limitations**

799 Question: Does the paper discuss the limitations of the work performed by the authors?

800 Answer: [Yes]

801 Justification: See Appendix B for discussion on limitations.

802 Guidelines:

- 803 • The answer [N/A] means that the paper has no limitation while the answer [No] means
804 that the paper has limitations, but those are not discussed in the paper.
- 805 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 806 • The paper should point out any strong assumptions and how robust the results are to
807 violations of these assumptions (e.g., independence assumptions, noiseless settings,
808 model well-specification, asymptotic approximations only holding locally). The authors
809 should reflect on how these assumptions might be violated in practice and what the
810 implications would be.
- 811 • The authors should reflect on the scope of the claims made, e.g., if the approach was
812 only tested on a few datasets or with a few runs. In general, empirical results often
813 depend on implicit assumptions, which should be articulated.
- 814 • The authors should reflect on the factors that influence the performance of the approach.
815 For example, a facial recognition algorithm may perform poorly when image resolution
816 is low or images are taken in low lighting. Or a speech-to-text system might not be
817 used reliably to provide closed captions for online lectures because it fails to handle
818 technical jargon.
- 819 • The authors should discuss the computational efficiency of the proposed algorithms
820 and how they scale with dataset size.
- 821 • If applicable, the authors should discuss possible limitations of their approach to
822 address problems of privacy and fairness.
- 823 • While the authors might fear that complete honesty about limitations might be used by
824 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
825 limitations that aren’t acknowledged in the paper. The authors should use their best
826 judgment and recognize that individual actions in favor of transparency play an impor-
827 tant role in developing norms that preserve the integrity of the community. Reviewers
828 will be specifically instructed to not penalize honesty concerning limitations.

829 3. Theory assumptions and proofs

830 Question: For each theoretical result, does the paper provide the full set of assumptions and
831 a complete (and correct) proof?

832 Answer: [Yes]

833 Justification: See Appendix C for the full set of assumptions and proofs.

834 Guidelines:

- 835 • The answer [N/A] means that the paper does not include theoretical results.
- 836 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
837 referenced.
- 838 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 839 • The proofs can either appear in the main paper or the supplemental material, but if
840 they appear in the supplemental material, the authors are encouraged to provide a short
841 proof sketch to provide intuition.
- 842 • Inversely, any informal proof provided in the core of the paper should be complemented
843 by formal proofs provided in appendix or supplemental material.
- 844 • Theorems and Lemmas that the proof relies upon should be properly referenced.

845 4. Experimental result reproducibility

846 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
847 perimental results of the paper to the extent that it affects the main claims and/or conclusions
848 of the paper (regardless of whether the code and data are provided or not)?

849 Answer: [Yes]

850 Justification: See Appendix E for the implementation details.

851 Guidelines:

- 852 • The answer [N/A] means that the paper does not include experiments.

- 853 • If the paper includes experiments, a [No] answer to this question will not be perceived
854 well by the reviewers: Making the paper reproducible is important, regardless of
855 whether the code and data are provided or not.
- 856 • If the contribution is a dataset and/or model, the authors should describe the steps taken
857 to make their results reproducible or verifiable.
- 858 • Depending on the contribution, reproducibility can be accomplished in various ways.
859 For example, if the contribution is a novel architecture, describing the architecture fully
860 might suffice, or if the contribution is a specific model and empirical evaluation, it may
861 be necessary to either make it possible for others to replicate the model with the same
862 dataset, or provide access to the model. In general, releasing code and data is often
863 one good way to accomplish this, but reproducibility can also be provided via detailed
864 instructions for how to replicate the results, access to a hosted model (e.g., in the case
865 of a large language model), releasing of a model checkpoint, or other means that are
866 appropriate to the research performed.
- 867 • While NeurIPS does not require releasing code, the conference does require all submissions
868 to provide some reasonable avenue for reproducibility, which may depend on the
869 nature of the contribution. For example
 - 870 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
871 to reproduce that algorithm.
 - 872 (b) If the contribution is primarily a new model architecture, the paper should describe
873 the architecture clearly and fully.
 - 874 (c) If the contribution is a new model (e.g., a large language model), then there should
875 either be a way to access this model for reproducing the results or a way to reproduce
876 the model (e.g., with an open-source dataset or instructions for how to construct
877 the dataset).
 - 878 (d) We recognize that reproducibility may be tricky in some cases, in which case
879 authors are welcome to describe the particular way they provide for reproducibility.
880 In the case of closed-source models, it may be that access to the model is limited in
881 some way (e.g., to registered users), but it should be possible for other researchers
882 to have some path to reproducing or verifying the results.

883 5. Open access to data and code

884 Question: Does the paper provide open access to the data and code, with sufficient instruc-
885 tions to faithfully reproduce the main experimental results, as described in supplemental
886 material?

887 Answer: [Yes]

888 Justification: I include the code in the supplemental material, data is publicly available.

889 Guidelines:

- 890 • The answer [N/A] means that paper does not include experiments requiring code.
- 891 • Please see the NeurIPS code and data submission guidelines ([https://neurips.cc/
892 public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 893 • While we encourage the release of code and data, we understand that this might not
894 be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not
895 including code, unless this is central to the contribution (e.g., for a new open-source
896 benchmark).
- 897 • The instructions should contain the exact command and environment needed to run to
898 reproduce the results. See the NeurIPS code and data submission guidelines ([https://
899 neurips.cc/public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 900 • The authors should provide instructions on data access and preparation, including how
901 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 902 • The authors should provide scripts to reproduce all experimental results for the new
903 proposed method and baselines. If only a subset of experiments are reproducible, they
904 should state which ones are omitted from the script and why.
- 905 • At submission time, to preserve anonymity, the authors should release anonymized
906 versions (if applicable).

- 907 • Providing as much information as possible in supplemental material (appended to the
908 paper) is recommended, but including URLs to data and code is permitted.

909 6. Experimental setting/details

910 Question: Does the paper specify all the training and test details (e.g., data splits, hyperpa-
911 rameters, how they were chosen, type of optimizer) necessary to understand the results?

912 Answer: [Yes]

913 Justification: See Appendix E and supplemental material for the implementation details.

914 Guidelines:

- 915 • The answer [N/A] means that the paper does not include experiments.
- 916 • The experimental setting should be presented in the core of the paper to a level of detail
917 that is necessary to appreciate the results and make sense of them.
- 918 • The full details can be provided either with the code, in appendix, or as supplemental
919 material.

920 7. Experiment statistical significance

921 Question: Does the paper report error bars suitably and correctly defined or other appropriate
922 information about the statistical significance of the experiments?

923 Answer: [Yes]

924 Justification: We run the experiments multiple times and report the mean.

925 Guidelines:

- 926 • The answer [N/A] means that the paper does not include experiments.
- 927 • The authors should answer [Yes] if the results are accompanied by error bars, confidence
928 intervals, or statistical significance tests, at least for the experiments that support the
929 main claims of the paper.
- 930 • The factors of variability that the error bars are capturing should be clearly stated (for
931 example, train/test split, initialization, random drawing of some parameter, or overall
932 run with given experimental conditions).
- 933 • The method for calculating the error bars should be explained (closed form formula,
934 call to a library function, bootstrap, etc.)
- 935 • The assumptions made should be given (e.g., Normally distributed errors).
- 936 • It should be clear whether the error bar is the standard deviation or the standard error
937 of the mean.
- 938 • It is OK to report 1-sigma error bars, but one should state it. The authors should
939 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
940 of Normality of errors is not verified.
- 941 • For asymmetric distributions, the authors should be careful not to show in tables or
942 figures symmetric error bars that would yield results that are out of range (e.g., negative
943 error rates).
- 944 • If error bars are reported in tables or plots, the authors should explain in the text how
945 they were calculated and reference the corresponding figures or tables in the text.

946 8. Experiments compute resources

947 Question: For each experiment, does the paper provide sufficient information on the com-
948 puter resources (type of compute workers, memory, time of execution) needed to reproduce
949 the experiments?

950 Answer: [Yes]

951 Justification: See Appendix E for the implementation details.

952 Guidelines:

- 953 • The answer [N/A] means that the paper does not include experiments.
- 954 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
955 or cloud provider, including relevant memory and storage.
- 956 • The paper should provide the amount of compute required for each of the individual
957 experimental runs as well as estimate the total compute.

- 958 • The paper should disclose whether the full research project required more compute
959 than the experiments reported in the paper (e.g., preliminary or failed experiments that
960 didn't make it into the paper).

961 **9. Code of ethics**

962 Question: Does the research conducted in the paper conform, in every respect, with the
963 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

964 Answer: [Yes]

965 Justification: The research in the manuscript is compliant with the NeurIPS Code of Ethics.

966 Guidelines:

- 967 • The answer [N/A] means that the authors have not reviewed the NeurIPS Code of
968 Ethics.
969 • If the authors answer [No], they should explain the special circumstances that require a
970 deviation from the Code of Ethics.
971 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
972 eration due to laws or regulations in their jurisdiction).

973 **10. Broader impacts**

974 Question: Does the paper discuss both potential positive societal impacts and negative
975 societal impacts of the work performed?

976 Answer: [Yes]

977 Justification: Appendix B discusses this verifier requirement and natural future directions
978 toward scientific-discovery applications.

979 Guidelines:

- 980 • The answer [N/A] means that there is no societal impact of the work performed.
981 • If the authors answer [N/A] or [No], they should explain why their work has no societal
982 impact or why the paper does not address societal impact.
983 • Examples of negative societal impacts include potential malicious or unintended uses
984 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
985 (e.g., deployment of technologies that could make decisions that unfairly impact specific
986 groups), privacy considerations, and security considerations.
987 • The conference expects that many papers will be foundational research and not tied
988 to particular applications, let alone deployments. However, if there is a direct path to
989 any negative applications, the authors should point it out. For example, it is legitimate
990 to point out that an improvement in the quality of generative models could be used to
991 generate Deepfakes for disinformation. On the other hand, it is not needed to point out
992 that a generic algorithm for optimizing neural networks could enable people to train
993 models that generate Deepfakes faster.
994 • The authors should consider possible harms that could arise when the technology is
995 being used as intended and functioning correctly, harms that could arise when the
996 technology is being used as intended but gives incorrect results, and harms following
997 from (intentional or unintentional) misuse of the technology.
998 • If there are negative societal impacts, the authors could also discuss possible mitigation
999 strategies (e.g., gated release of models, providing defenses in addition to attacks,
1000 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
1001 feedback over time, improving the efficiency and accessibility of ML).

1002 **11. Safeguards**

1003 Question: Does the paper describe safeguards that have been put in place for responsible
1004 release of data or models that have a high risk for misuse (e.g., pre-trained language models,
1005 image generators, or scraped datasets)?

1006 Answer: [N/A]

1007 Justification: We don't release any data or models that have a high risk for misuse.

1008 Guidelines:

- 1009 • The answer [N/A] means that the paper poses no such risks.

- 1010
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- 1011
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- 1012
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.
- 1013
- 1014
- 1015
- 1016
- 1017
- 1018

12. Licenses for existing assets

1019

1020 Question: Are the creators or original owners of assets (e.g., code, data, models), used in

1021 the paper, properly credited and are the license and terms of use explicitly mentioned and

1022 properly respected?

1023 Answer: [Yes]

1024 Justification: We credit the creators of the assets used in the paper.

1025 Guidelines:

- The answer [N/A] means that the paper does not use existing assets.
 - The authors should cite the original paper that produced the code package or dataset.
 - The authors should state which version of the asset is used and, if possible, include a URL.
 - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
 - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
 - If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, `paperswithcode.com/datasets` has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
 - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
 - If this information is not available online, the authors are encouraged to reach out to the asset's creators.
- 1026
- 1027
- 1028
- 1029
- 1030
- 1031
- 1032
- 1033
- 1034
- 1035
- 1036
- 1037
- 1038
- 1039
- 1040

13. New assets

1041

1042 Question: Are new assets introduced in the paper well documented and is the documentation

1043 provided alongside the assets?

1044 Answer: [N/A]

1045 Justification: This paper does not release new assets.

1046 Guidelines:

- The answer [N/A] means that the paper does not release new assets.
 - Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
 - The paper should discuss whether and how consent was obtained from people whose asset is used.
 - At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.
- 1047
- 1048
- 1049
- 1050
- 1051
- 1052
- 1053
- 1054

14. Crowdsourcing and research with human subjects

1055

1056 Question: For crowdsourcing experiments and research with human subjects, does the paper

1057 include the full text of instructions given to participants and screenshots, if applicable, as

1058 well as details about compensation (if any)?

1059 Answer: [N/A]

1060 Justification: This paper does not involve crowdsourcing nor research with human subjects.

1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100

Guidelines:

- The answer [N/A] means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [N/A]

Justification: This paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer [N/A] means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does *not* impact the core methodology, scientific rigor, or originality of the research, declaration is not required.

Answer: [N/A]

Justification: The core method development in this research does not involve LLMs.

Guidelines:

- The answer [N/A] means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy in the NeurIPS handbook for what should or should not be described.