

220 admit closed forms:

$$\mu_q = (1 + \lambda_{\text{neg}})p - \lambda_{\text{neg}}, \quad (5)$$

$$\sigma_q = (1 + \lambda_{\text{neg}})\sqrt{p(1 - p)}. \quad (6)$$

221 Thus, the normalized advantages for correct and incorrect
222 trajectories become

$$A^+ = \frac{(1 + \lambda_{\text{neg}})(1 - p)}{(1 + \lambda_{\text{neg}})\sqrt{p(1 - p)} + \varepsilon_{\text{std}}}, \quad (7)$$

$$A^- = \frac{-(1 + \lambda_{\text{neg}})p}{(1 + \lambda_{\text{neg}})\sqrt{p(1 - p)} + \varepsilon_{\text{std}}}.$$

223 When ε_{std} is small relative to $(1 + \lambda_{\text{neg}})\sqrt{p(1 - p)}$, the
224 geometry is dominated by p : rare successes ($p \ll 1$) receive
225 large positive advantages, while rare failures ($p \approx 1$) receive
226 large negative advantages. This yields an automatic *rare-
227 event amplification* effect driven by group normalization
228 rather than an explicit curriculum.

229 **Example with $G = 8$.** With $G = 8$ and $\varepsilon_{\text{std}} \ll \sqrt{p(1 - p)}$, the induced advantage geometry depends
230 mainly on p : hard groups with $p = 1/8$ yield $A^+ \approx \sqrt{7}$
231 and $A^- \approx -1/\sqrt{7}$, while easy groups with $p = 7/8$ yield
232 $A^+ \approx 1/\sqrt{7}$ and $A^- \approx -\sqrt{7}$. This illustrates the core
233 mechanism: hard prompts amplify rare successes (positive
234 anchors), whereas easy prompts amplify rare failures (negative
235 guidance), producing an adaptive curriculum without
236 explicit difficulty heuristics.

237 3.2. Instance Selection via Positive–Negative Pairing

238 **Core idea.** Instead of selecting training prompts solely
239 by a single scalar heuristic (e.g., historical-accuracy variance), we explicitly construct a *bidirectional* minibatch
240 consisting of (i) one prompt that yields a stable positive
241 anchor and (ii) one prompt that yields a stable negative
242 warning. Concretely, we select a two-example training set
243 $\mathcal{D}_{\pm} = \{q^+, q^-\}$, where q^+ is *hard-but-solvable* (rare
244 successes exist) and q^- is *easy-but-brittle* (rare failures exist).
245 Under WGRPO, these two regimes map directly to amplified
246 tail-event teaching signals (Sec. 3.1).

247 **Positive anchor: hard-but-solvable.** We choose q^+ such
248 that the current policy achieves a low but non-zero success
249 rate:

$$p(q^+) \in \left[\frac{1}{G}, \frac{c}{G} \right], \quad (8)$$

250 so that $0 < k < G$ and each group typically contains a
251 small number of correct rollouts. In this regime, WGRPO
252 assigns large positive advantages to rare correct trajectories,
253 concentrating updates on demonstrations of what the model
254 should do.

255 **Negative guidance: easy-but-brittle.** We choose q^- such
256 that the current policy achieves a high but not perfect success
257 rate:

$$p(q^-) \in \left[1 - \frac{c}{G}, 1 - \frac{1}{G} \right], \quad (9)$$

258 so that failures are rare but still occur. In this regime,
259 WGRPO assigns large-magnitude negative advantages to
260 rare failures, producing a sharp “do-not” signal that sup-
261 presses high-confidence failure modes while preserving al-
262 ternative plausible solutions under the model prior.

263 **Practical selection via lightweight probing.** To instanti-
264 ate positive–negative pairing with only two training prompts,
265 we perform a simple probing stage on two candidate pools
266 with different expected difficulty under the same base model.
267 We use an “easy” candidate pool \mathcal{C}^- and a “hard” candidate
268 pool \mathcal{C}^+ ; in our experiments \mathcal{C}^- is drawn from DeepScaleR-
269 sub and \mathcal{C}^+ is drawn from AIME 2025, but the procedure
270 is agnostic to the specific sources. For each candidate
271 $q \in \mathcal{C}^+ \cup \mathcal{C}^-$, we estimate its success rate under the current
272 policy by sampling M independent groups of size G and
273 averaging:

$$\bar{p}(q) = \frac{1}{M} \sum_{m=1}^M \hat{p}_m(q).$$

274 To ensure non-degenerate within-group variance, we discard
275 candidates with $\bar{p}(q) \notin [\delta, 1 - \delta]$, where we use $\delta = 1/G$ by
276 default. We then select one positive anchor and one negative
277 guidance prompt by targeting the two WGRPO regimes:

$$q^+ = \arg \min_{q \in \mathcal{C}^+} |\bar{p}(q) - p_{\text{hard}}|, \quad (10)$$

$$q^- = \arg \min_{q \in \mathcal{C}^-} |\bar{p}(q) - p_{\text{easy}}|,$$

278 where $p_{\text{hard}} \approx 1/G$ and $p_{\text{easy}} \approx 1 - 1/G$. This ensures
279 that q^+ operates in a low-but-nonzero success regime that
280 amplifies rare successes, while q^- operates in a high-but-
281 not-perfect success regime that amplifies rare failures. Over-
282 all, the selection is deliberately simple, uses only on-policy
283 probing (no historical training statistics), and directly instantiates
284 the rare-event amplification mechanism of WGRPO
285 with only two training examples.

4. Experimental Setup

286 **Models.** To study how different training-example selec-
287 tion strategies affect RLVR, we run our training pipeline
288 on several representative open-weight LLMs from differ-
289 ent families and scales. In particular, we train QWEN2.5-
290 MATH-7B, QWEN2.5-MATH-7B-INSTRUCT, and LLAMA-
291 3.1-8B-INSTRUCT.

292 **Training dataset.** The training examples we select come
293 from AIME 2025 (Art of Problem Solving, 2025a) and