
Cost-Aware Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We consider the problem of Cost-Aware Learning, where sampling different components of a finite-sum objective incurs different costs. The objective is to reach a target error while minimizing the total cost. We propose Cost-Aware SGD, which uses a distribution based on gradient norms and costs to sample components. We provide a thorough analysis of this algorithm, including cost-improvement bounds over baselines, a characterization of distribution proxy sub-optimality, and a lower bound. We apply our theoretical insights to reinforcement learning with language models, where the computational cost of sequence-level policy gradients varies with length. We find that the advantage magnitude serves as a high-fidelity proxy for gradient norms, and use this to introduce Cost-Aware GRPO. Empirical results on 1.5B, 4B, and 8B LLMs demonstrate that this algorithm significantly reduces the tokens used in policy optimization while matching or exceeding baseline accuracy.

13 1 Introduction

14 Stochastic gradient descent (SGD) [Nemirovsky et al., 1983] theory typically measures convergence in terms of the number of gradient steps required to reach a target error. This implicitly assumes that gradient evaluation on any data point incurs an identical computational cost. However, in modern language model training, the assumption of uniform cost is not always true. For example, in reinforcement learning for language models [Ouyang et al., 2022, Shao et al., 2024], the FLOPs required to compute a policy gradient scale linearly with sequence length, meaning a long reasoning trace can cost an order of magnitude more than a concise response.

21 Motivated by this computational heterogeneity, we study the problem of *Cost-Aware Learning*: given a finite-sum objective where each component has a known evaluation cost c_i , how should one sample gradients to reach ϵ -error at a minimum total cost? We propose Cost-Aware SGD, which leverages importance sampling [Kloek and Van Dijk, 1978, Beygelzimer et al., 2009] to derive the optimal sampling distribution $p_i^* \propto G_i/\sqrt{c_i}$, where G_i is the Lipschitz constant (gradient norm bound) of component i .

27 We analyze the cost-improvement over uniform and variance-reduction baselines, and validate this in a synthetic setting. Since exact gradient norms are often unavailable in practice, we characterize the sub-optimality of using a proxy distribution. We establish an information-theoretic lower bound for cost-aware learning, which motivates a subset selection algorithm. Finally, we show that the derived sampling strategy extends beyond convex objectives.

32 We apply these insights to the policy optimization phase of GRPO [Shao et al., 2024], where sequence length provides a natural cost. A key empirical finding, motivated by the form of the policy gradient, is that the advantage magnitude $|A_i|$ serves as a high-fidelity proxy for G_i based on our sub-optimality metrics. Equipped with this proxy, we introduce Cost-Aware GRPO, which samples model-generated sequences according to p^* with $|A_i|$ in place of G_i during policy optimization. Cost-Aware GRPO significantly reduces the FLOPs used in policy optimization while matching or exceeding baseline

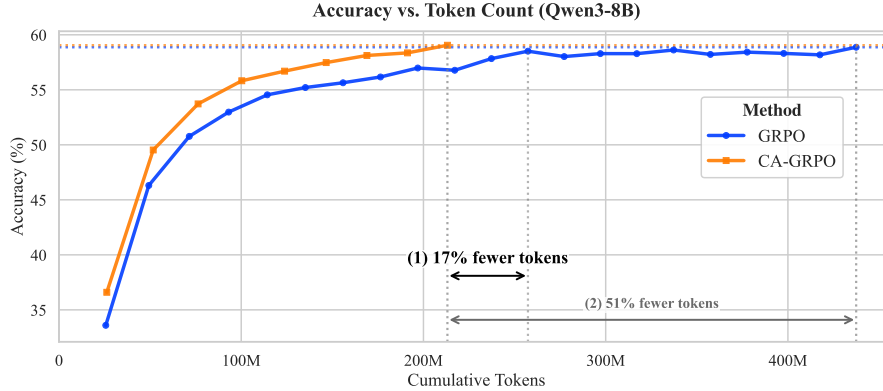


Figure 1: **Qwen3-8B training with GRPO and Cost-Aware GRPO (CA-GRPO).** We evaluate on AIME1983-2024 (pass@1/mean@32) and report cumulative tokens used in policy gradient computation. We plot every 20 training steps until peak accuracy is attained for both. CA-GRPO requires: (1) 17% fewer tokens to reach its absolute peak than GRPO requires to come within 2% of that performance, (2) 51% fewer tokens to match or exceed the peak accuracy of GRPO.

38 performance. Through extensive ablations, we demonstrate robustness to proxy noise, hyperparameter
 39 choices, and show that cost-aware sampling generalizes beyond GRPO to the CISPO objective [Chen
 40 et al., 2025].

41 The main contributions of this paper are as follows:

- 42 1. **Cost-Aware SGD (§4).** We derive the optimal cost-aware sampling distribution for convex
 43 and strongly convex finite-sum optimization, analyze its cost-improvement over baselines,
 44 characterize the sub-optimality of using a proxy distribution, and establish an information-
 45 theoretic lower bound.
- 46 2. **Cost-Aware GRPO (§5).** We identify $|A_i|$ as a high-fidelity proxy for gradient norms and
 47 apply cost-aware sampling to GRPO. Experiments with Qwen2.5-Math-1.5B-Instruct [Yang
 48 et al., 2024], Qwen3-4B Base and Qwen3-8B Base [Yang et al., 2025] demonstrate consistent
 49 improvements in token efficiency without compromising on performance.

50 2 Related Work

51 We describe the primary related work, and defer the discussion of additional works to Appendix A.

52 **Importance sampling for SGD.** A significant body of work focuses on accelerating SGD conver-
 53 gence via importance sampling, primarily through variance reduction [Zhao and Zhang, 2015, Needell
 54 et al., 2014]. A foundational result by Alain et al. [2015] established that the optimal sampling
 55 distribution for variance reduction assigns probabilities proportional to the per-sample gradient norms.
 56 In the context of deep learning, where exact gradient norms are expensive to compute, methods have
 57 been proposed to approximate these norms using upper bounds or loss values [Katharopoulos and
 58 Fleuret, 2018, Johnson and Guestrin, 2018]. However, these approaches operate under a standard
 59 optimization framework that treats iterations as the unit of resource consumption. They aim to maxi-
 60 mize the information gain *per step*, implicitly assuming uniform computational cost across samples.
 61 Our work generalizes this by considering the *cost-aware* setting where the aim is to minimize training
 62 cost, leading to a sampling distribution that balances gradient magnitude against sample cost.

63 In particular, the work of Chen et al. [2023b] shares similarities with our work by also studying
 64 a finite-sum setting with varying sample costs. However, our contributions differ significantly
 65 in analysis and scope: while their work focuses on asymptotic behavior, we provide finite-time
 66 guarantees complemented by a lower bound. Furthermore, we implement this cost-aware framework
 67 with language models via Cost-Aware GRPO.

68 **RL for Language Models.** While policy gradient methods [Sutton et al., 1999, Kakade, 2001,
69 Schulman et al., 2015] have been used in many settings, the works of Ouyang et al. [2022], Dai
70 et al. [2023] have garnered significant attention for their application to LLM post-training. Various
71 methods such as DPO [Rafailov et al., 2023], PPO [Schulman et al., 2017, Zheng et al., 2023], and
72 GRPO [Shao et al., 2024] have been shown to be extremely effective in improving accuracy on
73 standard benchmarks. However, these methods incur substantial computational and memory costs.
74 Several methods aim to improve data efficiency: [Wang et al., 2025, Fatemi et al., 2025, Li et al.,
75 2025, Yu et al., 2025, Xu et al., 2025, Zheng et al., 2025]. Yu et al. [2025] introduces, among other
76 contributions, a mechanism to continue resampling from a model until a sufficient amount of non-zero
77 advantage responses have been generated, showing that this data helps the model achieve better
78 performance. Xu et al. [2025] attempts to subsample the generations of a model in an effort to train
79 on less data. Zheng et al. [2025] keeps track of the history of various prompts in an effort to remove
80 low-signal prompts. Our work contrasts with these by using importance reweighting and taking into
81 account the computational cost of sequences when making decisions.

82 3 Problem Formulation

83 We consider minimizing a convex function f in the finite-sum setting where $f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x)$. We
84 assume that f is minimized at x^* , and that each component f_i is a convex and G_i -Lipschitz function.
85 In the cost-aware setting, each component i has an evaluation cost $c_i \geq 0$ for querying its gradient
86 oracle. Namely, upon querying a particular sample i at a point x , the algorithm obtains the exact
87 gradient $\nabla f_i(x)$ and incurs cost c_i . We assume that the learning algorithm has access to G_i and c_i for
88 all samples i .

89 For a given $\epsilon > 0$, the objective of a (possibly randomized) cost-aware learning algorithm is to incur
90 the minimum possible cost to obtain \hat{x} such that $\mathbb{E}[f(\hat{x}) - f(x^*)] \leq \epsilon$.

91 4 Cost-Aware SGD

92 In this section, we address the problem presented in Section 3. We propose Cost-Aware SGD, a
93 simple variant of standard SGD, which differs only in the use of a non-uniform sampling distribution
94 and importance weighting to account for bias. Pseudocode is provided in Appendix 1. While the
95 analysis in this section assumes convex loss functions in order to provide meaningful insights, in
96 Appendix F we show that the optimal sampling distribution derived does not strictly rely on convexity.

97 4.1 Cost Analysis

98 Consider a fixed sampling distribution $\mathbf{p} \in \Delta_n$ used in the Cost-Aware SGD algorithm. The total
99 cost of the algorithm over T iterations is the random variable $\mathcal{K}_T = \sum_{t=1}^T c_{i_t}$. By the linearity of
100 expectation, the expected total cost is:

$$\mathbb{E}[\mathcal{K}_T] = \sum_{t=1}^T \mathbb{E}_{i_t \sim \mathbf{p}} [c_{i_t}] = T \sum_{i=1}^n p_i c_i = T \cdot C(\mathbf{p}), \quad (1)$$

101 where $C(\mathbf{p}) = \sum_{i=1}^n p_i c_i$ represents the expected cost per iteration. Our goal is to minimize the total
102 expected cost required to reach a target error ϵ . Let $T(\epsilon, \mathbf{p})$ denote the number of iterations required
103 to achieve expected error ϵ using distribution \mathbf{p} . We seek to find \mathbf{p} which minimizes:

$$\mathcal{K}(\epsilon, \mathbf{p}) = T(\epsilon, \mathbf{p}) \cdot C(\mathbf{p}). \quad (2)$$

104 **Case 1: General Convex Functions** For general convex functions, standard convergence results
105 for SGD [Nemirovsky et al., 1983, Bubeck et al., 2015] state that with step size $\eta \propto 1/\sqrt{T}$, the
106 expected suboptimality satisfies: $\mathbb{E}[f(\bar{x}_T) - f(x^*)] \leq \frac{D\sqrt{S(\mathbf{p})}}{\sqrt{T}}$, where $S(\mathbf{p}) = \mathbb{E}_{i \sim \mathbf{p}}[\|\tilde{g}_i\|^2]$ is the
107 second moment of the gradient estimator, and D is the diameter of the parameter space. Expanding
108 the second moment for importance sampling: $S(\mathbf{p}) = \sum_{i=1}^n p_i \left\| \frac{\nabla f_i(x)}{np_i} \right\|^2 \leq \frac{1}{n^2} \sum_{i=1}^n \frac{G_i^2}{p_i}$, where G_i is
109 the Lipschitz constant for component f_i . To guarantee error ϵ , we require $T = \frac{D^2 S(\mathbf{p})}{\epsilon^2}$. Substituting

110 this into the cost formula:

$$\mathcal{K}(\epsilon, \mathbf{p}) = \left(\frac{D^2 S(\mathbf{p})}{\epsilon^2} \right) C(\mathbf{p}) \propto S(\mathbf{p}) C(\mathbf{p}). \quad (3)$$

111 **Case 2: Strongly Convex Functions** For f that is μ -strongly convex, SGD with decaying step
 112 size $\eta_t = 1/(\mu t)$ achieves a rate of $\mathcal{O}(S(\mathbf{p})/\mu T)$ [Rakhlin et al., 2011]. The expected total cost under
 113 distribution \mathbf{p} to attain error ϵ becomes:

$$\mathcal{K}(\epsilon, \mathbf{p}) = \left(\frac{4S(\mathbf{p})}{\mu\epsilon} \right) C(\mathbf{p}) \propto S(\mathbf{p}) C(\mathbf{p}). \quad (4)$$

114 In both cases, minimizing the total cost is equivalent to minimizing the product $J(\mathbf{p}) = S(\mathbf{p})C(\mathbf{p})$.

115 4.2 Derivation of the Optimal Distribution

116 **Theorem 4.1** (Optimal Cost-Weighted Distribution). *The sampling distribution \mathbf{p}^* that minimizes*
 117 $\mathcal{K}(\epsilon, \mathbf{p})$ is:

$$p_i^* = \frac{G_i/\sqrt{c_i}}{\sum_{j=1}^n G_j/\sqrt{c_j}}. \quad (5)$$

118 Under this distribution, the expected costs for general convex and μ -strongly convex functions are:

$$\mathcal{K}_{\text{cvx}}(\epsilon, \mathbf{p}^*) = \frac{D^2}{\epsilon^2} \left(\frac{1}{n} \sum_{i=1}^n G_i \sqrt{c_i} \right)^2, \quad \mathcal{K}_{\text{str-cvx}}(\epsilon, \mathbf{p}^*) = \frac{4}{\mu\epsilon} \left(\frac{1}{n} \sum_{i=1}^n G_i \sqrt{c_i} \right)^2. \quad (6)$$

119 *Proof.* We minimize the product term $J(\mathbf{p}) = S(\mathbf{p})C(\mathbf{p})$. By the Cauchy-Schwarz inequality, for
 120 vectors $u, v \in \mathbb{R}^n$ defined by $u_i = \sqrt{p_i c_i}$ and $v_i = G_i/\sqrt{p_i}$:

$$\left(\sum_{i=1}^n p_i c_i \right) \left(\sum_{i=1}^n \frac{G_i^2}{p_i} \right) = \|u\|^2 \|v\|^2 \geq \langle u, v \rangle^2 = \left(\sum_{i=1}^n G_i \sqrt{c_i} \right)^2. \quad (7)$$

121 Equality holds if and only if $u \propto v$, which implies $\sqrt{p_i c_i} \propto G_i/\sqrt{p_i}$, or $p_i \propto G_i/\sqrt{c_i}$. Substituting
 122 the minimum value of the product back into $J(\mathbf{p})$, we obtain $J(\mathbf{p}^*) = \frac{1}{n^2} (\sum G_i \sqrt{c_i})^2$. Substituting
 123 $J(\mathbf{p}^*)$ into Eq. 3 and Eq. 4 yields the stated costs. \square

124 4.3 Cost-Improvement over Baselines

125 We compare \mathbf{p}^* against two baselines. **Uniform sampling** sets $(p_{\text{unif}})_i = 1/n$. The **variance strategy**
 126 ignores costs and samples proportionally to gradient norms: $(p_{\text{var}})_i = G_i/\sum_j G_j$, which minimizes
 127 the variance of the gradient estimator [Alain et al., 2015, Zhao and Zhang, 2015, Needell et al., 2014]
 128 (see Appendix D.1 for a derivation).

129 Their expected costs to achieve ϵ error under general convex losses are: $\mathcal{K}(\epsilon, p_{\text{unif}}) =$
 130 $\frac{D^2}{\epsilon^2 n^2} (\sum_{i=1}^n G_i^2) (\sum_{i=1}^n c_i)$ and $\mathcal{K}(\epsilon, p_{\text{var}}) = \frac{D^2}{\epsilon^2 n^2} (\sum_{i=1}^n G_i) (\sum_{i=1}^n G_i c_i)$. The cost ratios relative to
 131 optimal sampling are:

$$\frac{\mathcal{K}(\epsilon, \mathbf{p}^*)}{\mathcal{K}(\epsilon, p_{\text{unif}})} = \frac{(\sum G_i \sqrt{c_i})^2}{(\sum G_i^2)(\sum c_i)} \leq 1, \quad \frac{\mathcal{K}(\epsilon, \mathbf{p}^*)}{\mathcal{K}(\epsilon, p_{\text{var}})} = \frac{(\sum G_i \sqrt{c_i})^2}{(\sum G_i)(\sum_i G_i c_i)} \leq 1. \quad (8)$$

132 When all costs are equal, $\mathcal{K}_{\text{opt}} = \mathcal{K}_{\text{var}}$, but \mathcal{K}_{opt} still improves over $\mathcal{K}_{\text{unif}}$ by minimizing variance (see
 133 Appendix D.1 for a detailed comparison). When costs and gradients are negatively correlated, the
 134 gains can be dramatic: setting $G_i = 2^{-i/4}$, $c_i = 2^i$ yields $\mathcal{K}_{\text{opt}}/\mathcal{K}_{\text{unif}} = \Theta(2^{-n/2})$ and $\mathcal{K}_{\text{opt}}/\mathcal{K}_{\text{var}} =$
 135 $\Theta(2^{-n/4})$.

136 4.4 Synthetic Validation

137 We evaluate the sampling strategies on a linear least squares task with $N = 3000$ points in dimension
 138 $d = 50$. The data vectors a_i are scaled such that their norms are bounded by $L = 10$. This geometry

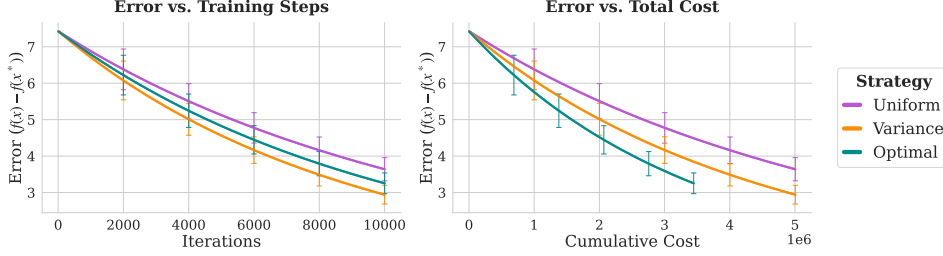


Figure 2: **Synthetic validation.** We compare the error with the total training steps and the total cost incurred by each sampling strategy.

139 ensures that the Lipschitz constants G_i are bounded but heterogeneous across the dataset. Evaluation
 140 costs c_i are assigned randomly from a uniform distribution between 1 and 1000. We compare three
 141 sampling approaches: *Uniform*, *Variance* ($p_i \propto G_i$), and *Optimal* ($p_i \propto G_i/\sqrt{c_i}$). The results are
 142 shown in Figure 2. We average over 1000 trials and report standard deviation error bars.

143 **Cost Efficiency:** For a fixed error target, the *Optimal* strategy incurs the least total cost. By weighting
 144 samples by the ratio of information to cost, it effectively identifies samples that provide significant
 145 variance reduction with a low cost. **Convergence Rate:** The *Variance* strategy converges the fastest
 146 in terms of iteration count. This is expected, as it minimizes the gradient estimator’s variance without
 147 regard for computational budget.

148 4.5 Proxy Sub-Optimality

149 In practice, it can be prohibitively costly to compute G_i . As such, one may not have access to the
 150 perfect sampling distribution p^* ; in fact, this is the case for later experiments in this paper. In this
 151 subsection, we consider the sub-optimality of using a proxy. We recall that $J(p) = S(p)C(p)$, and
 152 that $\mathcal{K}(\epsilon, p)$ is the expected cost to attain error ϵ under distribution p . In Section 4.1, we showed that
 153 for any ϵ , $\mathcal{K}(\epsilon, p) \propto J(p)$. $J(p')$ represents how cost scales with error for p' . Hence, the quantity
 154 $J(p')/J(p^*)$ represents the ratio between the cost incurred by p' and that of p^* .

155 The χ^2 -divergence is defined as follows for two distributions P, Q : $D_{\chi^2}(P||Q) = \sum \frac{P^2}{Q} - 1$. In
 156 Theorem 4.2, we exactly characterize the sub-optimality due to using a different sampling distribution
 157 p' in place of p^* .

158 **Theorem 4.2** (Sub-optimality gap for distribution proxy). *Let p' be any empirical sampling dis-*
 159 *tribution. We define the cost-biased distribution \tilde{p} corresponding to p as: $\tilde{p}_i = (p_i c_i)/C(p)$. The*
 160 *sub-optimality gap is: $J(p')/J(p^*) = 1 + D_{\chi^2}(\tilde{p}^*||\tilde{p}')$.*

161 To provide further intuition as to when one can use some proxy G'_i in place of G_i for a proxy sampling
 162 distribution, we provide Theorem D.5 (Appendix D.5). It implies that, when G_i and G'_i have near-
 163 1 Pearson correlation, $\mathbb{E}[J(p')] \approx J(p^*)$. The Pearson correlation between random variables X
 164 and Y with expected values μ_X, μ_Y , and standard deviations σ_X, σ_Y , respectively, is defined as:
 165 $\rho_{X,Y} = \frac{\text{Cov}(X,Y)}{\sigma_X \sigma_Y}$. Both proofs are deferred to Appendix D.5.

166 4.6 Lower Bound

167 Having established the upper bounds for Cost-Aware SGD, a natural question is whether it is possible
 168 to design a sampling strategy that achieves a lower total cost. In this section, we provide a lower
 169 bound for the class of convex functions.

170 **Theorem 4.3.** *Fix $\epsilon, G > 0$. For any number of components $n \geq G^2/\epsilon^2$, any set of nonnegative*
 171 *query costs $\{c_i\}_{i=1}^n$, and any (possibly randomized) algorithm with an expected error guarantee*
 172 *$\mathbb{E}[F(\hat{x}) - \min F(x)] \leq \epsilon$, there exists a convex, G -Lipschitz (i.e., with $G_i = G$ for all i) finite-sum*
 173 *problem instance over $X = [-1, 1]$ such that the algorithm’s expected total query cost is at least:*

$$\Omega\left(\frac{G^2}{\epsilon^2} \left(\frac{1}{n} \sum_{i \in S^*} \sqrt{c_i}\right)^2\right),$$

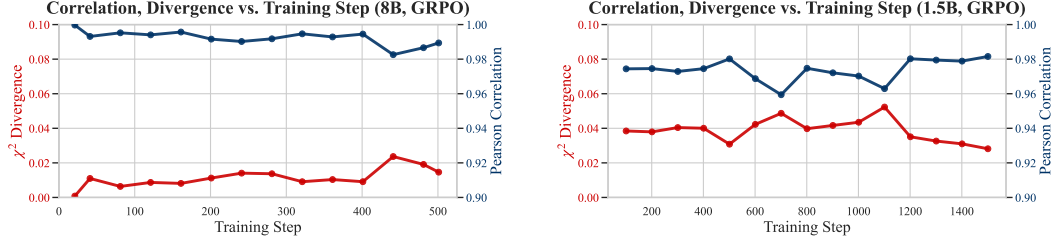


Figure 3: **Sub-optimality metrics** for using $|A_i|$ in place of G_i for GRPO experiments across two model sizes. The Pearson correlations are near 1. The cost-biased χ^2 -divergence between the true p^* and the distribution defined by this proxy is near 0.

174 where $S^* \subseteq [n]$ is any set of components such that following cost-uniformity condition holds:
 175 $(\max_{i \in S^*} \sqrt{c_i}) / (\frac{1}{n} \sum_{j \in S^*} \sqrt{c_j}) \leq \frac{G}{40\epsilon}$.

176 We defer the proof to Appendix E. The lower bound matches the upper bound from Theorem 4.1 in
 177 its $\frac{1}{n^2 \epsilon^2}$ dependence, but differs in how it aggregates costs: the upper bound scales with $(\sum_{i=1}^n \sqrt{c_i})^2$,
 178 while the lower bound scales with $(\sum_{i \in S^*} \sqrt{c_i})^2$. For uniform costs, $S^* = [n]$ and the bounds match
 179 up to constants. When costs are highly non-uniform, the gap arises because expensive components are
 180 excluded from S^* . This motivates combining importance-weighted sampling with *subset selection*:
 181 first choosing which components to include, then optimizing sampling probabilities over the selected
 182 subset. We develop this approach in Appendix C, where we show that the subset selection problem
 183 reduces to a min-cost knapsack problem, for which a polynomial-time 2-approximation exists that
 184 simply selects the cheapest components.

185 5 Cost-Aware GRPO

186 We apply our theoretical analysis to the empirical setting of GRPO [Shao et al., 2024]. In particular,
 187 we focus on *reducing the number of FLOPs used by policy optimization*.

188 5.1 Datasets and Models

189 We perform experiments with the Qwen2.5-Math-1.5B-Instruct [Yang et al., 2024], Qwen3-4B and
 190 Qwen3-8B [Yang et al., 2025] models. We use the DAPO dataset [Yu et al., 2025] for training,
 191 designed for improving accuracy on AIME problems. We evaluate on the following benchmarks:
 192 AIME1983-2024 [Veeraboina, 2023], AMC [Art of Problem Solving, n.d.], MATH500 [Hendrycks
 193 et al., 2021], and GSM8K [Cobbe et al., 2021]. All experiments are performed with the Verl
 194 framework [Sheng et al., 2025], and we make only slight modifications for the purpose of our
 195 algorithm. Hyperparameters for all experiments in Appendix G.

196 5.2 Cost-Aware Learning with GRPO

197 **Background on GRPO.** In GRPO, training alternates between two stages. The first stage generates
 198 many responses for each prompt in order to build a dataset for the second stage. The second stage is
 199 where policy gradient updates are made. The gradient updates are made on the current set of prompts
 200 and responses, which can be considered as the current dataset.

201 We assume the training dataset has “verifiable rewards,” such as a math dataset where the answers
 202 are known. The reward is 1 if the response is correct and 0 otherwise. Given M responses for a
 203 prompt, advantages are normalized: $A_i = (r_i - \frac{1}{M} \sum_{k=1}^M r_k) / (\text{std}(\{r_k\}_{k=1}^M))$. One of the key benefits
 204 of GRPO-style algorithms is that the advantage is very cheap to obtain.

205 **Defining c_i .** The notion of cost fits seamlessly in this setting of GRPO because different prompt
 206 and response pairs have different lengths. The amount of FLOPs required in the policy gradient
 207 step depends linearly on the total number of tokens in the prompt and response. Hence, for each

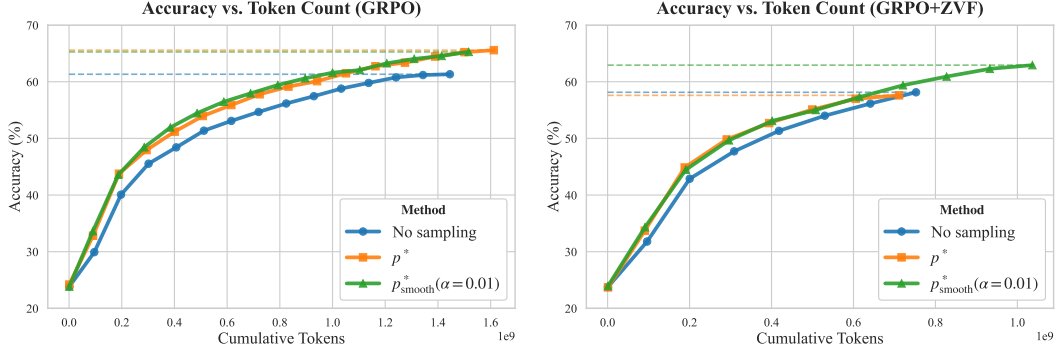


Figure 4: **Accuracy vs. token count** on AIME throughout training for both GRPO and GRPO+ZVF settings for the 1.5B model. We plot the accuracy on the y-axis and the cumulative number of tokens used in policy optimization on the x-axis, to compare the number of tokens used for a fixed accuracy. We evaluate every 100 steps.

208 prompt and response pair, we define c_i to be the sum of the length of the prompt and the length of the
 209 response.

210 **Proxy for G_i .** It remains to find a strong proxy for G_i ; recall that G_i represents the norm of the
 211 gradient of a sample. In language model training, it would not be logical to compute the norm of the
 212 gradient of each sample in order to obtain an optimal sampling distribution, as this has a very high
 213 computational burden and defeats the purpose. Crucially, we estimate G_i with $|A_i|$, the magnitude of
 214 the advantage of sample i . This is motivated by the expression for the policy gradient, in which the
 215 magnitude of the advantage appears to play a large role.

216 We justify this choice by observing low cost-biased χ^2 -divergence between the sampling proxy
 217 obtained this way and the true p^* distribution (invoking Theorem 4.2), as well as high Pearson
 218 correlation between G_i and $|A_i|$ (invoking Theorem D.5). In Figure 3, we report the Pearson
 219 correlation between $|A_i|$ and G_i across training for GRPO with the 8B and 1.5B models. In order
 220 to obtain G_i , we compute true sequence-level gradients. Upon computing the true gradient norms,
 221 we compute the cost-biased χ^2 divergence between the distributions obtained using $|A_i|$ and G_i
 222 respectively. We report sub-optimality metrics for all training settings in Appendix G.8.

223 **Algorithm.** Equipped with notions of cost and a gradient proxy, we combine cost-aware learning
 224 with GRPO in the policy-gradient stage of training. When we perform the policy gradient steps, we
 225 fill mini-batches according to the p^* distribution rather than simply using all samples once. The p^*
 226 distribution is obtained by letting $p_i^* \propto |A_i|/\sqrt{c_i}$. Thus, computing p^* simply requires obtaining the
 227 advantage for all samples. We use importance sampling to mitigate bias. Aside from this, we make
 228 no changes to the GRPO pipeline. The exact procedure is detailed in Algorithm 2.

229 5.3 Settings and Methods

230 We evaluate combinations of two *training settings* and three *sampling methods*.

231 **Training settings.** (1) **GRPO:** Standard GRPO training as described in Section 5.2. (2) **GRPO +**
 232 **ZVF:** GRPO with zero-variance filtering (ZVF) [Yu et al., 2025], which removes zero-advantage
 233 prompts. To maintain a fixed batch size, rollouts are resampled until sufficient non-zero-advantage
 234 samples are obtained.

235 **Sampling methods.** (1) **No sampling:** Standard training on all generated rollouts without resampling
 236 or importance weighting. (2) **p^* sampling:** We sample rollouts according to p^* (Section 5.2,
 237 Algorithm 2) and apply importance weighting to correct for bias. (3) **Smoothed p^* sampling:** We
 238 also consider $p_{\text{smooth}}^*(\alpha) = (1 - \alpha)p^* + \alpha\mathcal{U}$, where \mathcal{U} is the uniform distribution.

239 For the 1.5B model, we evaluate all settings and sampling methods. For the 4B and 8B models, we
 240 restrict our attention to the GRPO setting because it outperforms GRPO+ZVF in the smaller scale
 241 setting. Implementation and hyperparameter details are provided in Appendix G.

Model	Setting	Method	AIME	AMC	MATH	GSM8K	Avg. Accuracy
Math-1.5B-IT	GRPO	No sampling	61.3	64.1	73.2	86.2	71.2
	GRPO	p^*	65.6	71.2	73.2	86.0	74.0
	GRPO	$p^*_{\text{smooth}}(\alpha = 0.01)$	65.3	68.1	72.3	86.0	72.9
	ZVF	No sampling	58.1	64.7	73.2	86.1	70.6
	ZVF	p^*	57.6	63.0	72.9	85.9	69.8
	ZVF	$p^*_{\text{smooth}}(\alpha = 0.01)$	62.9	68.0	73.3	85.8	72.5
4B-Base	GRPO	No sampling	56.4	72.1	86.5	94.5	77.4
	GRPO	p^*	56.7	74.5	86.3	94.2	77.9
	GRPO	$p^*_{\text{smooth}}(\alpha = 0.01)$	56.9	72.0	86.8	94.2	77.4
8B-Base	GRPO	No sampling	58.3	72.8	87.0	96.0	78.5
	GRPO	p^*	58.5	74.7	86.7	95.2	78.8
	GRPO	$p^*_{\text{smooth}}(\alpha = 0.01)$	58.5	74.0	86.7	94.9	78.5

Table 1: **Cost-aware sampling preserves model performance.** For each setting and method, we report the best averaged checkpoint accuracy across four benchmarks. p^* sampling and variants consistently preserve overall model performance.

242 5.4 Results

243 We evaluate all methods across two primary axes: benchmark accuracy and the cumulative tokens
 244 used in policy optimization. Since token count is directly proportional to FLOPs in the policy
 245 gradient step, this measures the cost-accuracy tradeoff of each strategy. All reported accuracies
 246 denote pass@1/mean@32, the average correctness of querying the model 32 times. In all figures, we
 247 plot up to peak accuracy. We focus on AIME accuracy, as this is the intended downstream benchmark
 248 for the DAPO training set.

249 **Qwen2.5-Math-1.5B-Instruct Results.** Figure 4 plots the cumulative policy optimization tokens
 250 against AIME accuracy. Under standard GRPO, both p^* and $p^*_{\text{smooth}}(\alpha = 0.01)$ match the baseline’s
 251 peak accuracy using 28% and 32% fewer tokens respectively and ultimately surpass the baseline’s
 252 maximum accuracy by 5 percentage points (pp).

253 Under GRPO+ZVF, $p^*_{\text{smooth}}(\alpha = 0.01)$ maintains an advantage, achieving the baseline’s peak accu-
 254 racy using 13% fewer tokens before ultimately outperforming it by 5pp.

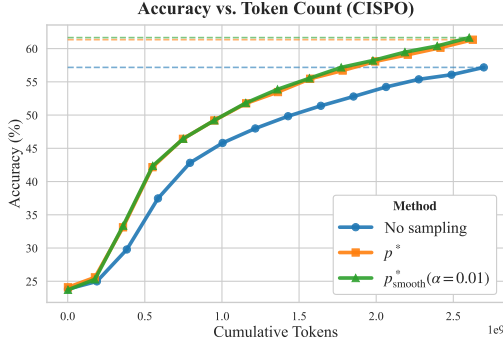
255 Results on AMC benchmark (visualized in Figure 12) show a similar pattern: (1) p^* sampling with
 256 GPRO requires 34% fewer tokens than the baseline to achieve the same accuracy and improves peak
 257 accuracy by 7pp, and (2) $p^*_{\text{smooth}}(0.01)$ sampling with GRPO+ZVF requires 5% fewer tokens than
 258 the baseline to achieve the same accuracy and improves peak accuracy by 3pp.

259 **Qwen3-4B Results.** On AIME (Figure 13), 35% fewer tokens are needed for p^* sampling to reach its
 260 absolute peak than no sampling requires to come within 1% of this accuracy. On AMC (Figure 14),
 261 p^* sampling exceeds the baseline by 2pp, and requires ~24% fewer tokens to reach the peak accuracy.

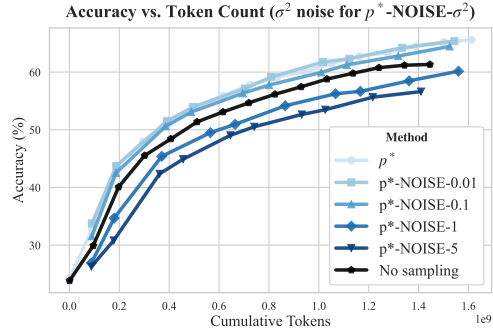
262 **Qwen3-8B Results.** Figure 1 shows AIME accuracy. p^* sampling achieves the baseline’s peak
 263 accuracy using 51% fewer tokens. Because the baseline’s late-stage improvement is marginal, we
 264 also report an asymmetric comparison: matching p^* when it reaches 99% of its peak accuracy. Even
 265 under this relaxed threshold, p^* requires 17% fewer tokens than the baseline. This efficiency trend
 266 extends to the AMC benchmark (Figure 17), where p^* sampling reaches the baseline’s maximum
 267 accuracy using 47% fewer tokens. We also explored smoothed variants (Figure 16). These improve
 268 upon the baseline, but p^* performs best, possibly due to the high fidelity of $|A_i|$ as a proxy (Figure 3).

269 Lastly, Table 1 reports the best checkpoint accuracy for all models and methods. Crucially, these
 270 results demonstrate that cost-aware sampling preserves model performance with respect to standard
 271 GRPO training; in fact, it is even able to improve upon baseline performance.

272 Across the board, cost-aware sampling preserves, and often improves upon, the accuracy of standard
 273 GRPO, while substantially reducing the token cost of policy optimization, as seen in the results above.
 274 Across model scales, we showed that $|A_i|$ is a reliable proxy for G_i , and that cost-aware sampling
 275 reduces the cost of policy optimization. While our primary focus is on the cost of policy optimization,
 276 results also indicate that fewer overall training steps are needed for cost-aware sampling.



(a) **CISPO.** AIME accuracy (pass@1/mean@32) for training Qwen2.5-Math-1.5B-Instruct with CISPO. Cost-aware learning is robust to training objectives.



(b) **Noisy proxies.** Math-1.5B with GRPO when adding zero-mean, σ^2 -variance noise to the proxy $|A_i|$ prior to computing the sampling distribution.

277 5.5 Ablations

278 **Training objective.** We consider a variant of the GRPO objective, CISPO Chen et al. [2025]. CISPO
 279 applies a stop-gradient to the importance weight of the objective. We base our hyperparameters on
 280 findings from Khatri et al. [2025]. Results, reported in Figure 5a, show that our method is not specific
 281 to the GRPO objective; rather, it generalizes across training objectives. p^* sampling and $p^*_{\text{smooth}}(0.01)$
 282 sampling require 31% and 34% fewer tokens respectively to attain the baseline’s peak accuracy,
 283 and ultimately outperform by about 4pp. This finding is further supported by strong sub-optimality
 284 metrics, performance on additional benchmarks, and smoothed variants in Appendix G.7.

285 **Proxy sub-optimality.** We ablate the noise level associated with our proxy $|A_i|$. In particular, we
 286 add zero-mean noise with variance σ^2 to $|A_i|$ prior to computing the sampling distribution. Results
 287 are shown in Figure 5b. We observe that Cost-Aware GRPO is robust to noise up until $\sigma^2 = 0.1$, but
 288 as noise approaches $\sigma^2 = 1.0$, it deteriorates with respect to the baseline.

289 **Choice of α and proxy.** In Appendix G.6, we ablate our choice of α and report results for p^*_{smooth}
 290 with $\alpha = 0.05$ and $\alpha = 0.1$ as well. We also ablate our choice of approximating G_i with $|A_i|$ and
 291 consider using $p_i^* = 1/\sqrt{c_i}$, which corresponds to substituting G_i with 1 for all samples. We observe
 292 that it does not perform as well, which implies a need for correctly approximating the gradient norm
 293 of a sample. In Appendix G.11, we report smoothed variants for the 4B and 8B model as well.

294 6 Conclusion

295 In this paper, we formalize the problem of cost-aware learning and propose Cost-Aware SGD.
 296 Our analysis derives the optimal sampling distribution for convex and strongly convex objectives,
 297 providing theoretical guarantees on the expected cost to attain an error of ϵ . We extensively analyze
 298 the properties of this algorithm.

299 We apply our theoretical insights to GRPO and introduce Cost-Aware GRPO. Our use of $|A_i|$ as
 300 a proxy for G_i in the optimal sampling distribution is supported by strong sub-optimality metrics
 301 across scales. Our experiments demonstrate the broad effectiveness of this approach; it significantly
 302 reduces the tokens used in policy optimization, while maintaining or exceeding the peak performance
 303 of standard baselines.

304 **Future Directions.** In Cost-Aware GRPO, we primarily rely on the magnitude of the advantage
 305 as a proxy for the gradient norm; identifying superior proxies remains an interesting direction for
 306 future research. Furthermore, having demonstrated that cost-aware sampling generalizes successfully
 307 from standard GRPO to the CISPO objective, a natural next step is to explore its integration into a
 308 broader class of RLVR and RLHF algorithms. It would also be of interest to dynamically set α in
 309 smoothed variants of Cost-Aware GRPO in order to obtain superior performance. Lastly, closing the
 310 gap between the upper and lower bounds in Section 4 would be of theoretical significance.

311 References

- 312 Guillaume Alain, Alex Lamb, Chinnadhurai Sankar, Aaron Courville, and Yoshua Bengio. Variance
313 reduction in sgd by distributed importance sampling. *arXiv preprint arXiv:1511.06481*, 2015.
- 314 Art of Problem Solving. Amc problems and solutions. [https://artofproblemsolving.com/
315 wiki/index.php/AMC_Problems_and_Solutions](https://artofproblemsolving.com/wiki/index.php/AMC_Problems_and_Solutions), n.d. Accessed: 2026-01-28.
- 316 Alina Beygelzimer, Sanjoy Dasgupta, and John Langford. Importance weighted active learning. In
317 *Proceedings of the 26th annual international conference on machine learning*, pages 49–56, 2009.
- 318 Sébastien Bubeck et al. Convex optimization: Algorithms and complexity. *Foundations and Trends®
319 in Machine Learning*, 8(3-4):231–357, 2015.
- 320 Aili Chen, Aonian Li, Bangwei Gong, Binyang Jiang, Bo Fei, Bo Yang, Boji Shan, Changqing Yu,
321 Chao Wang, Cheng Zhu, et al. Minimax-m1: Scaling test-time compute efficiently with lightning
322 attention. *arXiv preprint arXiv:2506.13585*, 2025.
- 323 Lichang Chen, Shiyang Li, Jun Yan, Hai Wang, Kalpa Gunaratna, Vikas Yadav, Zheng Tang, Vijay
324 Srinivasan, Tianyi Zhou, Heng Huang, et al. Alpapasus: Training a better alpaca with fewer data.
325 *arXiv preprint arXiv:2307.08701*, 2023a.
- 326 Xiangyi Chen, Sijia Liu, Ruoyu Sun, and Mingyi Hong. On the convergence of a class of adam-type
327 algorithms for non-convex optimization. *arXiv preprint arXiv:1808.02941*, 2018.
- 328 Ziang Chen, Jianfeng Lu, Huajie Qian, Xinshang Wang, and Wotao Yin. Hetersgd: Tackling
329 heterogeneous sampling costs via optimal reweighted stochastic gradient descent. In Francisco
330 Ruiz, Jennifer Dy, and Jan-Willem van de Meent, editors, *Proceedings of The 26th International
331 Conference on Artificial Intelligence and Statistics*, volume 206 of *Proceedings of Machine
332 Learning Research*, pages 10732–10781. PMLR, 25–27 Apr 2023b.
- 333 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
334 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve
335 math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- 336 Corinna Cortes, Giulia DeSalvo, Claudio Gentile, Mehryar Mohri, and Ningshan Zhang. Region-
337 based active learning. In *The 22nd International Conference on Artificial Intelligence and Statistics*,
338 pages 2801–2809. PMLR, 2019.
- 339 Corinna Cortes, Giulia Desalvo, Claudio Gentile, Mehryar Mohri, and Ningshan Zhang. Adaptive
340 region-based active learning. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the
341 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine
342 Learning Research*, pages 2144–2153. PMLR, 13–18 Jul 2020. URL [https://proceedings.
343 mlr.press/v119/cortes20a.html](https://proceedings.mlr.press/v119/cortes20a.html).
- 344 János Csirik, Johannes Bartholomeus Gerardus Frenk, Martine Labbé, and Shuzhong Zhang. Heuris-
345 tics for the 0-1 min-knapsack problem. *Acta Cybernetica*, 10(1-2):15–20, 1991.
- 346 Josef Dai, Xuehai Pan, Ruiyang Sun, Jiaming Ji, Xinbo Xu, Mickel Liu, Yizhou Wang, and
347 Yaodong Yang. Safe rlhf: Safe reinforcement learning from human feedback. *arXiv preprint
348 arXiv:2310.12773*, 2023.
- 349 Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-
350 efficient exact attention with io-awareness. *Advances in neural information processing systems*, 35:
351 16344–16359, 2022.
- 352 Sanjoy Dasgupta. Two faces of active learning. *Theoretical computer science*, 412(19):1767–1781,
353 2011.
- 354 Olivier Devolder, Francois Glineur, and Yurii Nesterov. First-order methods of smooth convex
355 optimization with inexact oracle. *Mathematical Programming*, 146(1):37–75, 2014.
- 356 Mehdi Fatemi, Banafsheh Rafiee, Mingjie Tang, and Kartik Talamadupula. Concise reasoning via
357 reinforcement learning. *arXiv preprint arXiv:2504.05185*, 2025.

- 358 Saeed Ghadimi and Guanghui Lan. Stochastic first-and zeroth-order methods for nonconvex stochastic
359 programming. *SIAM journal on optimization*, 23(4):2341–2368, 2013.
- 360 Nicholas J. A. Harvey, Christopher Liaw, Yaniv Plan, and Sikander Randhawa. Tight analyses for non-
361 smooth stochastic gradient descent. In Alina Beygelzimer and Daniel Hsu, editors, *Proceedings*
362 *of the Thirty-Second Conference on Learning Theory*, volume 99 of *Proceedings of Machine*
363 *Learning Research*, pages 1579–1613. PMLR, 25–28 Jun 2019. URL [https://proceedings.](https://proceedings.mlr.press/v99/harvey19a.html)
364 [mlr.press/v99/harvey19a.html](https://proceedings.mlr.press/v99/harvey19a.html).
- 365 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,
366 and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv*
367 *preprint arXiv:2103.03874*, 2021.
- 368 Prateek Jain, Dheeraj Nagaraj, and Praneeth Netrapalli. Making the last iterate of sgd information
369 theoretically optimal. In Alina Beygelzimer and Daniel Hsu, editors, *Proceedings of the Thirty-*
370 *Second Conference on Learning Theory*, volume 99 of *Proceedings of Machine Learning Research*,
371 pages 1752–1755. PMLR, 25–28 Jun 2019. URL [https://proceedings.mlr.press/v99/](https://proceedings.mlr.press/v99/jain19a.html)
372 [jain19a.html](https://proceedings.mlr.press/v99/jain19a.html).
- 373 Tyler B Johnson and Carlos Guestrin. Training deep models faster with robust, approximate impor-
374 tance sampling. *Advances in Neural Information Processing Systems*, 31, 2018.
- 375 Sham M Kakade. A natural policy gradient. *Advances in neural information processing systems*, 14,
376 2001.
- 377 Angelos Katharopoulos and Francois Fleuret. Not all samples are created equal: Deep learning with
378 importance sampling. In *International conference on machine learning*, pages 2525–2534. PMLR,
379 2018.
- 380 Devvrit Khatri, Lovish Madaan, Rishabh Tiwari, Rachit Bansal, Sai Surya Duvvuri, Manzil Zaheer,
381 Inderjit S Dhillon, David Brandfonbrener, and Rishabh Agarwal. The art of scaling reinforcement
382 learning compute for llms. *arXiv preprint arXiv:2510.13786*, 2025.
- 383 Teun Kloek and Herman K Van Dijk. Bayesian estimates of equation system parameters: an
384 application of integration by monte carlo. *Econometrica: Journal of the Econometric Society*,
385 pages 1–19, 1978.
- 386 Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E.
387 Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model
388 serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating*
389 *Systems Principles*, 2023.
- 390 Xuefeng Li, Haoyang Zou, and Pengfei Liu. Limr: Less is more for rl scaling. *arXiv preprint*
391 *arXiv:2502.11886*, 2025.
- 392 Deanna Needell, Nathan Srebro, and Rachel Ward. Stochastic gradient descent, weighted sampling,
393 and the randomized kaczmarz algorithm. *Advances in neural information processing systems*, 27,
394 2014.
- 395 AS Nemirovsky, DB Yudin, and E Dawson. Wiley-interscience series in discrete mathematics, 1983.
- 396 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong
397 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow
398 instructions with human feedback. *Advances in neural information processing systems*, 35:27730–
399 27744, 2022.
- 400 Mansheej Paul, Surya Ganguli, and Gintare Karolina Dziugaite. Deep learning on a data diet: Finding
401 important examples early in training. *Advances in neural information processing systems*, 34:
402 20596–20607, 2021.
- 403 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea
404 Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances*
405 *in neural information processing systems*, 36:53728–53741, 2023.

- 406 Alexander Rakhlin, Ohad Shamir, and Karthik Sridharan. Making gradient descent optimal for
407 strongly convex stochastic optimization. *arXiv preprint arXiv:1109.5647*, 2011.
- 408 John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region
409 policy optimization. In *International conference on machine learning*, pages 1889–1897. PMLR,
410 2015.
- 411 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
412 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 413 Ohad Shamir and Tong Zhang. Stochastic gradient descent for non-smooth optimization: Convergence
414 results and optimal averaging schemes. In *International conference on machine learning*, pages
415 71–79. PMLR, 2013.
- 416 Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang,
417 Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathemat-
418 ical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- 419 Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng,
420 Haibin Lin, and Chuan Wu. Hybridflow: A flexible and efficient rlhf framework. In *Proceedings
421 of the Twentieth European Conference on Computer Systems*, pages 1279–1297, 2025.
- 422 Ben Sorscher, Robert Geirhos, Shashank Shekhar, Surya Ganguli, and Ari Morcos. Beyond neural
423 scaling laws: beating power law scaling via data pruning. *Advances in Neural Information
424 Processing Systems*, 35:19523–19536, 2022.
- 425 Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient
426 methods for reinforcement learning with function approximation. In S. Solla, T. Leen, and
427 K. Müller, editors, *Advances in Neural Information Processing Systems*, volume 12. MIT
428 Press, 1999. URL [https://proceedings.neurips.cc/paper_files/paper/1999/file/
429 464d828b85b0bed98e80ade0a5c43b0f-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/1999/file/464d828b85b0bed98e80ade0a5c43b0f-Paper.pdf).
- 430 Hemish Veeraboina. Aime problem set 1983-2024, 2023. URL [https://www.kaggle.com/
431 datasets/hemishveeraboina/aime-problem-set-1983-2024](https://www.kaggle.com/datasets/hemishveeraboina/aime-problem-set-1983-2024).
- 432 Yiping Wang, Qing Yang, Zhiyuan Zeng, Liliang Ren, Liyuan Liu, Baolin Peng, Hao Cheng, Xuehai
433 He, Kuan Wang, Jianfeng Gao, et al. Reinforcement learning for reasoning in large language
434 models with one training example. *arXiv preprint arXiv:2504.20571*, 2025.
- 435 Sang Michael Xie, Shibani Santurkar, Tengyu Ma, and Percy S Liang. Data selection for language
436 models via importance resampling. *Advances in Neural Information Processing Systems*, 36:
437 34201–34227, 2023.
- 438 Yixuan Even Xu, Yash Savani, Fei Fang, and Zico Kolter. Not all rollouts are useful: Down-sampling
439 rollouts in llm reinforcement learning. *arXiv preprint arXiv:2504.13818*, 2025.
- 440 An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu, Jian-
441 hong Tu, Jingren Zhou, Junyang Lin, et al. Qwen2. 5-math technical report: Toward mathematical
442 expert model via self-improvement. *arXiv preprint arXiv:2409.12122*, 2024.
- 443 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang
444 Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*,
445 2025.
- 446 Qiying Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian
447 Fan, Gaohong Liu, Lingjun Liu, et al. Dapo: An open-source llm reinforcement learning system at
448 scale. *arXiv preprint arXiv:2503.14476*, 2025.
- 449 Jingzhao Zhang, Tianxing He, Suvrit Sra, and Ali Jadbabaie. Why gradient clipping accelerates
450 training: A theoretical justification for adaptivity, 2020. URL [https://arxiv.org/abs/1905.
451 11881](https://arxiv.org/abs/1905.11881).
- 452 Peilin Zhao and Tong Zhang. Stochastic optimization with importance sampling for regularized loss
453 minimization. In *international conference on machine learning*, pages 1–9. PMLR, 2015.

- 454 Haizhong Zheng, Yang Zhou, Brian R Bartoldson, Bhavya Kailkhura, Fan Lai, Jiawei Zhao, and
455 Beidi Chen. Act only when it pays: Efficient reinforcement learning for llm reasoning via selective
456 rollouts. *arXiv preprint arXiv:2506.02177*, 2025.
- 457 Rui Zheng, Shihan Dou, Songyang Gao, Yuan Hua, Wei Shen, Binghai Wang, Yan Liu, Senjie Jin,
458 Qin Liu, Yuhao Zhou, et al. Secrets of rlhf in large language models part i: Ppo. *arXiv preprint*
459 *arXiv:2307.04964*, 2023.

460	Contents	
461	1 Introduction	1
462	2 Related Work	2
463	3 Problem Formulation	3
464	4 Cost-Aware SGD	3
465	4.1 Cost Analysis	3
466	4.2 Derivation of the Optimal Distribution	4
467	4.3 Cost-Improvement over Baselines	4
468	4.4 Synthetic Validation	4
469	4.5 Proxy Sub-Optimality	5
470	4.6 Lower Bound	5
471	5 Cost-Aware GRPO	6
472	5.1 Datasets and Models	6
473	5.2 Cost-Aware Learning with GRPO	6
474	5.3 Settings and Methods	7
475	5.4 Results	8
476	5.5 Ablations	9
477	6 Conclusion	9
478	A Related Work	15
479	B Algorithms	16
480	B.1 Cost-Aware SGD	16
481	B.2 Cost-Aware GRPO	16
482	C Subset Selection	17
483	D Proofs	18
484	D.1 Variance reduction strategy derivation	18
485	D.2 Cost-Improvement & Comparison Analysis	18
486	D.3 Min Cost Knapsack	19
487	D.4 Strongly Convex Subset Selection	20
488	D.5 Sub-Optimality Proofs	20
489	E Lower bound	23
490	F Extensions and Generality of the Theoretical Analysis	26
491	F.1 Differences between Theory and Practice	26
492	F.2 Generality of the Cost Objective in Non-Convex Settings	26
493	F.3 Average Iterate (\bar{x}_T) vs. Last Iterate (x_T)	26
494	G Experimental Details and Ablations	27
495	G.1 Experimental Details	27
496	G.2 Implementation	27
497	G.3 Importance Weighting with GRPO	27
498	G.4 Hyperparameters	28
499	G.5 Computational Requirements	29
500	G.6 Ablation on variants of p^*	30
501	G.7 Full CISPO results	31
502	G.8 Full sub-optimality results	32
503	G.9 Additional results for Qwen2.5-Math-1.5B-Instruct	33
504	G.10 Additional results for Qwen3-4B	34
505	G.11 Additional results for Qwen3-8B	35

506 **A Related Work**

507 We describe further related work.

508 **Training Data Pruning/ Selection.** There is a large literature focused on pruning a training dataset.
509 Typically, the objective is to obtain a better predictor or maintain performance with fewer samples.
510 Related works include [Sorscher et al., 2022, Paul et al., 2021, Xie et al., 2023, Chen et al., 2023a].
511 Our work differs from these as it explicitly considers a variable cost per sample.

512 **Active learning.** There is a significant body of work studying the problem of active learning, which
513 considers when to request more information about training samples. We refer the reader to [Dasgupta,
514 2011] for an in-depth overview of the literature. Most related to our work, Beygelzimer et al. [2009],
515 Cortes et al. [2020, 2019] use importance weighting to probabilistically select samples for which
516 to query a label. Generally, these algorithms fix a hypothesis set a priori, in which hypotheses are
517 eliminated as more data is acquired. Our work shares the similarity of deciding when to request more
518 information, however, in contrast, our work requests the *gradient* of a sample. Furthermore, we do
519 not use a current hypothesis set in order to make these decisions.

520 **B Algorithms**

521 **B.1 Cost-Aware SGD**

Algorithm 1 Cost-Aware SGD

- 1: **Input:** Initial point x_1 , iterations T , step size η , sampling distribution \mathbf{p}^* .
 - 2: **for** $t = 1$ **to** T **do**
 - 3: Sample index $i_t \sim \mathbf{p}^*$
 - 4: Query $\nabla f_{i_t}(x_t)$, incur cost c_{i_t} .
 - 5: Compute importance-weighted stochastic gradient:
 - 6: $\tilde{g}_t = \frac{1}{n p_{i_t}} \nabla f_{i_t}(x_t)$
 - 7: Update parameter: $x_{t+1} = \Pi_{\mathcal{X}}(x_t - \eta \tilde{g}_t)$
 - 8: **end for**
 - 9: **Output:** $\bar{x}_T = \frac{1}{T} \sum_{t=1}^T x_t$ (or suffix averaging solution for strongly convex case [Rakhlin et al., 2011])
-

522 **B.2 Cost-Aware GRPO**

Algorithm 2 Cost-Aware GRPO

- Input:** Initial policy π_θ , Dataset \mathcal{D} , Prompt batch size n , Rollouts per prompt M , Mini-batch size B , Iterations K
- for** iteration $k = 1$ **to** K **do**
- // Phase 1: Data Collection*
- Sample prompt batch $X = \{x_1, \dots, x_n\} \sim \mathcal{D}$
- Generate M outputs per prompt: $y_{i,j} \sim \pi_{\theta_{\text{old}}}(\cdot | x_i)$
- Form rollout dataset $\mathcal{D}_{\text{step}}$ of size $N = n \times M$
- Compute advantages for all samples in $\mathcal{D}_{\text{step}}$
- // Phase 2: Compute Sampling Distribution*
- Estimate weights $s_u \leftarrow G_u / \sqrt{c_u}$ for all $u \in \mathcal{D}_{\text{step}}$
- Compute probabilities $p_u^* \leftarrow s_u / \sum_{v \in \mathcal{D}_{\text{step}}} s_v$
- // Phase 3: Importance Weighted Updates*
- Set number of updates $T \leftarrow \lceil N/B \rceil$
- for** update step $t = 1$ **to** T **do**
- // Sample mini-batch and optimize policy*
- Sample mini-batch $\sim \text{Multinomial}(\text{support} = \mathcal{D}_{\text{step}}, \text{probs} = p^*, \text{size} = B)$
- Update π_θ by maximizing $\hat{\mathcal{J}}_{\text{GRPO}}^{(p^*)}(\theta)$ over mini-batch
- end for**
- end for**
-

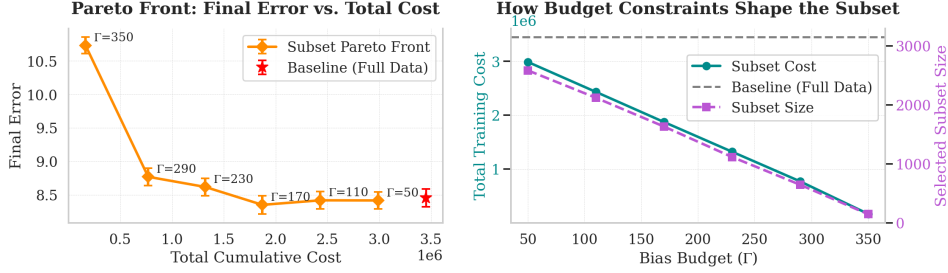


Figure 6: Synthetic validation of the greedy subset selection algorithm. Γ is a tunable parameter representing the allowable bias budget. As it is increased, the greedy algorithm chooses fewer (and cheaper) samples to train on.

523 C Subset Selection

524 The importance sampling strategy \mathbf{p}^* achieves cost $(\sum_i G_i \sqrt{c_i})^2$, but may remain expensive when
 525 high-cost samples have large gradients. To reduce cost, we consider optimizing over a subset
 526 $\pi \subseteq \{1, \dots, n\}$, defining $f_\pi(x) = \frac{1}{n} \sum_{i \in \pi} f_i(x)$. This introduces bias, as optimization converges to
 527 x_π^* instead of x^* , but reduces computational cost.

528 **Biased SGD on a subset.** Let $\|\nabla f_i(x)\| \leq G_i$. For a subset π and sampling distribution \mathbf{q}
 529 supported on π , define the bias $\beta_\pi = \frac{1}{n} \sum_{i \notin \pi} G_i$ and note $\|\nabla f(x) - \nabla f_\pi(x)\| \leq \beta_\pi$. Stan-
 530 dard SGD results [Devolder et al., 2014] imply that with $\eta_t \propto 1/\sqrt{t}$, $\mathbb{E}[f(\bar{x}_T) - f(x^*)] \leq$
 531 $\frac{D\sqrt{\sigma_\pi^2(\mathbf{q})}}{\sqrt{T}} + D\beta_\pi$, where $\sigma_\pi^2(\mathbf{q}) \leq \frac{1}{n^2} \sum_{i \in \pi} \frac{G_i^2}{q_i}$. Thus, the cost to achieve error $\epsilon > D\beta_\pi$ is:
 532 $\mathcal{K}_\pi(\epsilon) = \frac{D^2}{n^2(\epsilon - D\beta_\pi)^2} (\sum_{i \in \pi} G_i \sqrt{c_i})^2$, where we used the optimal sampling distribution restricted to
 533 π .

534 **Optimal subset selection.** We minimize $\mathcal{K}_\pi(\epsilon)$ subject to a bias budget $D\beta_\pi \leq \Gamma$.

535 **Theorem C.1** (Knapsack equivalence). *This problem is equivalent to a min-cost knapsack:*

$$\min_{z_i \in \{0,1\}} \sum_i z_i G_i \sqrt{c_i} \quad \text{s.t.} \quad \sum_i z_i G_i \geq \sum_i G_i - \frac{n\Gamma}{D}.$$

536 We defer the proof to Appendix D.2. A similar result for strongly convex objectives is in Ap-
 537 pendix D.4.

538 **Greedy approximation.** A standard greedy algorithm selects items by density $\rho_i = v_i/\omega_i = 1/\sqrt{c_i}$,
 539 i.e., simply choosing the lowest-cost samples until the constraint is met, yielding a 2-approximation
 540 [Csirik et al., 1991]. Remarkably, the optimal greedy strategy is simply to *select samples with the*
 541 *lowest cost c_i* , regardless of their Lipschitz constant G_i . Under a similar setting to Section 4.4, we
 542 implement this greedy approximation. Within a smaller bias budget, subset selection is able to attain
 543 similar error at a significantly reduced cost. As the selected subset becomes much smaller, the error
 544 increases due to bias.

545 D Proofs

546 In this section we provide the proofs and derivations for statements made in the paper.

547 D.1 Variance reduction strategy derivation

548 In this case, we have:

$$J(\mathbf{p}) = S(\mathbf{p})C(\mathbf{p}) = \left(\frac{1}{n^2} \sum_{j=1}^n \frac{G_j^2}{p_j} \right) \left(\sum_{i=1}^n p_i \right),$$

549 where $\mathbf{p} \in \Delta_n$. From Cauchy-Schwarz, we can lower bound $J(p)$:

$$J(\mathbf{p}) \geq \frac{1}{n^2} \left(\sum_{i=1}^n G_i \right)^2.$$

550 This lower bound is attained when $G_i/\sqrt{p_i} \propto \sqrt{p_i}$, or equivalently, $p_i \propto G_i$. Hence, when costs are
551 uniform, we focus on minimizing the second moment. This gives:

$$(p_{\text{var}})_i = \frac{G_i}{\sum_j G_j}. \quad (9)$$

552 When we apply p_{var} to the non-uniform cost setting, the total expected cost to attain error ϵ becomes:

$$\begin{aligned} K_{\text{var}}(\epsilon) &= \frac{D^2}{\epsilon^2 n^2} \left(\sum_{i=1}^n G_i \right)^2 \left(\sum_{i=1}^n (p_{\text{var}})_i c_i \right) \\ &= \frac{D^2}{\epsilon^2 n^2} \left(\sum_{i=1}^n G_i \right)^2 \left(\sum_{i=1}^n \frac{G_i}{\sum_j G_j} c_i \right) \\ &= \frac{D^2}{\epsilon^2 n^2} \left(\sum_{i=1}^n G_i \right) \left(\sum_{i=1}^n G_i c_i \right). \end{aligned}$$

553 D.2 Cost-Improvement & Comparison Analysis

554 **Theorem D.1** (Comparison of Adaptive vs. Uniform). *Unlike the optimal strategy, the adaptive*
555 *strategy p_{var} is not guaranteed to outperform uniform sampling. The ratio of their costs is given by:*

$$\frac{K_{\text{var}}(\epsilon)}{K_{\text{unif}}(\epsilon)} = \frac{(\sum_{i=1}^n G_i)(\sum_{i=1}^n G_i c_i)}{(\sum_{i=1}^n G_i^2)(\sum_{i=1}^n c_i)}. \quad (10)$$

556 Let $\mathbb{E}[\cdot]$ denote the expectation taken over a uniform distribution of indices $i \sim \text{Uniform}(1, \dots, n)$.
557 Then $K_{\text{var}} \leq K_{\text{unif}}$ holds if and only if:

$$\mathbb{E}[G] \text{Cov}(G, c) \leq \mathbb{E}[c] \text{Var}(G). \quad (11)$$

558 *Consequently, if gradients and costs are negatively correlated or uncorrelated, adaptive sampling is*
559 *superior. However, if gradients and costs are strongly positively correlated (i.e., samples with large*
560 *gradients are disproportionately expensive), uniform sampling may yield a lower total cost.*

561 *Proof.* Recall that $K_{\text{unif}} \propto (\sum G_i^2)(\sum c_i)$ and $K_{\text{var}} \propto (\sum G_i)(\sum G_i c_i)$. We can rewrite these sums in
562 terms of expectations over uniformly sampled indices. For any vector x , $\sum x_i = n \mathbb{E}[x]$. Thus:

$$\frac{K_{\text{var}}}{K_{\text{unif}}} = \frac{(n \mathbb{E}[G])(n \mathbb{E}[Gc])}{(n \mathbb{E}[G^2])(n \mathbb{E}[c])} = \frac{\mathbb{E}[G] \mathbb{E}[Gc]}{\mathbb{E}[G^2] \mathbb{E}[c]}.$$

563 We seek the condition under which this ratio is ≤ 1 , which implies $\mathbb{E}[G] \mathbb{E}[Gc] \leq \mathbb{E}[G^2] \mathbb{E}[c]$. We
564 apply the definitions $\mathbb{E}[Gc] = \text{Cov}(G, c) + \mathbb{E}[G] \mathbb{E}[c]$ and $\mathbb{E}[G^2] = \text{Var}(G) + \mathbb{E}[G]^2$. Substituting
565 these into the inequality:

$$\mathbb{E}[G](\text{Cov}(G, c) + \mathbb{E}[G] \mathbb{E}[c]) \leq (\text{Var}(G) + \mathbb{E}[G]^2) \mathbb{E}[c].$$

566 Expanding both sides:

$$\mathbb{E}[G] \text{Cov}(G, c) + \mathbb{E}[G]^2 \mathbb{E}[c] \leq \mathbb{E}[c] \text{Var}(G) + \mathbb{E}[G]^2 \mathbb{E}[c].$$

567 Subtracting the common term $\mathbb{E}[G]^2 \mathbb{E}[c]$ from both sides yields the final condition:

$$\mathbb{E}[G] \text{Cov}(G, c) \leq \mathbb{E}[c] \text{Var}(G).$$

568 Since $\mathbb{E}[G]$, $\mathbb{E}[c]$, and $\text{Var}(G)$ are non-negative, this inequality always holds if $\text{Cov}(G, c) \leq 0$
 569 (negative or zero correlation). It is violated only if $\text{Cov}(G, c)$ is sufficiently large and positive. \square

570 D.3 Min Cost Knapsack

571 **Theorem D.2** (Equivalence to Min-Knapsack Problem). *Let Γ be the maximum allowable bias budget*
 572 *such that convergence is possible. Finding the optimal subset π^* that minimizes the expected total*
 573 *cost subject to this budget is equivalent to solving the **Min-Cost Knapsack Problem**.*

574 *Specifically, let $z_i \in \{0, 1\}$ be the indicator variable where $z_i = 1$ if $i \in \pi$. The problem is:*

$$\min_z \sum_{i=1}^n z_i \omega_i \quad \text{subject to} \quad \sum_{i=1}^n z_i v_i \geq V_{\text{req}}, \quad (12)$$

575 where the ‘‘item cost’’ is $\omega_i = G_i \sqrt{c_i}$, the ‘‘item value’’ is $v_i = G_i$, and the required value coverage
 576 is $V_{\text{req}} = \sum_{j=1}^n G_j - \frac{n\Gamma}{D}$.

577 *Proof.* We aim to minimize the cost complexity factor derived for the optimal distribution restricted to
 578 π : $J(\pi) = (\sum_{i \in \pi} G_i \sqrt{c_i})^2$. The constraint is on the bias error contribution $D\beta_\pi \leq \Gamma < \epsilon$. Recalling
 579 the definition of bias:

$$\beta_\pi \leq \frac{1}{n} \sum_{j \notin \pi} G_j \implies D \left(\frac{1}{n} \sum_{j \notin \pi} G_j \right) \leq \Gamma.$$

580 Minimizing the sum of gradients *excluded* from π is equivalent to maximizing the sum of gradients
 581 *included* in π . We rewrite the constraint:

$$\sum_{j \notin \pi} G_j \leq \frac{n\Gamma}{D} \implies \sum_{j=1}^n G_j - \sum_{i \in \pi} G_i \leq \frac{n\Gamma}{D} \implies \sum_{i \in \pi} G_i \geq \sum_{j=1}^n G_j - \frac{n\Gamma}{D}.$$

582 Therefore, if we let $V_{\text{req}} = \sum_{j=1}^n G_j - \frac{n\Gamma}{D}$, let $\omega_i = G_i \sqrt{c_i}$ as the penalty for including item i , and let
 583 $v_i = G_i$ as the gain toward satisfying the constraint, we recover the formulation in Equation (12). \square

584 Using this algorithm, the worst case $\mathbb{E}[K(\epsilon)]$, where $\epsilon > \Gamma$, for $D\beta_\pi \leq \Gamma$ is:

$$\min_{\{\pi: \beta_\pi \leq \Gamma/D\}} \frac{D^2 J(\pi)}{n^2 (\epsilon - \Gamma)^2}.$$

585 The bias constraint $D\beta_\pi \leq \Gamma$ determines the required gradient coverage $V_{\text{req}} = \sum_i G_i - n\Gamma/D$, and
 586 thus defines a specific Min-Cost Knapsack instance as in Theorem C.1. By varying Γ over its feasible
 587 range and solving the associated knapsack problem for each choice, we obtain an optimal subset
 588 $\pi(\Gamma)$ minimizing the total cost subject to that bias budget. Since the total expected complexity

$$\mathbb{E}[K_\pi(\epsilon)] = \frac{D^2 J(\pi)}{n^2 (\epsilon - D\beta_\pi)^2}$$

589 depends on both the numerator (the sampling cost) and the denominator (the bias floor), enumerating
 590 over Γ allows us to evaluate the achievable error–cost tradeoff for every subset. Consequently, we
 591 can compute the optimal target accuracy ϵ by selecting the bias level Γ that yields the smallest overall
 592 expected cost.

593 **D.4 Strongly Convex Subset Selection**

594 **Theorem D.3** (Subset Selection for Strongly Convex Functions). *Assume f is μ -strongly convex. Let*
 595 *$z_i \in \{0, 1\}$ be the indicator variable for including data point i in subset π . We seek to minimize the cost*
 596 *complexity $\sum z_i G_i \sqrt{c_i}$ subject to a budget Γ on the function value suboptimality $f(x_\pi^*) - f(x^*) \leq \Gamma$.*

597 *This problem is structurally equivalent to the Min-Cost Knapsack formulation in Theorem C.1, with a*
 598 *modified requirement V_{req} . The optimal solution to the LP relaxation is obtained by selecting indices*
 599 *i strictly based on increasing cost c_i .*

600 *Proof.* Using the Polyak-Lojasiewicz (PL) inequality and the fact that $\nabla f_S(x_S^*) = 0$, the gradient at
 601 the proxy minimizer is determined solely by the excluded gradients:

$$\nabla f(x_\pi^*) = \frac{1}{n} \sum_{j \notin \pi} \nabla f_j(x_\pi^*).$$

602 We bound the norm of this gradient using the triangle inequality ($\|\sum v\| \leq \sum \|v\|$) and the gradient
 603 upper bound definition ($\|\nabla f_j(x)\| \leq G_j$):

$$\begin{aligned} \|\nabla f(x_\pi^*)\| &= \left\| \frac{1}{n} \sum_{j \notin \pi} \nabla f_j(x_\pi^*) \right\| \\ &\leq \frac{1}{n} \sum_{j \notin \pi} \|\nabla f_j(x_\pi^*)\| \\ &\leq \frac{1}{n} \sum_{j \notin \pi} G_j. \end{aligned}$$

604 Substituting this back into the PL inequality ($f(x_\pi^*) - f(x^*) \leq \frac{1}{2\mu} \|\nabla f(x_\pi^*)\|^2$) yields:

$$f(x_\pi^*) - f(x^*) \leq \frac{1}{2\mu n^2} \left(\sum_{j \notin \pi} G_j \right)^2 = \frac{1}{2\mu n^2} \left(\sum_{i=1}^n (1 - z_i) G_i \right)^2.$$

605 We enforce the constraint $f(x_\pi^*) - f(x^*) \leq \Gamma$. Because the gradient norms G_i are non-negative, we
 606 can take the square root of the inequality to linearize the constraint without altering the optimization
 607 region:

$$\sum_{i=1}^n (1 - z_i) G_i \leq n \sqrt{2\mu\Gamma}.$$

608 This is now a linear constraint on the excluded mass. Let $B = n \sqrt{2\mu\Gamma}$ be the effective budget for the
 609 excluded gradients. We rewrite this in terms of the included variables z_i :

$$\sum_{j=1}^n G_j - \sum_{i=1}^n z_i G_i \leq B \implies \sum_{i=1}^n z_i G_i \geq \sum_{j=1}^n G_j - B.$$

610 The optimization problem is to minimize $\sum z_i (G_i \sqrt{c_i})$ subject to $\sum z_i G_i \geq V_{\text{req}}$. This is identical
 611 to Equation 12. This problem is structurally equivalent to the Min-Cost Knapsack formulation in
 612 Theorem C.1, with a modified requirement V_{req} . Its LP relaxation is a fractional knapsack problem
 613 whose optimal solution is obtained by ordering indices by increasing evaluation cost c_i (equivalently,
 614 decreasing density v_i/ω_i). \square

615 **D.5 Sub-Optimality Proofs**

Theorem D.4 (Sub-optimality Gap via Cost-Biased Divergence). *Let p^* be the optimal cost-aware*
sampling distribution and p' be any empirical sampling distribution. Let the expected cost per iteration
for a distribution p be $C(p) = \sum_{i=1}^n p_i c_i$. We define the cost-biased distribution \tilde{p} corresponding
to p as:

$$\tilde{p}_i = \frac{p_i c_i}{C(p)}.$$

The optimality gap is exactly determined by the Pearson χ^2 -divergence between the cost-biased optimal distribution and the cost-biased empirical distribution:

$$\frac{J(p')}{J(p^*)} = 1 + D_{\chi^2}(\tilde{p}^* || \tilde{p}').$$

Proof. Recall that minimizing the total expected cost is equivalent to minimizing the product $J(p) = S(p)C(p)$. The optimal distribution is defined as $p_i^* = \frac{G_i/\sqrt{c_i}}{Z_{den}}$, where $Z_{den} = \sum_{j=1}^n \frac{G_j}{\sqrt{c_j}}$. We can rearrange this to express the true gradient norms in terms of p_i^* :

$$G_i = Z_{den} p_i^* \sqrt{c_i}.$$

Substituting this expression for G_i into the variance term $S(p') = \frac{1}{n^2} \sum_{i=1}^n \frac{G_i^2}{p_i'}$ yields the objective for the empirical distribution:

$$J(p') = \frac{Z_{den}^2}{n^2} C(p') \left(\sum_{i=1}^n \frac{(p_i^*)^2 c_i}{p_i'} \right).$$

Evaluating this same expression for the optimal distribution p^* (where $p_i' = p_i^*$) gives $J(p^*) = \frac{Z_{den}^2}{n^2} C(p^*)^2$. Taking the ratio of the two costs isolates the gap:

$$\frac{J(p')}{J(p^*)} = \frac{C(p') \sum_{i=1}^n c_i \frac{(p_i^*)^2}{p_i'}}{C(p^*)^2}$$

By the definition of the cost-biased distribution, we have $p_i = \tilde{p}_i \frac{C(p)}{c_i}$. Substituting this expression for both p^* and p' into the summation term:

$$\sum_{i=1}^n c_i \frac{(p_i^*)^2}{p_i'} = \sum_{i=1}^n c_i \frac{\left(\tilde{p}_i^* \frac{C(p^*)}{c_i} \right)^2}{\tilde{p}_i' \frac{C(p')}{c_i}}.$$

Simplifying the terms inside the sum:

$$= \sum_{i=1}^n c_i \frac{(\tilde{p}_i^*)^2 \frac{C(p^*)^2}{c_i^2}}{\tilde{p}_i' \frac{C(p')}{c_i}} = \frac{C(p^*)^2}{C(p')} \sum_{i=1}^n \frac{(\tilde{p}_i^*)^2}{\tilde{p}_i'}.$$

Substituting this simplified summation back into our original cost ratio:

$$\frac{J(p')}{J(p^*)} = \frac{C(p')}{C(p^*)^2} \left(\frac{C(p^*)^2}{C(p')} \sum_{i=1}^n \frac{(\tilde{p}_i^*)^2}{\tilde{p}_i'} \right) = \sum_{i=1}^n \frac{(\tilde{p}_i^*)^2}{\tilde{p}_i'}.$$

616 Using the definition of the Pearson χ^2 -divergence, $D_{\chi^2}(P||Q) = \sum \frac{P^2}{Q} - 1$, we arrive at the identity.

Theorem D.5. Let p^* be the optimal cost-aware sampling distribution using the true gradient norms G_i , and let p' be the empirical sampling distribution using estimated gradient norms G'_i , such that $p'_i \propto G'_i/\sqrt{c_i}$. Assume the estimates follow an additive noise model $G'_i = G_i + \epsilon_i$, where ϵ_i are independent, zero-mean noise terms with variance σ_ϵ^2 , independent of c_i . Let ρ be the Pearson correlation coefficient between the true norms G_i and the estimates G'_i . Assuming the noise ϵ_i is relatively small compared to G_i , the expected optimality gap is approximately bounded by:

$$\frac{\mathbb{E}[J(p')]}{J(p^*)} \approx 1 + \left(\frac{1 - \rho^2}{\rho^2} \right) \sigma_G^2 \frac{\sum_{i=1}^n \frac{\sqrt{c_i}}{G_i}}{\sum_{i=1}^n G_i \sqrt{c_i}}$$

617 where σ_G^2 is the variance of the true gradient norms.

Proof. The cost objective for a distribution p is $J(p) = S(p)C(p)$. For the estimated distribution p' , we have $p'_i = \frac{G'_i/\sqrt{c_i}}{\sum_{j=1}^n G'_j/\sqrt{c_j}}$. Substituting this into the cost objective yields:

$$J(p') = \frac{1}{n^2} \left(\sum_{i=1}^n \frac{G_i^2 \sqrt{c_i}}{G'_i} \right) \left(\sum_{i=1}^n G'_i \sqrt{c_i} \right).$$

Under the model $G'_i = G_i + \epsilon_i$, the covariance between G and G' is strictly the variance of G , namely $Cov(G, G') = \sigma_G^2$. The variance of G' is $Var(G') = \sigma_G^2 + \sigma_\epsilon^2$. The Pearson correlation coefficient ρ is defined as:

$$\rho = \frac{Cov(G, G')}{\sigma_G \sigma_{G'}} = \frac{\sigma_G^2}{\sigma_G \sqrt{\sigma_G^2 + \sigma_\epsilon^2}} = \frac{\sigma_G}{\sqrt{\sigma_G^2 + \sigma_\epsilon^2}}$$

Solving for the noise variance σ_ϵ^2 in terms of ρ :

$$\sigma_\epsilon^2 = \sigma_G^2 \left(\frac{1 - \rho^2}{\rho^2} \right).$$

Because G'_i appears in the denominator of the second moment, we apply a second-order Taylor expansion of $1/G'_i$ centered around G_i :

$$\frac{1}{G_i + \epsilon_i} \approx \frac{1}{G_i} - \frac{\epsilon_i}{G_i^2} + \frac{\epsilon_i^2}{G_i^3}$$

618 Substituting this into the first term of $J(p')$ and taking the expectation with respect to the noise ϵ :

$$\mathbb{E} \left[\frac{G_i^2 \sqrt{c_i}}{G'_i} \right] \approx G_i^2 \sqrt{c_i} \left(\frac{1}{G_i} - \frac{\mathbb{E}[\epsilon_i]}{G_i^2} + \frac{\mathbb{E}[\epsilon_i^2]}{G_i^3} \right) = G_i \sqrt{c_i} + \sigma_\epsilon^2 \frac{\sqrt{c_i}}{G_i}.$$

Taking the expectation of the full cost $J(p')$, noting that the two sums are weakly dependent for large n and $\mathbb{E}[G'_i] = G_i$:

$$\begin{aligned} \mathbb{E}[J(p')] &\approx \frac{1}{n^2} \left(\sum_{i=1}^n \left(G_i \sqrt{c_i} + \sigma_\epsilon^2 \frac{\sqrt{c_i}}{G_i} \right) \right) \left(\sum_{i=1}^n G_i \sqrt{c_i} \right) \\ \mathbb{E}[J(p')] &\approx \frac{1}{n^2} \left(\sum_{i=1}^n G_i \sqrt{c_i} \right)^2 + \frac{\sigma_\epsilon^2}{n^2} \left(\sum_{i=1}^n \frac{\sqrt{c_i}}{G_i} \right) \left(\sum_{i=1}^n G_i \sqrt{c_i} \right) \end{aligned}$$

Recognizing that $J(p^*) = \frac{1}{n^2} (\sum_{i=1}^n G_i \sqrt{c_i})^2$, we divide by $J(p^*)$:

$$\frac{\mathbb{E}[J(p')]}{J(p^*)} \approx 1 + \sigma_\epsilon^2 \frac{\sum_{i=1}^n \frac{\sqrt{c_i}}{G_i}}{\sum_{i=1}^n G_i \sqrt{c_i}}.$$

619 Finally, substituting $\sigma_\epsilon^2 = \sigma_G^2 \left(\frac{1 - \rho^2}{\rho^2} \right)$ completes the proof.

620 **E Lower bound**

621 We consider the minimization of a finite-sum objective:

$$\min_{x \in \mathcal{X}} F(x) = \frac{1}{n} \sum_{i=1}^n f_i(x).$$

622 where each $f_i : \mathcal{X} \rightarrow \mathbb{R}$ is convex and G -Lipschitz continuous. We assume access to an exact
 623 first-order oracle: when queried at index i and point x , the oracle returns the exact gradient $\nabla f_i(x)$.
 624 Crucially, we assume there is a variable cost $c_i > 0$ to querying the i -th component (at any point x).

625 The goal is to minimize the total expected accumulated cost required to output a point \hat{x} such that
 626 $\mathbb{E}[F(\hat{x}) - F(x^*)] \leq \epsilon$, where the expectation is taken over any randomization in the algorithm
 627 producing \hat{x} .

628 **Theorem E.1.** Fix $\epsilon, G > 0$. For any number of components $n \geq G^2/\epsilon^2$, any set of nonnegative
 629 query costs $\{c_i\}_{i=1}^n$, and any (possibly randomized) algorithm with an expected error guarantee
 630 $\mathbb{E}[F(\hat{x}) - \min F(x)] \leq \epsilon$, there exists a convex, G -Lipschitz finite-sum problem instance over
 631 $\mathcal{X} = [-1, 1]$ such that the algorithm's expected total query cost is at least:

$$\Omega\left(\frac{G^2}{\epsilon^2} \left(\frac{1}{n} \sum_{i \in S^*} \sqrt{c_i}\right)^2\right),$$

632 where $S^* \subseteq [n]$ is any set of components such that following cost-uniformity condition holds:

$$\frac{\max_{i \in S^*} \sqrt{c_i}}{\frac{1}{n} \sum_{j \in S^*} \sqrt{c_j}} \leq \frac{G}{40\epsilon}. \quad (13)$$

633 **Remark.** If we let the costs be sorted $c_1 \leq c_2 \leq \dots \leq c_n$, then the set S^* that maximizes the lower
 634 bound (if it is nonempty) must be of the form $S = [k]$ for some $1 \leq k \leq n$, namely, the set of k
 635 cheapest costs. Indeed, if $j \in S^*$ then also $i \in S^*$ for any $i < j$, since otherwise including i into S^*
 636 increases the lower bound, while still satisfying Eq. (13): it can only increase the denominator on the
 637 LHS, while not affecting the numerator.

638 For reasonably uniform costs, say perfectly uniform up to a small (constant) fraction of arbitrarily
 639 large outliers, the condition in Eq. (13) is satisfied by fixing S^* to contain all components except for
 640 the outliers, as long as G/ϵ is larger than a constant. It could be, however, that the maximal S^* is
 641 empty (and the lower bound is zero) if the costs are highly non-uniform. For example, if the costs
 642 are exponentially increasing, say $c_i = 4^i$ for all i , then the ratio between the maximum and the sum
 643 of $\sqrt{c_i}$ over any set of the form $S^* = [k]$ is between $\frac{1}{2}$ and 1, and so the cost-uniformity condition
 644 implies that $\epsilon \leq \frac{G}{20n}$. Thus, for any larger value of ϵ there is no set S^* for which the condition is
 645 satisfied.

646 *Proof of Theorem 4.3.* The proof proceeds by constructing a hard distribution of instances based on
 647 linear functions and applying information-theoretic lower bounds. We break the argument into four
 648 steps.

649 **Step 1: Definition of hard instance.** Consider the 1-dimensional domain $\mathcal{X} = [-1, 1]$. We define
 650 the component functions as linear functions:

$$\forall i \in [n] : \quad f_i(x) = -b_i x,$$

651 where $b_i \in \mathbb{R}$ is a slope parameter. The global objective is $F(x) = -x \cdot \bar{b}$, where $\bar{b} = \frac{1}{n} \sum_{i=1}^n b_i$.
 652 The minimizer x^* depends entirely on the sign of the average slope \bar{b} ; namely, $x^* = \text{sign}(\bar{b})$. The
 653 value of the minimum is $F(x^*) = -|\bar{b}|$, and for any candidate $\hat{x} \in [-1, 1]$, the optimality gap is
 654 $F(\hat{x}) - F(x^*) = \bar{b}(x^* - \hat{x})$. In particular, for any algorithm producing \hat{x} ,

$$\mathbb{E}[F(\hat{x}) - F(x^*)] \geq \bar{b} \cdot \mathbb{P}[\text{sign}(\hat{x}) \neq \text{sign}(\bar{b})]. \quad (14)$$

655 We define the hard distribution \mathcal{D} as the uniform mixture of two priors, \mathcal{D}_0 and \mathcal{D}_1 , parameterized
 656 by $\Delta_1, \dots, \Delta_n \geq 0$:

- 657 • Under \mathcal{D}_0 : $b_i = \text{clip}(z_i)$, where $z_i \sim \mathcal{N}(+\Delta_i, \sigma^2)$ are i.i.d., $\sigma = G/2$;
- 658 • Under \mathcal{D}_1 : $b_i = \text{clip}(z_i)$, where $z_i \sim \mathcal{N}(-\Delta_i, \sigma^2)$ are i.i.d., $\sigma = G/2$.

659 Here, $\text{clip}(z_i) = \max(-G, \min(G, z_i))$ is a truncation to $[-G, G]$ that strictly enforces the Lipschitz
660 condition (deterministically). The b_i are then distributed according to a truncated Gaussian distribu-
661 tion with variance σ^2 strictly smaller than G^2 . For reasons that will become apparent in later steps,
662 we will set the parameters Δ_i so that the following conditions are met:

$$\frac{1}{n} \sum_{i=1}^n \Delta_i = 20\epsilon, \quad \text{and} \quad \Delta_i \leq \frac{G}{2} \quad \text{for all } i \in [n]. \quad (15)$$

663 **Step 2: Reduction to hypothesis testing.** Define the hypothesis testing error probability as

$$p_{err} = \frac{1}{2}p_{err}^0 + \frac{1}{2}p_{err}^1, \quad \text{where} \quad p_{err}^0 = \mathbb{P}_{\mathcal{D}_0}(\hat{x} < 0) \quad \text{and} \quad p_{err}^1 = \mathbb{P}_{\mathcal{D}_1}(\hat{x} > 0).$$

664 We will next show that the expected error (Eq. (14)) is lower bounded in terms of p_{err} .

665 Conditioning on \mathcal{D}_0 , a sufficient condition for the algorithm to fail is if $\hat{x} \leq 0$ and $\bar{b} > 0$, in which
666 case the optimization error is at least $|\bar{b}|$; thus

$$\mathbb{E}_{\mathcal{D}_0}[F(\hat{x}) - F(x^*)] \geq \mathbb{E}_{\mathcal{D}_0}[|\bar{b}| \cdot \mathbb{1}(\hat{x} \leq 0) \cdot \mathbb{1}(\bar{b} > 0)]. \quad (16)$$

667 We next bound the expected value (under \mathcal{D}_0) of the mean of truncate variables. Let $g(\delta) =$
668 $\mathbb{E}[\text{clip}(Z + \delta, -G, G)]$ where $Z \sim \mathcal{N}(0, \sigma^2)$. The mean of the clipped variable is $\mathbb{E}[b_i] = g(\Delta_i)$. By
669 the mean value theorem, $g(\Delta_i) = g(0) + g'(\xi)\Delta_i = g'(\xi)\Delta_i$ for some $\xi \in (0, \Delta_i)$. The derivative is
670 $g'(t) = \mathbb{P}(-G \leq Z + t \leq G)$. With $\sigma = G/2$ and $\Delta_i \leq G/2$, the interval $[-G - \Delta_i, G - \Delta_i]$ always
671 contains $[-G/2, G/2] = [-\sigma, \sigma]$. The Gaussian mass in this interval is > 0.5 , thus $\mathbb{E}[b_i] \geq 0.5\Delta_i$.
672 The aggregate mean then satisfies $\mathbb{E}[\bar{b}] \geq 0.5(\frac{1}{n} \sum_{i=1}^n \Delta_i) = 10\epsilon$.

673 The truncated variables b_i are independent and bounded in $[-G, G]$. By Hoeffding's inequality,
674 \bar{b} concentrates around its expectation: defining the event $A = \{\bar{b} \geq 5\epsilon\}$, the probability of the
675 complement $\delta' = \mathbb{P}_{\mathcal{D}_0}(A^c)$ is bounded by

$$\delta' \leq \mathbb{P}_{\mathcal{D}_0}(\bar{b} \leq \mathbb{E}_{\mathcal{D}_0}[\bar{b}] - 5\epsilon) \leq \exp\left(-\frac{2n^2(5\epsilon)^2}{n(2L)^2}\right) \leq \exp\left(-\frac{3n\epsilon^2}{G^2}\right) \leq 0.05,$$

676 where the final inequality is true since we assume $n \geq G^2/\epsilon^2$. Conditioning on the event A (where
677 $\bar{b} \geq 5\epsilon > 0$) and recalling Eq. (16), we have

$$\begin{aligned} \mathbb{E}_{\mathcal{D}_0}[F(\hat{x}) - F(x^*)] &\geq \mathbb{E}_{\mathcal{D}_0}[|\bar{b}| \cdot \mathbb{1}(\hat{x} \leq 0) \cdot \mathbb{1}(A)] \\ &\geq 5\epsilon \mathbb{P}_{\mathcal{D}_0}(\{\hat{x} \leq 0\} \cap A) \\ &\geq 5\epsilon(\mathbb{P}_{\mathcal{D}_0}(\hat{x} \leq 0) - \mathbb{P}_{\mathcal{D}_0}(A^c)) \\ &= 5\epsilon(p_{err}^0 - \delta'). \end{aligned}$$

678 Similarly, $\mathbb{E}_{\mathcal{D}_1}[F(\hat{x}) - F(x^*)] \geq 5\epsilon(p_{err}^1 - \delta')$, thus $\mathbb{E}_{\mathcal{D}}[F(\hat{x}) - F(x^*)] \geq 5\epsilon(p_{err} - \delta')$. The
679 guarantee $\mathbb{E}[F(\hat{x}) - F(x^*)] \leq \epsilon$ implies $5\epsilon(p_{err} - 0.05) \leq \epsilon$, which leads to $p_{err} \leq 0.2 + 0.05 = 0.25$.
680 Thus, the algorithm must identify the correct hypothesis with probability at least 0.75.

681 **Step 3: Information-theoretic lower bound.** Let P and Q denote the probability distributions of
682 the entire sequence of oracle responses observed by the algorithm when the instance is drawn from
683 \mathcal{D}_0 and \mathcal{D}_1 , respectively. Let P' and Q' denote the distributions of the underlying unclipped Gaussian
684 vectors $z = (z_1, \dots, z_n)$ under \mathcal{D}_0 and \mathcal{D}_1 . Since the observed gradients $-b_i$ are deterministic
685 functions of z_i (specifically, $b_i = \text{clip}(z_i)$), the data processing inequality for KL divergences implies:

$$D_{KL}(P \parallel Q) \leq D_{KL}(P' \parallel Q').$$

686 The variables z_i are independent, and $z_i \sim \mathcal{N}(\Delta_i, \sigma^2)$ under \mathcal{D}_0 , while $z_i \sim \mathcal{N}(-\Delta_i, \sigma^2)$ under \mathcal{D}_1 .
687 The KL divergence between two univariate Gaussians with means μ_1, μ_2 and common variance σ^2 is
688 $\frac{1}{2}(\mu_1 - \mu_2)^2/\sigma^2$. Here, the separation is $2\Delta_i$ and $\sigma^2 = G^2/4$, so the KL divergence for z_i is upper
689 bounded by $8\Delta_i^2/G^2$. Since the gradient oracles are exact and the objectives are linear, repeated
690 queries to the same component yield identical results. Thus, information is gained only from the set

691 of distinct components queried. Let q_i be the probability (under \mathcal{D}_0) that component i is queried at
 692 least once by the algorithm. By the chain rule for KL divergence (and linearity of expectation), the
 693 total divergence is bounded by:

$$D_{KL}(P \parallel Q) \leq D_{KL}(P' \parallel Q') \leq \sum_{i=1}^n q_i \frac{8\Delta_i^2}{G^2}. \quad (17)$$

694 To relate the divergence to the error probability p_{err} , we use Le Cam's argument, which states that
 695 for any hypothesis testing problem with uniform prior, $p_{err} \geq \frac{1}{2}(1 - TV(P, Q))$, where $TV(P, Q)$
 696 is the total variation distance. Given our requirement that $p_{err} \leq \frac{1}{4}$, we must have:

$$\frac{1}{4} \geq \frac{1}{2}(1 - TV(P, Q)) \implies TV(P, Q) \geq \frac{1}{2}.$$

697 By Pinsker's inequality, $TV(P, Q) \leq \sqrt{\frac{1}{2}D_{KL}(P \parallel Q)}$, which yields the required lower bound on the
 698 KL divergence:

$$\frac{1}{2} \leq \sqrt{\frac{1}{2}D_{KL}(P \parallel Q)} \implies D_{KL}(P \parallel Q) \geq \frac{1}{2}.$$

699 Substituting the upper bound from Eq. (17) yields the condition:

$$\sum_{i=1}^n q_i \Delta_i^2 \geq \frac{G^2}{16}. \quad (18)$$

700 **Step 4: Optimizing the lower bound.** From Eq. (18), any algorithm solving the problem must
 701 satisfy the information constraint $\sum_{i=1}^n q_i \Delta_i^2 \geq G^2/16$. Since the total expected query cost C (under
 702 the mixture distribution \mathcal{D}) is lower bounded by half times the expected cost under \mathcal{D}_0 , we also have
 703 $C \geq \frac{1}{2} \sum_{i=1}^n q_i c_i$. For a fixed set of parameters Δ_i , we can lower bound the algorithm's expected cost
 704 in terms of the "information density" $\rho_i = \Delta_i^2/c_i$ of the components; indeed,

$$\frac{G^2}{16} \leq \sum_{i=1}^n q_i \Delta_i^2 = \sum_{i=1}^n (q_i c_i) \rho_i \leq \left(\sum_{i=1}^n q_i c_i \right) \max_i \rho_i \implies C \geq \frac{G^2}{32 \max_i \rho_i}.$$

705 To maximize this lower bound, the adversary minimizes $\max_i \rho_i$ subject to the constraints $\frac{1}{n} \sum_{i=1}^n \Delta_i =$
 706 20ϵ and $\Delta_i \leq G/2$ for all i (recall Eq. (15)). The minimum is achieved when the ratios ρ_i are
 707 equalized across the active components. Thus, we set $\Delta_i = \alpha \sqrt{c_i}$ for components i in an "active"
 708 subset $S^* \subseteq [n]$, and $\Delta_i = 0$ otherwise. The scaling factor α is determined using the constraint
 709 $\frac{1}{n} \sum_{i=1}^n \Delta_i = 20\epsilon$:

$$\frac{1}{n} \sum_{i \in S^*} (\alpha \sqrt{c_i}) = 20\epsilon \implies \alpha = \frac{20n\epsilon}{\sum_{i \in S^*} \sqrt{c_i}}. \quad (19)$$

710 To certify the additional condition $\Delta_i \leq G/2$ for all i , we require for all $i \in S^*$:

$$\alpha \sqrt{c_i} \leq \frac{G}{2} \iff \frac{20n\epsilon \sqrt{c_i}}{\sum_{j \in S^*} \sqrt{c_j}} \leq \frac{G}{2} \iff \frac{\max_{i \in S^*} \sqrt{c_i}}{\frac{1}{n} \sum_{j \in S^*} \sqrt{c_j}} \leq \frac{G}{40\epsilon}.$$

711 This is the precisely the uniformity condition given in the theorem statement. The subset S^* is defined
 712 as the largest set of cheapest components satisfying this condition. Finally, substituting α back into
 713 the cost lower bound, we conclude:

$$C \geq \frac{G^2}{32} \left(\frac{\sum_{i \in S^*} \sqrt{c_i}}{20n\epsilon} \right)^2 = \frac{G^2}{12800\epsilon^2} \left(\frac{1}{n} \sum_{i \in S^*} \sqrt{c_i} \right)^2.$$

714

□

715 F Extensions and Generality of the Theoretical Analysis

716 In this section, we discuss the generality of the theory in Sections 4 and C. In particular, for the
717 main paper we assume convexity in order to derive meaningful theoretical statements. Given that
718 our empirical setting does not optimize a convex loss function, we address the differences here.
719 Importantly, we emphasize that our theory does not strictly rely on convexity.

720 F.1 Differences between Theory and Practice

721 In the GRPO experiments, our training objective is not convex. Additionally, we do not use the
722 averaged iterate but instead the last iterate for reporting final error. Furthermore, in GRPO, the
723 objective changes every time a new set of rollouts is collected. As such, one can think of our
724 implementation as many rounds of cost-aware learning. With every new set of rollouts, we define a
725 new sampling distribution over the newly generated pool. The policy is then optimized by iterating
726 over mini-batches sampled from this reweighted pool, ensuring the cost-aware distribution actively
727 shapes the variance of the gradient updates at each step.

728 F.2 Generality of the Cost Objective in Non-Convex Settings

729 The mathematical optimality of p^* is not strictly limited to the convex setting. Its derivation depends
730 solely on the relationship between iteration complexity (T) and estimator variance ($S(p)$). Across
731 various smoothness conditions, the number of iterations T required to reach an ϵ -stationary point is
732 fundamentally bottlenecked by the stochastic gradient variance. Specifically:

- 733 • **Non-Convex SGD:** Reaching an ϵ -stationary point requires an iteration complexity of
734 $T \propto \frac{S(p)}{\epsilon^2}$ [Ghadimi and Lan, 2013].
- 735 • **Adaptive Optimizers (Adam):** Recent theoretical analyses of Adam-type algorithms
736 demonstrate a convergence rate of $\tilde{O}\left(\frac{\sqrt{S(p)}}{\sqrt{T}}\right)$ [Chen et al., 2018], implying $T \propto \frac{S(p)}{\epsilon^2}$
737 (modulo sub-polynomial logarithmic factors).
- 738 • **Generalized Smoothness:** LLM loss landscapes exhibit generalized smoothness, where
739 local curvature grows proportionally with the gradient norm. Optimization theory establishes
740 that gradient clipping and adaptive methods guarantee convergence in this regime [Zhang
741 et al., 2020]. Crucially, this theory establishes that the dominant stochastic term in the
742 iteration complexity explicitly retains the $O\left(\frac{S(p)}{\epsilon^2}\right)$ dependence on the variance $S(p)$.

Because the total expected training cost \mathcal{K} is the product of iterations and the expected per-step cost ($T \cdot C(p)$), the objective to minimize reduces to:

$$J(p) = S(p)C(p)$$

743 Since this objective remains structurally identical across these diverse optimization settings, minimiz-
744 ing the variance-cost product via p^* continues to be the precise optimal strategy. Thus, the optimal
745 distribution $p_i^* \propto G_i/\sqrt{c_i}$ derived in Theorem 4.1 remains justified in non-convex regimes.

746 F.3 Average Iterate (\bar{x}_T) vs. Last Iterate (x_T)

747 While Algorithm 1 returns the average iterate \bar{x}_T following standard conventions in the SGD literature,
748 the properties of the last iterate x_T have also been extensively studied [Shamir and Zhang, 2013,
749 Jain et al., 2019, Harvey et al., 2019]. As shown by Harvey et al. [2019], for a fixed step size, the
750 last-iterate of SGD is a factor of $\log(T)$ worse than the average. While the average iterate generally
751 achieves superior performance—sometimes strictly so—specific step-size modification schemes
752 [Jain et al., 2019] allow the expected suboptimality of the last iterate to enjoy the same convergence
753 guarantees. Consequently, sampling with the optimal distribution p^* derived in Theorem 4.1 provides
754 similar theoretical guarantees for the last iterate as we establish for the averaged iterate.

755 **G Experimental Details and Ablations**

756 **G.1 Experimental Details**

757 To ensure the reproducibility of our experiments, we provide a comprehensive overview of the
 758 hyperparameter configurations and the computational framework employed. Our implementation
 759 is built upon the Verl framework, utilizing vLLM [Kwon et al., 2023] for efficient inference and
 760 Flash-Attention [Dao et al., 2022] for optimized memory throughput.

761 **G.2 Implementation**

762 The primary modifications to the Verl codebase involve the construction of mini-batches for the
 763 policy gradient update and the integration of importance sampling. To maintain training stability, we
 764 re-center the importance weights around 1. This normalization allows us to maintain a consistent
 765 gradient norm across experiments. Crucially, this simply multiplies the training objective by a
 766 constant.

767 **G.3 Importance Weighting with GRPO**

Let $P(Q)$ be the distribution over prompts and $\pi_{\theta_{\text{old}}}(q)$ be the distribution over responses from policy $\pi_{\theta_{\text{old}}}$ given prompt q . The objective contribution of a prompt q and response o is as follows:

$$\ell(o; \theta, q) = \frac{1}{|o|} \sum_{t=1}^{|o|} \min \left(\frac{\pi_{\theta}(o_t | q, o_{<t})}{\pi_{\theta_{\text{old}}}(o_t | q, o_{<t})} A^{(t)}, \text{clip} \left(\frac{\pi_{\theta}(o_t | q, o_{<t})}{\pi_{\theta_{\text{old}}}(o_t | q, o_{<t})}, 1 - \epsilon, 1 + \epsilon \right) A^{(t)} \right) - \beta \mathbb{D}_{\text{KL}}(\pi_{\theta} \| \pi_{\text{ref}}).$$

768 The GRPO objective for parameters θ is:

$$\mathcal{J}_{\text{GRPO}}(\theta) = \mathbb{E}_{\substack{q \sim P(Q) \\ \{o_i\}_{i=1}^M \sim \pi_{\theta_{\text{old}}}(q)}} \left[\frac{1}{M} \sum_{g=1}^M \ell(o_i; \theta) \right]. \quad (20)$$

769 Given n prompts and M rollouts per prompt, and letting $o_{i,g}$ be the g 'th rollout from the i 'th prompt,
 770 the empirical objective is:

$$\hat{\mathcal{J}}_{\text{GRPO}}(\theta) = \frac{1}{nM} \sum_{i=1}^n \sum_{g=1}^M \ell(o_{i,g}; \theta, q_i). \quad (21)$$

771 In effect, the empirical GRPO objective is a finite sum objective over prompts and rollouts. If we
 772 sample according to another distribution $p \in \Delta^{nM}$, where we sample rollout pair $(q_i, o_{i,g})$ with
 773 probability $p_{(i,g)}$, the importance-weighted objective is:

$$\hat{\mathcal{J}}_{\text{GRPO}}^{(p)}(\theta) = \frac{1}{nM} \sum_{i=1}^n \sum_{g=1}^M \frac{1}{p_{(i,g)}} \ell(o_{i,g}; \theta, q_i). \quad (22)$$

774 **G.4 Hyperparameters**

775 Table 2 summarizes the key hyperparameters used for both the GRPO baseline and our GRPO + ZVF
 776 method for the 1.5B experiments. For the Qwen3-8B experiments, we detail hyperparameters in the
 777 Table 4. We report Scale-RL hyperparameters in Table 5, largely influenced by the findings of [Khatri
 778 et al., 2025].

Hyperparameter	Value
Base Model	Qwen/Qwen2.5-Math-1.5B-Instruct
Algorithm	GRPO
Learning Rate	1×10^{-6}
Weight Decay	0.1
Train Batch Size	256
Rollout Samples per Prompt (n)	16
Rollout Temperature	1.0
PPO Mini-Batch Size	64
KL Loss Coefficient	0.001
Clip Ratio Low	0.2
Clip Ratio High	0.28

Table 2: Key Training Hyperparameters for 1.5B training, GRPO and GRPO + ZVF.

Hyperparameter	Value
Base Model	Qwen/Qwen3-4B
Algorithm	GRPO
Learning Rate	1×10^{-6}
Total Epochs	15
Train Batch Size	512
PPO Mini-Batch Size	256
Rollout Samples per Prompt (n)	8
Rollout Temperature	1.0
Max Prompt Length	1024
Max Response Length	4096
KL Loss Coefficient	0.001

Table 3: Key Training Hyperparameters for Qwen3-4B Training.

Hyperparameter	Value
Base Model	Qwen/Qwen3-8B
Algorithm	GRPO
Learning Rate	1×10^{-6}
Total Epochs	15
Train Batch Size	512
PPO Mini-Batch Size	256
Rollout Samples per Prompt (n)	5
Rollout Temperature	1.0
Max Prompt Length	1024
Max Response Length	4096
KL Loss Coefficient	0.001

Table 4: Key Training Hyperparameters for Qwen3-8B Training.

Hyperparameter	Value
Base Model	Qwen/Qwen2.5-Math-1.5B-Instruct
Advantage Estimator	GRPO
Policy Loss Type	CISPO
Learning Rate	5×10^{-7}
LR Warmup Steps	100
Weight Decay	0.01
Train Batch Size	512
PPO Mini-Batch Size	128
Gradient Clipping	1.0
KL Loss Coefficient	0.001
Entropy Coefficient	0
Clip Ratio Low	0.0
Clip Ratio High	5.0
Clip Ratio C	3.0
LM Head Precision	FP32

Table 5: Key Training Hyperparameters for ScaleRL/CISPO Optimized Script

779 G.5 Computational Requirements

780 All experiments were conducted using the hardware configurations standard to the Verl framework.
781 The training and evaluation durations are as follows:

- 782 • **GRPO**: Approximately 48 hours.
- 783 • **GRPO + ZVF**: Approximately 72 hours.
- 784 • **Scale RL**: Approximately 72 hours.

785 For the Qwen3-4B experiments, we used 8 NVIDIA H100 GPUs per training run. The training runs
786 required around 2 days each.

787 For the Qwen3-8B experiments, we used 8 NVIDIA H100 GPUs per training run. The training runs
788 required around 5 days each.

789 **G.6 Ablation on variants of p^***

790 In this section we ablate the different variants of p^* . In particular, we also report p^*_{smooth} for
 791 $\alpha = 0.05, \alpha = 0.1$. Furthermore, we consider p^*_{LEN} , in which $(p^*_{\text{LEN}})_i = \frac{1}{\sqrt{c_i}}$. This sampling
 792 distribution only considers the length of sequences in its distribution. We report these variants both
 793 for GRPO and GRPO+ZVF. In particular with the GRPO setting, we see that p^*_{LEN} does not attain very
 794 strong performance. We hypothesize that this is because many shorter responses have 0 advantage,
 795 and hence there is little training signal. As in the main paper, we also report final performance for all
 796 methods.

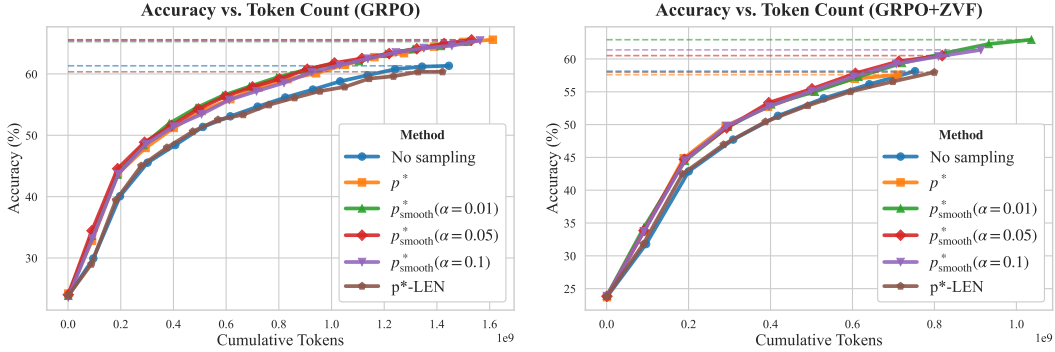


Figure 7: Cumulative tokens compared with AIME pass@1/mean@32 accuracy throughout training for both GRPO and GRPO+ZVF methods.

Setting	Method	AIME	AMC	MATH500	GSM8K	Avg. Accuracy
GRPO	No sampling	61.3	64.1	73.2	86.2	71.2
GRPO	p^*	65.6	71.2	73.2	86.0	74.0
GRPO	$p^*_{\text{smooth}}(\alpha = 0.01)$	65.3	68.1	72.3	86.0	72.9
GRPO	$p^*_{\text{smooth}}(\alpha = 0.05)$	65.6	68.3	72.8	86.1	73.2
GRPO	$p^*_{\text{smooth}}(\alpha = 0.1)$	65.4	66.0	72.6	85.8	72.5
GRPO	p^*_{LEN}	60.3	64.8	73.0	86.0	71.0
ZVF	No sampling	58.1	64.7	73.2	86.1	70.6
ZVF	p^*	57.6	63.0	72.9	85.9	69.8
ZVF	$p^*_{\text{smooth}}(\alpha = 0.01)$	62.9	68.0	73.3	85.8	72.5
ZVF	$p^*_{\text{smooth}}(\alpha = 0.05)$	60.5	65.9	73.0	86.4	71.4
ZVF	$p^*_{\text{smooth}}(\alpha = 0.1)$	61.4	68.3	72.9	86.0	72.2
ZVF	p^*_{LEN}	58.0	63.3	73.9	86.2	70.3

Table 6: (Math-1.5B-Instruct, pass@1/mean@32): Complete results.

Setting	Method	AIME	AMC	MATH500	GSM8K	Avg. Accuracy
GRPO	No sampling	58.3	72.8	87.0	96.0	78.5
GRPO	p^*	58.5	74.7	86.7	95.2	78.8
GRPO	$p^*_{\text{smooth}}(\alpha = 0.01)$	58.5	74.0	86.7	94.9	78.5
GRPO	$p^*_{\text{smooth}}(\alpha = 0.05)$	58.3	73.6	86.6	95.6	78.5

Table 7: (8B-Base, pass@1/mean@32): Complete results.

797 **G.7 Full CISPO results**

798 We report performance across benchmarks. The CISPO curves are reported in Figure 8 and the best
 799 checkpoint performance is reported in Table 8.

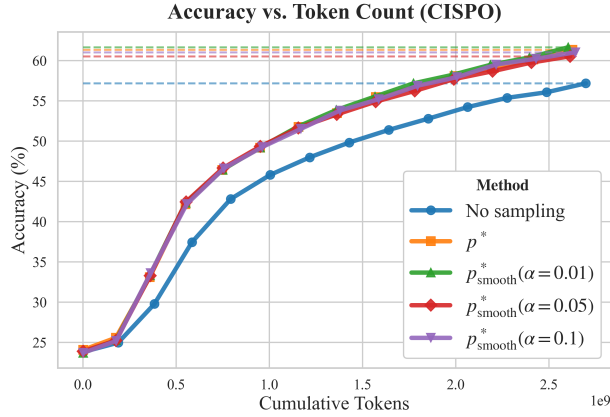


Figure 8: Full CISPO objective results for Qwen2.5-Math-1.5B-Instruct on AIME. We see robustness to smoothed cost-aware learning.

Setting	Method	AIME	AMC	MATH500	GSM8K	Avg. Accuracy
GRPO	No sampling	57.2	61.8	73.0	85.9	69.5
GRPO	p^*	61.3	65.7	73.6	86.2	71.7
GRPO	p_{smooth}^* ($\alpha = 0.01$)	61.7	65.4	73.5	86.0	71.6
GRPO	p_{smooth}^* ($\alpha = 0.05$)	60.5	66.0	73.0	86.1	71.4
GRPO	p_{smooth}^* ($\alpha = 0.1$)	61.0	65.2	73.6	86.0	71.5

Table 8: (Math-1.5B-Instruct, pass@1/mean@32): Complete results for CISPO.

800 **G.8 Full sub-optimality results**

801 In the main paper, we report the sub-optimality metrics for the 8B and 1.5B models under the GRPO
 802 setting. Here, we also include the sub-optimality metrics for the GRPO+ZVF setting and the CISPO
 803 settings in Figures 9 and 10 respectively.

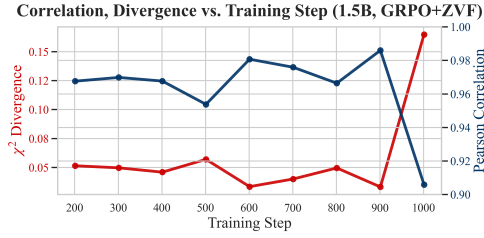


Figure 9: Sub-optimality metrics for 1.5B GRPO+ZVF training run.

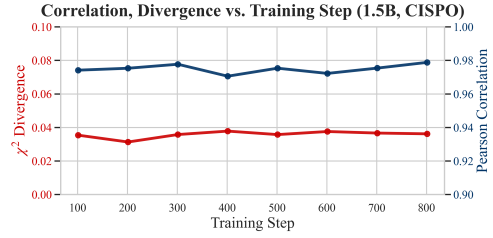


Figure 10: Sub-optimality metrics for 1.5B CISPO training objective run.

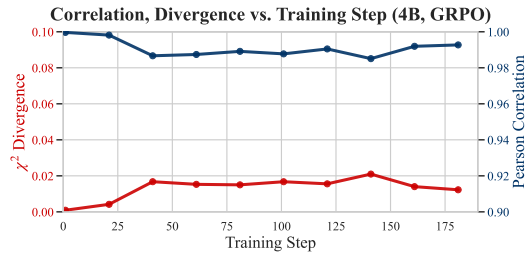


Figure 11: Sub-optimality metrics for 4B run.

804 **G.9 Additional results for Qwen2.5-Math-1.5B-Instruct**

We report AMC performance as a function of cumulative tokens in Figure 12.

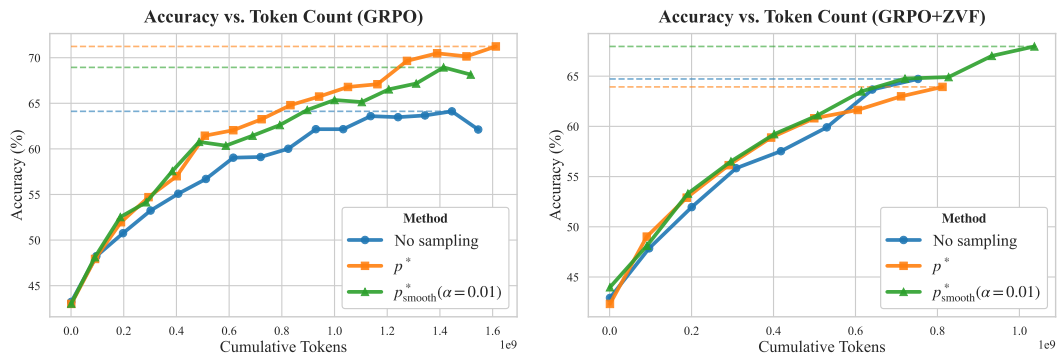


Figure 12: Cumulative tokens compared with AMC pass@1/mean@32 accuracy throughout training for both GRPO and GRPO+ZVF settings for the 1.5B model.

805

806 **G.10 Additional results for Qwen3-4B**

807 We report results for Qwen3-4B on AIME and AMC below. We plot up until peak accuracy is attained.
 808 We also report smoothed variants.

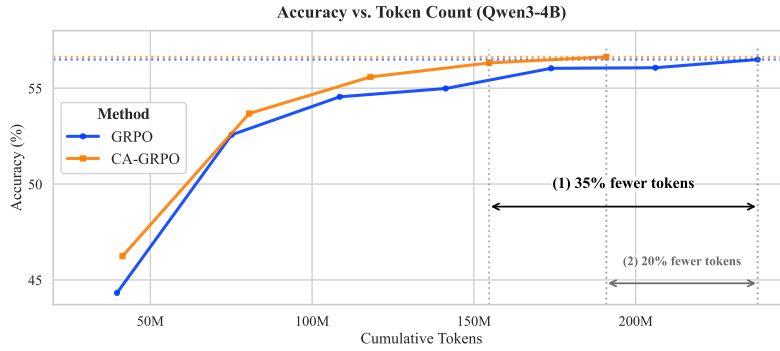


Figure 13: Qwen3-4B on AIME.

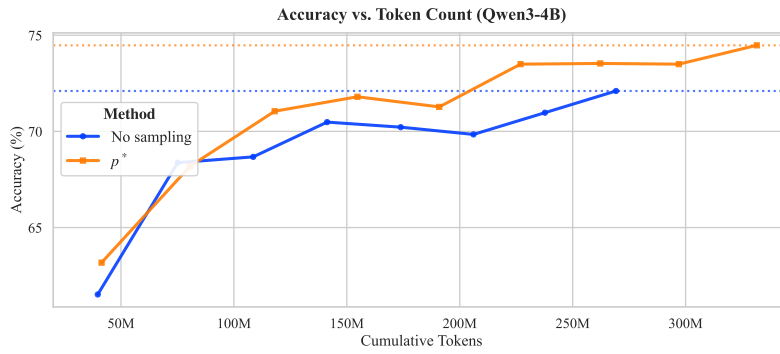


Figure 14: Qwen3-4B on AMC.

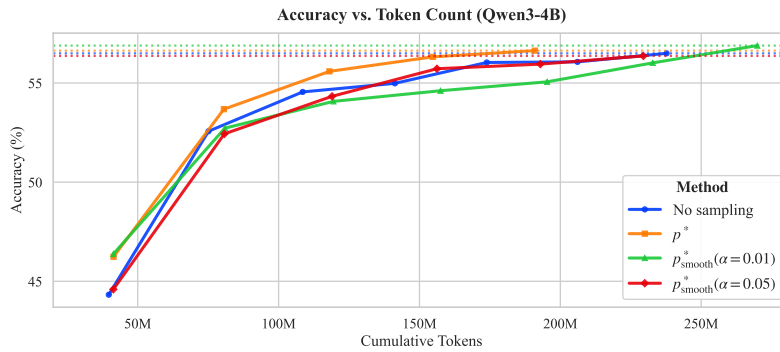


Figure 15: Qwen3-4B and smoothed variants on AIME.

809 **G.11 Additional results for Qwen3-8B**

810 We present the results for the smoothed variants of the 8B model on AIME in Figure 16. Additionally,
 811 we plot AMC performance as a function of cumulative tokens in Figure 17. In Table 7, we report full
 812 results on best checkpoint accuracy for smoothed variants on four benchmarks.

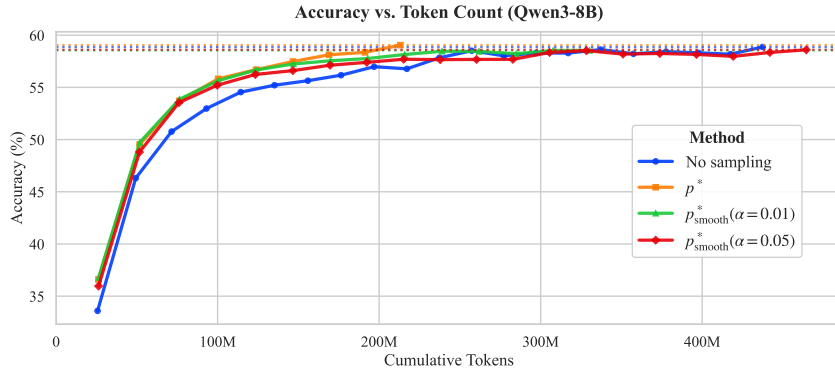


Figure 16: Qwen3-8B AIME results for all variants.

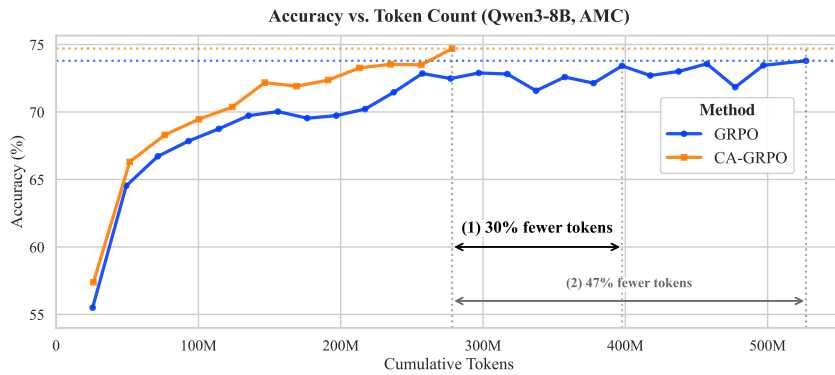


Figure 17: AMC results for Qwen3-8B Base.

813 **NeurIPS Paper Checklist**

814 **1. Claims**

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

817 Answer: [\[Yes\]](#)

818 Justification: Claims are supported empirically and theoretically.

819 Guidelines:

- 820 • The answer [\[N/A\]](#) means that the abstract and introduction do not include the claims
821 made in the paper.
- 822 • The abstract and/or introduction should clearly state the claims made, including the
823 contributions made in the paper and important assumptions and limitations. A [\[No\]](#) or
824 [\[N/A\]](#) answer to this question will not be perceived well by the reviewers.
- 825 • The claims made should match theoretical and experimental results, and reflect how
826 much the results can be expected to generalize to other settings.
- 827 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
828 are not attained by the paper.

829 **2. Limitations**

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

831 Answer: [\[Yes\]](#)

832 Justification: Throughout the work, we discuss limitations.

833 Guidelines:

- 834 • The answer [\[N/A\]](#) means that the paper has no limitation while the answer [\[No\]](#) means
835 that the paper has limitations, but those are not discussed in the paper.
- 836 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 837 • The paper should point out any strong assumptions and how robust the results are to
838 violations of these assumptions (e.g., independence assumptions, noiseless settings,
839 model well-specification, asymptotic approximations only holding locally). The authors
840 should reflect on how these assumptions might be violated in practice and what the
841 implications would be.
- 842 • The authors should reflect on the scope of the claims made, e.g., if the approach was
843 only tested on a few datasets or with a few runs. In general, empirical results often
844 depend on implicit assumptions, which should be articulated.
- 845 • The authors should reflect on the factors that influence the performance of the approach.
846 For example, a facial recognition algorithm may perform poorly when image resolution
847 is low or images are taken in low lighting. Or a speech-to-text system might not be
848 used reliably to provide closed captions for online lectures because it fails to handle
849 technical jargon.
- 850 • The authors should discuss the computational efficiency of the proposed algorithms
851 and how they scale with dataset size.
- 852 • If applicable, the authors should discuss possible limitations of their approach to
853 address problems of privacy and fairness.
- 854 • While the authors might fear that complete honesty about limitations might be used by
855 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
856 limitations that aren’t acknowledged in the paper. The authors should use their best
857 judgment and recognize that individual actions in favor of transparency play an impor-
858 tant role in developing norms that preserve the integrity of the community. Reviewers
859 will be specifically instructed to not penalize honesty concerning limitations.

860 **3. Theory assumptions and proofs**

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

863 Answer: [\[Yes\]](#)

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

Justification: All assumptions are stated in theorem statements and problem descriptions.

Guidelines:

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

4. Experimental result reproducibility

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

Answer: [Yes]

Justification: Empirical details, pseudocode, and hyperparameters are provided.

Guidelines:

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

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968

Answer: [No]

Justification: We mention specifically which training datasets and models we use. We provide pseudocode as well as the library on top of which the algorithm is implemented rather than code.

Guidelines:

- The answer [N/A] means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

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

Answer: [Yes]

Justification: Hyperparameters and training details are discussed in Section 2 and Appendix G.

Guidelines:

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

7. Experiment statistical significance

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

Answer: [Yes]

Justification: For accuracy, we report mean@32 which is an average of prompting the model 32 times. Additionally, we report careful ablations showing the significance of noise in our experiments. For synthetic results, we report error bars.

Guidelines:

- The answer [N/A] means that the paper does not include experiments.
- The authors should answer [Yes] if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).

- 969 • The method for calculating the error bars should be explained (closed form formula,
970 call to a library function, bootstrap, etc.)
- 971 • The assumptions made should be given (e.g., Normally distributed errors).
- 972 • It should be clear whether the error bar is the standard deviation or the standard error
973 of the mean.
- 974 • It is OK to report 1-sigma error bars, but one should state it. The authors should
975 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
976 of Normality of errors is not verified.
- 977 • For asymmetric distributions, the authors should be careful not to show in tables or
978 figures symmetric error bars that would yield results that are out of range (e.g., negative
979 error rates).
- 980 • If error bars are reported in tables or plots, the authors should explain in the text how
981 they were calculated and reference the corresponding figures or tables in the text.

982 8. Experiments compute resources

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

986 Answer: [Yes]

987 Justification: This information is discussed in Appendix G.

988 Guidelines:

- 989 • The answer [N/A] means that the paper does not include experiments.
- 990 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
991 or cloud provider, including relevant memory and storage.
- 992 • The paper should provide the amount of compute required for each of the individual
993 experimental runs as well as estimate the total compute.
- 994 • The paper should disclose whether the full research project required more compute
995 than the experiments reported in the paper (e.g., preliminary or failed experiments that
996 didn't make it into the paper).

997 9. Code of ethics

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

1000 Answer: [Yes]

1001 Justification: We have carefully followed the ethics guidelines provided by NeurIPS.

1002 Guidelines:

- 1003 • The answer [N/A] means that the authors have not reviewed the NeurIPS Code of
1004 Ethics.
- 1005 • If the authors answer [No], they should explain the special circumstances that require a
1006 deviation from the Code of Ethics.
- 1007 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
1008 eration due to laws or regulations in their jurisdiction).

1009 10. Broader impacts

1010 Question: Does the paper discuss both potential positive societal impacts and negative
1011 societal impacts of the work performed?

1012 Answer: [N/A]

1013 Justification: We do not believe that this work has any societal impact, as it is a theoretical
1014 and algorithmic contribution.

1015 Guidelines:

- 1016 • The answer [N/A] means that there is no societal impact of the work performed.
- 1017 • If the authors answer [N/A] or [No], they should explain why their work has no societal
1018 impact or why the paper does not address societal impact.

- 1019
- 1020
- 1021
- 1022
- 1023
- 1024
- 1025
- 1026
- 1027
- 1028
- 1029
- 1030
- 1031
- 1032
- 1033
- 1034
- 1035
- 1036
- 1037
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
 - The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate Deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
 - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
 - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

1038 11. Safeguards

1039 Question: Does the paper describe safeguards that have been put in place for responsible
1040 release of data or models that have a high risk for misuse (e.g., pre-trained language models,
1041 image generators, or scraped datasets)?

1042 Answer: [N/A]

1043 Justification: The paper poses no such risks.

1044 Guidelines:

- 1045
- 1046
- 1047
- 1048
- 1049
- 1050
- 1051
- 1052
- 1053
- 1054
- The answer [N/A] means that the paper poses no such risks.
 - 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.
 - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
 - 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.

1055 12. Licenses for existing assets

1056 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
1057 the paper, properly credited and are the license and terms of use explicitly mentioned and
1058 properly respected?

1059 Answer: [Yes]

1060 Justification: All models, datasets, and frameworks used are properly cited. We use the Qwen
1061 model family (Apache 2.0 license), the DAPO training dataset, and standard evaluation
1062 benchmarks (MATH500, AMC, GSM8K, AIME1983-2024). Our implementation builds on
1063 the Verl framework with vLLM and Flash-Attention. All assets are used in accordance with
1064 their respective licenses.

1065 Guidelines:

- 1066
- 1067
- 1068
- 1069
- 1070
- 1071
- 1072
- 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.

- 1073
- 1074
- 1075
- 1076
- 1077
- 1078
- 1079
- 1080
- 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.

1081 **13. New assets**

1082 Question: Are new assets introduced in the paper well documented and is the documentation
1083 provided alongside the assets?

1084 Answer: [N/A]

1085 Justification: The paper does not release new assets.

1086 Guidelines:

- 1087
- 1088
- 1089
- 1090
- 1091
- 1092
- 1093
- 1094
- 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.

1095 **14. Crowdsourcing and research with human subjects**

1096 Question: For crowdsourcing experiments and research with human subjects, does the paper
1097 include the full text of instructions given to participants and screenshots, if applicable, as
1098 well as details about compensation (if any)?

1099 Answer: [N/A]

1100 Justification: The paper does not involve crowdsourcing nor research with human subjects.

1101 Guidelines:

- 1102
- 1103
- 1104
- 1105
- 1106
- 1107
- 1108
- 1109
- 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.

1110 **15. Institutional review board (IRB) approvals or equivalent for research with human
1111 subjects**

1112 Question: Does the paper describe potential risks incurred by study participants, whether
1113 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
1114 approvals (or an equivalent approval/review based on the requirements of your country or
1115 institution) were obtained?

1116 Answer: [N/A]

1117 Justification: The paper does not involve crowdsourcing nor research with human subjects.

1118 Guidelines:

- 1119
- 1120
- 1121
- 1122
- 1123
- 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.

- 1124 • We recognize that the procedures for this may vary significantly between institutions
1125 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
1126 guidelines for their institution.
- 1127 • For initial submissions, do not include any information that would break anonymity (if
1128 applicable), such as the institution conducting the review.

1129 **16. Declaration of LLM usage**

1130 Question: Does the paper describe the usage of LLMs if it is an important, original, or
1131 non-standard component of the core methods in this research? Note that if the LLM is used
1132 only for writing, editing, or formatting purposes and does *not* impact the core methodology,
1133 scientific rigor, or originality of the research, declaration is not required.

1134 Answer: [N/A]

1135 Justification: The core method development in this research does not involve LLMs as any
1136 important, original, or non-standard components.

1137 Guidelines:

- 1138 • The answer [N/A] means that the core method development in this research does not
1139 involve LLMs as any important, original, or non-standard components.
- 1140 • Please refer to our LLM policy in the NeurIPS handbook for what should or should not
1141 be described.