

Exploiting Hardness and Diversity for Data-Efficient Fine-Tuning

Anonymous ACL submission

Abstract

Fine-tuning large language models for mathematical reasoning is often done using large training sets, despite that many training examples are redundant once a model is already instruction-tuned. Under practical compute and time constraints, it is therefore important to understand which training examples actually matter. We investigate this problem on GSM8K by fine-tuning Gemma-2-2B-it with LoRA under a *fixed* data budget. We compare uniform random sampling with two structured data selection methods. A taxonomy-based method, **Skill-Balanced Sampling (SBS)**, enforces balanced coverage across predefined skill categories but shows only modest and inconsistent gains. We then propose **Hardness-Weighted Diversity (HWD)**, which explicitly controls the proportion of easy, medium, and hard examples while promoting semantic diversity. Our empirical results indicate clear performance saturation well before the full dataset is utilized. Moreover, HWD reaches the best performance using only 9% of the GSM8K training data, outperforming both random sampling and SBS with substantially fewer training examples.

1 Introduction

Fine-tuning large language models for mathematical reasoning is often approached as a “more data is better” problem. However, as datasets grow, so do training costs. For many practitioners operating under tight compute or time constraints, the linear scaling of training time becomes a significant bottleneck. This raises a practical but underexplored question: how much data is actually necessary, and which examples contribute most to performance gains?

We study this question from a fixed-budget perspective. When only a small fraction of the available data ($K \ll N$) can be used, the composition of the training subset becomes critical. Redundant examples waste limited capacity, while

highly skewed subsets risk over-specialization and degraded generalization.

GSM8K provides a natural testbed for this analysis. Despite covering a wide range of mathematical skills, the dataset contains substantial redundancy: many problems share similar underlying solution structure with minor surface-level changes. For an already instruction-tuned model such as Gemma-2-2B-it, fine-tuning on the full dataset may therefore be inefficient.

To address this challenge, we first show that uniform random sampling is a strong baseline and performs well when data is abundant. However, its reliability degrades as the training budget shrinks. We then show that a simple taxonomy-based approach, **Skill-Balanced Sampling (SBS)**, yields inconsistent improvements. This motivates **Hardness-Weighted Diversity (HWD)**, a novel data selection method designed to construct a single high-quality subset prior to training. HWD explicitly controls the mixture of difficulty levels in the training data while leveraging semantic embeddings to reduce redundancy.

Our results lead to three main observations. First, reasoning performance on GSM8K saturates early, with limited gains from additional training data. Second, while random sampling remains competitive at larger budgets, it becomes unreliable in low-data regimes. Finally, our proposed HWD achieves the best performance using less than 10% of the training data, outperforming both random and skill-based sampling approaches. Overall, these findings suggest that for mathematical reasoning, careful control over data difficulty and diversity coverage is more important than dataset size alone.

2 Related Work

2.1 Data Selection for Efficient Training

Selecting informative subsets for efficient training has a long history in machine learning. Clas-

sical approaches include active learning and sub-modular optimization, which aim to construct representative subsets that preserve coverage of the full dataset (Wei et al., 2015; Mirzasoaleiman et al., 2020). While effective for representativeness, these methods typically treat all selected examples as equally useful, ignoring variation in example difficulty—an issue that becomes critical under tight data budgets for reasoning tasks.

More recent work explores model-dependent signals to estimate example importance. Gradient-based influence methods trace the effect of individual training points on model predictions (Koh and Liang, 2017; Pruthi et al., 2020), while early-training loss and gradient norms have been shown to identify influential examples early in training (Paul et al., 2021). Related work studies long-tail effects and memorization through influence estimation (Feldman and Zhang, 2020). Although effective, these approaches generally require partial or repeated training runs, making them expensive for large models. In contrast, our method operates in a *static selection* regime, where the subset is chosen once prior to fine-tuning.

2.2 Training Dynamics and Difficulty

Several studies analyze training dynamics to estimate dataset difficulty or example utility. Forgetting-based analyses identify examples that are learned and later misclassified as particularly influential for generalization (Toneva et al., 2019), while information-theoretic perspectives characterize dataset difficulty through usable information (Ethayarajh et al., 2021). Recent work on data valuation explicitly studies which examples are most useful for fine-tuning large language models (Xia et al., 2024).

Despite their insights, these methods typically rely on dynamic signals obtained during training. Our work instead focuses on low-cost, static hardness signals that can be computed without repeated optimization, making them suitable for parameter-efficient fine-tuning with LoRA.

2.3 Curriculum Learning under Fixed Budget

Curriculum learning studies how ordering training examples from easy to hard can improve optimization and generalization (Bengio et al., 2009). Extensions such as self-paced learning adaptively expand the training set based on model competence (Kumar et al., 2010). Recent work demonstrates the effectiveness of difficulty-aware curricula for

mathematical reasoning with large language models (Zhou et al., 2023b; Liu et al., 2023).

Unlike curriculum learning, which assumes access to the full dataset and focuses on example ordering, our setting considers *fixed-budget selection*, where only a small subset of the data can be used. In this regime, difficulty alone is insufficient, as redundancy across similar problems can significantly reduce the effective diversity of the subset.

2.4 Redundancy and Scaling in Reasoning Datasets

Scaling laws predict diminishing returns as dataset size increases (Hestness et al., 2017; Kaplan et al., 2020), motivating careful data selection when training compute is limited. For language models, compute-optimal training further emphasizes the importance of efficient data usage (Hoffmann et al., 2022). In reasoning benchmarks such as GSM8K (Cobbe et al., 2021), structural redundancy across problems exacerbates this effect, making uniform random sampling a strong and surprisingly competitive baseline (Chung et al., 2024).

At the same time, recent work such as LIMA shows that small, carefully curated datasets can outperform much larger ones (Zhou et al., 2023a). Our results build on these insights by demonstrating that explicitly controlling both difficulty and semantic diversity enables more data-efficient fine-tuning under fixed budgets.

3 Problem Setup

Our goal is to identify effective ways to select a small training subset for mathematical reasoning tasks. Given a large dataset $\mathcal{D} = \{x_i\}_{i=1}^N$ and a budget $K \ll N$, we aim to choose a subset $\mathcal{S} \subset \mathcal{D}$ with $|\mathcal{S}| = K$ such that fine-tuning on \mathcal{S} leads to the best possible performance on the test set.

Formally, this can be written as

$$\mathcal{S}^* = \arg \max_{\mathcal{S} \subset \mathcal{D}, |\mathcal{S}|=K} \text{Acc}(\mathcal{M}_{\mathcal{S}}),$$

where $\mathcal{M}_{\mathcal{S}}$ denotes the model fine-tuned on subset \mathcal{S} . Since it is infeasible to evaluate all possible subsets of size K , we rely on proxy signals to guide subset construction.

Experimental framework. We use GSM8K as our primary benchmark and fine-tune **Gemma-2-2B-it** using Low-Rank Adaptation (LoRA). To ensure that differences in performance reflect the data

selection strategy rather than changes in optimization, we keep the model architecture and training configuration fixed across all experiments. This includes the LoRA setup, batch size, number of epochs, and learning-rate schedule. Performance is evaluated using strict pass@1 accuracy on the GSM8K test split.

Selection signals. To guide subset selection, each training example is annotated with three auxiliary signals. These signals are computed *offline* and are used only during subset construction; they are not accessed during fine-tuning or evaluation. Details on how these signals are computed are provided in Appendix A.

- **Hardness (H_i):** A scalar difficulty score in $[0, 1]$, where higher values indicate harder problems. We estimate hardness using a fixed, unfine-tuned reference model by evaluating each example over multiple independent decoding attempts. Specifically, each problem is queried M times under stochastic decoding, and H_i is defined as one minus the fraction of correct model outputs.
- **Semantic embeddings (E_i):** Vector representations of each problem statement obtained from a frozen sentence encoder. All embeddings are ℓ_2 -normalized so that cosine similarity can be used to measure redundancy and semantic novelty.
- **Skill labels ($s(i)$):** Discrete categories derived from the standard GSM8K skill taxonomy (e.g., arithmetic, geometry, proportions). Problems may be associated with multiple skills.

Evaluation protocol. We evaluate each subset selection method across multiple budget sizes K . To account for variability in low-data fine-tuning, we average results over multiple random seeds and learning rates. For each (method, K) pair, we report performance under the learning rate that achieves the highest mean accuracy across seeds. This protocol helps separate the effect of data selection from noise introduced by optimization.

4 Method

We study two structured data selection methods under a fixed training budget: a taxonomy-based baseline, **Skill-Balanced Sampling (SBS)**, and our proposed method, **Hardness-Weighted Diversity (HWD)**. SBS serves as a simple skill-coverage

baseline, while HWD explicitly balances difficulty, semantic diversity, and skill coverage.

4.1 Skill-Balanced Sampling (SBS)

SBS is a taxonomy-driven method that aims to ensure balanced coverage over a predefined set of skills. In this work, we consider a taxonomy consisting of seventeen skill categories (Appendix A.3). Each training example i is annotated with a (possibly multi-valued) skill label $s(i)$ drawn from a fixed skill taxonomy.

Per-skill targets. Let f_s denote the number of training examples associated with skill s in the full dataset. Given a budget K , SBS assigns a target count K_s to each skill s via a fixed global sampling rate ρ , chosen such that the expected total number of selected examples is approximately K . Specifically, we define

$$K_s = \begin{cases} 0, & f_s = 0, \\ \text{clip}(\lceil \rho f_s \rceil; 1, f_s), & f_s > 0. \end{cases}$$

Here $\rho \in (0, 1)$ is a fixed sampling rate (set to $\rho = 0.1$ in our experiments), and $\text{clip}(x; a, b) = \min(\max(x, a), b)$. In practice, if $\sum_s K_s \neq K$, we uniformly subsample or pad across skills to match the exact budget.

Greedy multi-cover selection. Let $C_s(S)$ denote the number of selected examples in subset S that contain skill s . Starting from $S = \emptyset$, SBS iteratively adds one example at a time by selecting

$$i^* = \arg \max_{i \notin S} \text{key}(i \mid S),$$

where the greedy key is the lexicographic tuple

$$\text{key}(i \mid S) = (\text{gain}(i), \text{rarity}(i), |s(i)|, -i).$$

Here

$$\text{gain}(i) = \sum_{s \in s(i)} \mathbf{1}_s, \quad \mathbf{1}_s = \begin{cases} 1, & C_s(S) < K_s, \\ 0, & \text{otherwise.} \end{cases}$$

Ties are broken deterministically by the smallest example index (captured by $-i$). After selecting i^* , the counts $C_s(S)$ are updated for all $s \in s(i^*)$.

Stopping criterion. The procedure terminates when either all skill targets are met ($C_s(S) \geq K_s$ for all s) or the budget $|S| = K$ is reached. SBS does not incorporate example difficulty or semantic similarity, and treats all examples within a skill category as equally informative. We use SBS as an interpretable baseline to assess the benefits of data selection based solely on skill diversity.

4.2 Hardness-Weighted Diversity (HWD)

4.2.1 Overview

Figure 1 gives an overview of the proposed HWD pipeline, showing how hardness control, semantic novelty, and soft constraints interact during subset construction for data-efficient fine-tuning. Given a budget $K \ll N$, HWD aims to select a subset S that maximizes downstream performance by balancing difficulty, coverage, and redundancy.

HWD jointly considers two forms of diversity: *semantic diversity*, enforced through embedding-based novelty to reduce redundancy, and *skill diversity*, encouraged via soft coverage constraints over discrete skill labels. Both are important in low-data settings, where redundant examples or skill collapse can severely limit generalization.

HWD builds the subset using a greedy selection with soft regularization, followed by a local refinement step. This refinement consistently improves stability by correcting early greedy decisions.

The overall HWD procedure is as follows:

1. Normalize embeddings and estimate global skill priors.
2. Track subset composition using the hardness bins.
3. Greedily construct a size- K subset by trading off difficulty and semantic novelty, while applying soft penalties for skill coverage and hardness balance.
4. Apply swap-based refinement to improve the objective without changing the subset size.

Notations. Each training example i is annotated with three auxiliary signals used only for subset selection: a hardness score $H_i \in [0, 1]$, a normalized semantic embedding $E_i \in \mathbb{R}^d$, and a discrete skill label $s(i)$. These signals are not accessed during fine-tuning. We denote by $\mathcal{U}_b = \{i : b(i) = b\}$ the set of items in hardness bin b . The method uses weights $\lambda_H, \lambda_D, \lambda_S, \lambda_{\text{Mix}}$, a slack parameter δ , and a skill margin α . Throughout the paper, *semantic diversity* refers to embedding-based novelty, while *skill diversity* denotes balanced coverage over discrete skill labels.

4.2.2 Preprocessing

We first normalize embeddings so dot products correspond to cosine similarity:

$$E_i \leftarrow \frac{E_i}{\|E_i\|_2}.$$

We also estimate the global skill prior p_s from the full dataset and define a soft per-skill target count:

$$T_s = K \cdot p_s.$$

This target is not treated as a hard constraint; instead it anchors a penalty that discourages extreme concentration in any single skill.

4.2.3 Hardness discretization

We discretize hardness levels using a set of thresholds $\tau = (\tau_0, \tau_1, \tau_2, \tau_3)$, where $0 = \tau_0 < \tau_1 < \tau_2 < \tau_3 = 1$. For example, setting $(\tau_1, \tau_2) = (0.5, 0.8)$ partitions the data into three bins corresponding to easy, medium, and hard examples, spanning hardness score ranges $[0, 0.5)$, $[0.5, 0.8)$, and $[0.8, 1]$, respectively.

Additionally, we associate each hardness bin with a target proportion \mathbf{t} . For example, setting $\mathbf{t} = [0.10, 0.60, 0.30]$ allocates 10%, 60%, and 30% of the budget K for the easy, medium, and hard bins, respectively.

4.2.4 Top- M Candidate Pool

After hardness discretization, the next step is to construct the training subset by using the hardness bins. However, performing a full greedy scan over all N items at each selection step is unnecessary and inefficient. Therefore, we restrict selection to a candidate pool of size M :

$$M = \text{clamp}(\text{multiplier} \cdot K, M_{\min}, M_{\max}),$$

with fixed hyperparameters $\text{multiplier} > 1$, M_{\min} , and M_{\max} .

The candidate pool is constructed by globally sorting all eligible items by hardness score and retaining the top M examples with successfully computed hardness scores and embeddings (i.e., excluding items with missing or undefined values). The resulting Top- M set serves only as a search-space restriction and does not impose any ordering during greedy selection.

Hardness bins are *not* used to stratify the candidate pool; they are instead used during greedy construction to compute soft penalties that discourage early dominance of any single difficulty regime.

4.2.5 Diversity-Aware Greedy Selection

We construct the subset S iteratively, starting from a single anchor item and adding one element at a time until $|S| = K$.

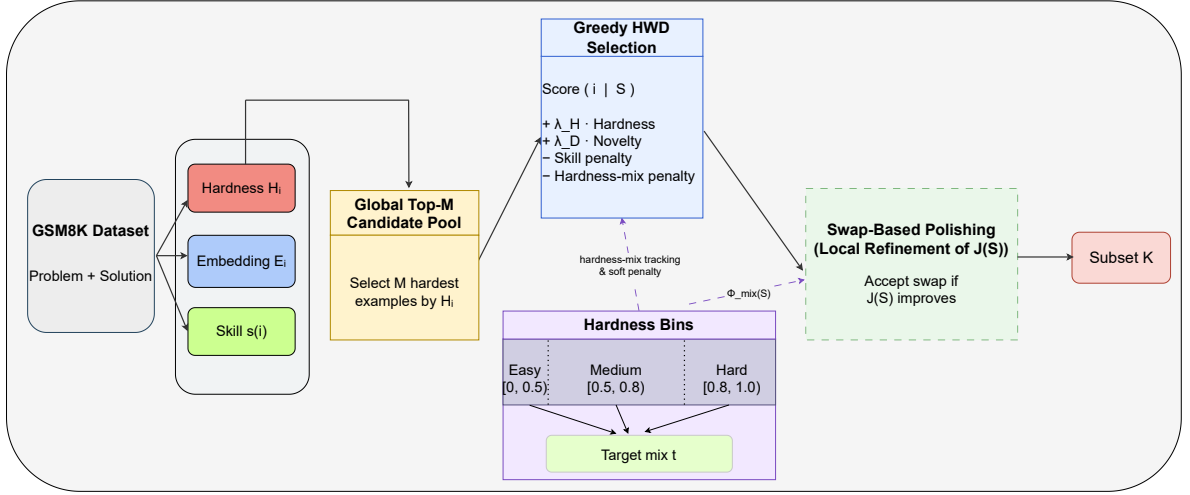


Figure 1: **The HWD Selection Pipeline.** Selection is restricted to the Top- M hardest examples, from which a size- K subset is greedily constructed by balancing hardness and semantic novelty under soft skill and difficulty penalties. An optional swap-based step further improves the objective $J(S)$.

Initialization. We initialize the selection with the single hardest item:

$$S = \left\{ \arg \max_i H_i \right\}.$$

Novelty score. At each step, every candidate item $i \notin S$ receives a novelty score relative to the current set:

$$D_i(S) = 1 - \max_{j \in S} \langle E_i, E_j \rangle.$$

This conservative max-similarity formulation penalizes the closest near-duplicate in S , encouraging coverage of distinct semantic regions.

Skill handling. Although a problem may be annotated with multiple skills, HWD assigns each example a single *primary skill* and counts only this label during greedy selection. This avoids over-weighting multi-skill examples in the objective and leads to more stable skill-coverage penalties.

Skill coverage penalty. Let C_s denote the number of selected items in S with skill s . To encourage *skill diversity* and prevent the subset from collapsing into a small number of dominant skills, we impose a soft coverage constraint based on the empirical skill prior. If selecting item i would cause the count of its skill $s(i)$ to exceed the tolerated budget $\alpha T_{s(i)}$, we apply the asymmetric penalty:

$$\text{pen}_S(i) = \frac{[C_{s(i)} + 1 - \alpha T_{s(i)}]_+}{\max(1, T_{s(i)})}.$$

Hardness-mix penalty. Let C_b denote the number of selected items in hardness bin b . At iteration $t = |S| + 1$, the expected count for bin b is $\tau_b^{(t)} = t t_b$. We define

$$\text{dev}(i) = \frac{(C_{b(i)} + 1) - \tau_{b(i)}^{(t)}}{\max(1, \tau_{b(i)}^{(t)})}.$$

Deviations within a slack band δ are tolerated; larger deviations incur a quadratic penalty:

$$\text{pen}_M(i) = [|\text{dev}(i)| - \delta]_+^2.$$

Selection score. The per-step selection score is:

$$\begin{aligned} \text{Score}(i | S) &= \lambda_H H_i + \lambda_D D_i(S) \\ &\quad - \lambda_S \text{pen}_S(i) - \lambda_{\text{Mix}} \text{pen}_M(i). \end{aligned}$$

At each iteration we select

$$i^* = \arg \max_{i \notin S} \text{Score}(i | S)$$

that maximizes the overall score. We then update S and all associated counts, and repeat until $|S| = K$.

4.2.6 Swap-Based Refinement

Greedy selection is order-dependent: once an item is chosen early, it affects all future novelty computations. To reduce sensitivity to these early decisions, we perform up to N_{swaps} random one-for-one swaps between the current set S and the candidate pool.

Leave-one-out novelty. For an item $i \in S$, we define the leave-one-out novelty:

$$D_i^{\text{LOO}}(S) = 1 - \max_{j \in S \setminus \{i\}} \langle E_i, E_j \rangle,$$

with $D_i^{\text{LOO}}(S) = 1$ when $|S| = 1$.

Global objective. Using $[x]_+ = \max(0, x)$ and $m(x) = \max(1, x)$, we define:

$$\begin{aligned} J(S) = & \lambda_H \sum_{i \in S} H_i + \lambda_D \sum_{i \in S} D_i^{\text{LOO}}(S) \\ & - \lambda_S \sum_{s \in S} \frac{[C_s(S) - \alpha T_s]_+}{m(T_s)} \\ & - \lambda_{\text{Mix}} \Phi_{\text{mix}}(S). \end{aligned}$$

Mix penalty term. Let i_t be the t -th selected item and $C_b^{(t)}$ the number of items from bin b after t steps. With $\tau_b^{(t)} = t t_b$, we penalize overshoot beyond a slack band:

$$\phi_{\text{over}}(c, \tau) = \left(\frac{[c - (1 + \delta)\tau]_+}{\max(1, (1 + \delta)\tau)} \right)^2,$$

and accumulate:

$$\Phi_{\text{mix}}(S) = \sum_{t=1}^{|S|} \phi_{\text{over}}\left(C_{b(i_t)}^{(t)}, \tau_{b(i_t)}^{(t)}\right).$$

5 Performance Evaluation

5.1 Experimental Setup

Training and Hyperparameters. We fine-tune Gemma-2-2B-it on GSM8K using LoRA. GSM8K consists of 7,500 training samples and 1,000 test samples. All runs share the same optimizer, batch size, sequence length, and LoRA settings. Models are trained for a single epoch, which we found sufficient for convergence; longer training often led to overfitting at small subset sizes. For HWD, we evaluate three hardness composition targets \mathbf{t} : a medium-dominant setting (0.10, 0.60, 0.30), a hard-heavy setting (0.05, 0.45, 0.50), and a near-uniform setting (0.34, 0.33, 0.33).

Because LoRA is sensitive to learning rate, each subset is evaluated over five learning rates, while all other hyperparameters remain fixed. The five learning rates are derived from a size-specific anchor multiplied by $\{0.7, 0.9, 1.0, 1.2, 1.5\}$. The subset selection seed is fixed (42), while fine-tuning is repeated with five random seeds. We use a mini-batch size of 8, sequence length 1024, and bf16 training.

LoRA uses rank $r=16$, $\alpha=32$, dropout 0.1, and targets `q/k/v/o_proj`. We use strict GSM8K pass@1 with a decoding limit of 512 tokens for performance comparison. The full hyperparameters are reported in Appendix B.3.

Baselines. We compare our method with three baselines: (1) pretrained Gemma-2-2B-it without fine-tuning, (2) uniform random sampling, and (3) skill-based sampling. The pretrained Gemma-2-2B-it baseline without any fine-tuning yields 62.02% pass@1 accuracy on GSM8K.

Subset Sizes. We vary the training budget $K \in \{690, 960, 1200, 1600, 2400, 3200, 4800, 7473\}$, covering low-data to full-data regimes and enabling analysis of saturation and data efficiency.

5.2 Performance Results

Tables 1 to 4 provide detailed results, and Figure 2 illustrates the overall performance trends.

We first observe that the performance of **Uniform Random Sampling** saturates early. Accuracy improves slightly up to $K = 1200$ (peaking at 63.97%) and then remains nearly flat around 64%, even when using the full dataset. This indicates substantial redundancy in GSM8K for fine-tuning.

We also observe that **Skill-Based Sampling** yields modest but inconsistent gains. Its best result occurs at $K = 2400$ with 64.41% accuracy, with performance closely matching uniform random sampling at smaller budgets and converging at larger budgets. Skill balancing alone is therefore insufficient for reliable data efficiency.

In contrast, **HWD** consistently achieves strong performance with fewer training examples by explicitly controlling difficulty and redundancy. For medium-dominant setting (0.1, 0.6, 0.3), HWD reaches its best result at $K = 690$, achieving **64.47%** accuracy using under 10% of the training data. For the hard-heavy setting (0.05, 0.45, 0.50), peak performance shifts to $K = 1600$ (**64.44%**), suggesting harder subsets require more data for stability. The near-uniform setting (0.34, 0.33, 0.33) behaves as a middle ground, reaching approximately 64.3% at moderate subset sizes.

Across all settings, HWD peaks at smaller K than the baselines, producing a clear leftward shift in the accuracy curves shown in Figure 2.

5.3 Stability

All methods exhibit low variance across random seeds, ranging from 0.1% to 0.6%. Notably,

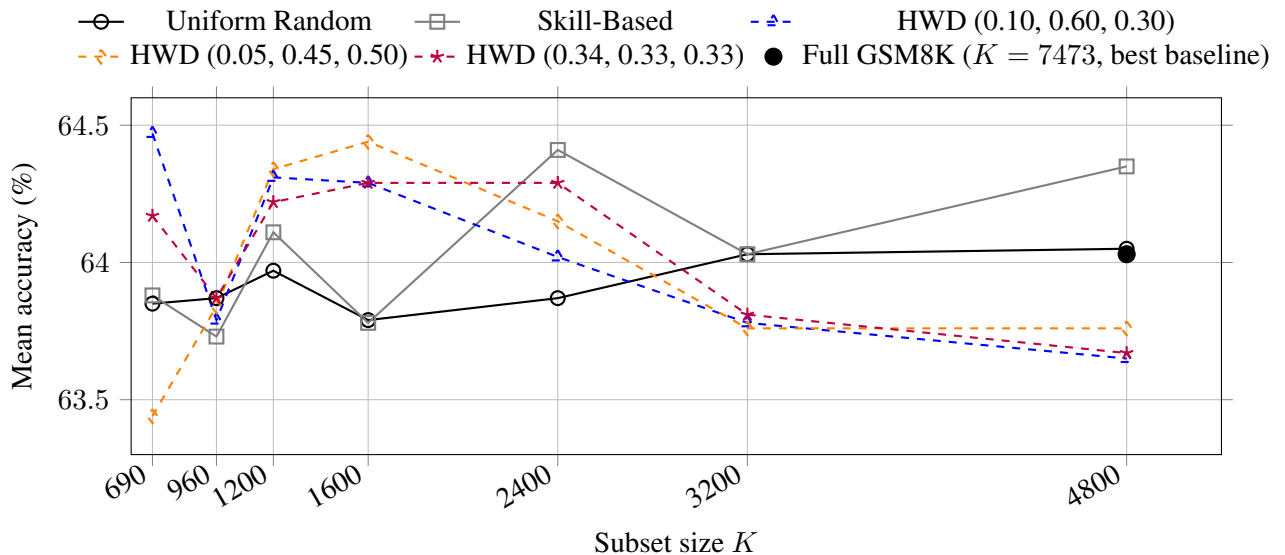


Figure 2: Accuracy trends across subset sizes K . Curves show performance up to $K = 4800$. The single marker indicates performance using the full GSM8K training set ($K = 7473$), highlighting early saturation and the data efficiency of HWD.

HWD exhibits more stable performance in low-data regimes compared with the uniform sampling method.

5.4 Summary

Overall, our findings indicate that:

- Fine-tuning on GSM8K exhibits clear performance saturation before the full training dataset is utilized, reinforcing that small, carefully selected datasets can outperform much larger ones.
- Uniform random sampling performs competitively at larger data budgets but is unreliable in low-data regimes.
- The skill-based sampling method yields occasional improvements but lacks consistent performance gains.
- Our data selection method guided by hardness-weighted diversity achieves comparable or superior performance while requiring substantially fewer training examples.

6 Analysis and Discussion

We examine why **Hardness-Weighted Diversity** outperforms random sampling by analyzing difficulty control, coverage, and stability.

Difficulty control. GSM8K is skewed toward medium-difficulty problems. HWD closely matches its target hardness distributions, preventing early over-concentration in dominant difficulty ranges—a common issue with uniform sampling.

Coverage and diversity. Ablations show that neither skill balancing nor semantic diversity alone is sufficient. Skill-only selection produces redundant reasoning patterns, while diversity-only selection underrepresents important skills. Effective performance requires controlling both.

Refinement and stability. Swap-based refinement yields small but consistent gains, suggesting that greedy selection benefits from lightweight local correction. HWD also exhibits lower variance across seeds, particularly in low-data regimes.

7 Conclusion and Future Work

This work studies data-efficient fine-tuning for mathematical reasoning and shows that performance on GSM8K saturates well before using the full dataset. In low-data regimes, subset composition plays a critical role. We introduce **Hardness-Weighted Diversity (HWD)**, an effective data selection method that jointly controls difficulty and semantic redundancy while preserving skill coverage. Across multiple subset sizes, learning rates, and random seeds, HWD matches or exceeds baseline performance using far fewer training examples.

Several directions remain open. More powerful search strategies, such as evolutionary or structured local methods, may further improve subset quality under fixed budgets. Another promising direction is adaptive hardness targeting, where the difficulty mixture evolves based on training dynamics rather than being fixed in advance.

Table 1: GSM8K fine-tuning accuracy under Uniform Random Sampling and Skill-Based Sampling. For each K , we report the best-performing learning rate averaged over five random seeds.

Strategy	K	Best LR	Mean Acc.	Std Dev.	Seeds
Uniform Random Sampling	690	6e-6	63.85%	0.40%	5/5
Uniform Random Sampling	960	3.6e-6	63.87%	0.11%	5/5
Uniform Random Sampling	1200	3e-6	63.97%	0.55%	5/5
Uniform Random Sampling	1600	2.4e-6	63.79%	0.33%	5/5
Uniform Random Sampling	2400	1.5e-6	63.87%	0.33%	5/5
Uniform Random Sampling	3200	1.1e-6	64.03%	0.31%	5/5
Uniform Random Sampling	4800	7.2e-7	64.05%	0.11%	5/5
Uniform Random Sampling	7473	4.5e-7	63.90%	0.23%	5/5
Skill-Based Sampling	690	6e-6	63.88%	0.31%	5/5
Skill-Based Sampling	960	3.6e-6	63.73%	0.38%	5/5
Skill-Based Sampling	1200	3e-6	64.11%	0.34%	5/5
Skill-Based Sampling	1600	1.8e-6	63.78%	0.28%	5/5
Skill-Based Sampling	2400	1.5e-6	64.41%	0.31%	5/5
Skill-Based Sampling	3200	1.1e-6	64.03%	0.39%	5/5
Skill-Based Sampling	4800	7.2e-7	64.35%	0.51%	5/5
Skill-Based Sampling	7473	4.5e-7	64.03%	0.27%	5/5

Table 2: Best-performing HWD configurations for the medium-dominant hardness mix (0.1, 0.6, 0.3). For each subset size K , we report the best mean accuracy across learning rates and (λ_H, λ_D) values.

HWD Mix	K	(λ_H, λ_D)	Best LR	Mean Acc.	Std
(0.1,0.6,0.3)	690	(0.8, 1.6)	6e-6	64.47%	0.46%
(0.1,0.6,0.3)	960	(0.8, 1.0)	3e-6	63.79%	0.37%
(0.1,0.6,0.3)	1200	(1.2, 1.0)	3.3e-6	64.31%	0.68%
(0.1,0.6,0.3)	1600	(0.8, 2.20)	2.2e-6	64.29%	0.30%
(0.1,0.6,0.3)	2400	(1.2, 1.6)	1.4e-6	64.02%	0.32%
(0.1,0.6,0.3)	3200	(1.2, 1.6)	1.2e-6	63.78%	0.39%
(0.1,0.6,0.3)	4800	(1.0, 2.20)	8e-7	63.65%	0.23%

Table 3: Best-performing HWD configurations for the hard-heavy hardness mix (0.05, 0.45, 0.50). Peak performance shifts to a larger subset size compared to the medium-dominant mix.

HWD Mix	K	(λ_H, λ_D)	Best LR	Mean Acc.	Std
(0.05,0.45,0.50)	690	(1.0, 1.0)	3.6e-6	63.44%	0.64%
(0.05,0.45,0.50)	960	(1.0, 1.6)	3e-6	63.84%	0.48%
(0.05,0.45,0.50)	1200	(1.2, 1.6)	3e-6	64.34%	0.39%
(0.05,0.45,0.50)	1600	(0.8, 1.0)	2.2e-6	64.44%	0.47%
(0.05,0.45,0.50)	2400	(1.2, 1.0)	1.4e-6	64.15%	0.26%
(0.05,0.45,0.50)	3200	(0.8, 1.0)	1e-6	63.76%	0.17%
(0.05,0.45,0.50)	4800	(1.2, 1.0)	8e-7	63.76%	0.34%

Table 4: Best-performing HWD configurations for the near-uniform hardness mix (0.34, 0.33, 0.33). For each subset size K , we report the best mean accuracy across learning rates and (λ_H, λ_D) values.

HWD Mix	K	(λ_H, λ_D)	Best LR	Mean Acc.	Std
(0.34,0.33,0.33)	690	(0.8, 2.20)	6e-6	64.17%	0.24%
(0.34,0.33,0.33)	960	(0.8, 1.6)	3.6e-6	63.87%	0.20%
(0.34,0.33,0.33)	1200	(1.2, 1.0)	3e-6	64.22%	0.43%
(0.34,0.33,0.33)	1600	(0.8, 2.20)	2.2e-6	64.29%	0.26%
(0.34,0.33,0.33)	2400	(1.0, 2.20)	1.4e-6	64.29%	0.35%
(0.34,0.33,0.33)	3200	(1.2, 2.20)	1e-6	63.81%	0.45%
(0.34,0.33,0.33)	4800	(0.8, 1.0)	8e-7	63.67%	0.19%

552 Limitations

553 Our proposed HWD approach relies on precom-
554 puted hardness scores and embeddings, which may
555 be noisy or imperfect estimates of true example dif-
556 ficulty. In addition, the subset selection objective
557 is optimized only approximately via greedy con-
558 struction and lightweight refinement, rather than
559 exact global optimization. Finally, difficulty is dis-
560 cretized into a fixed set of bins, which improves
561 interpretability but limits flexibility for dynamic or
562 continuously evolving settings.

563 References

564 Yoshua Bengio, Jérôme Louradour, Ronan Collobert,
565 and Jason Weston. 2009. [Curriculum learning](#). In
566 *Proceedings of the 26th International Conference on*
567 *Machine Learning*, pages 41–48. ACM.

568 Hyung Won Chung, Le Hou, Shayne Longpre, Barret
569 Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi
570 Wang, Mostafa Dehghani, Siddhartha Brahma, and
571 1 others. 2024. [Scaling instruction-finetuned lan-
572 guage models](#). *Journal of Machine Learning Re-
573 search*, 25(70):1–53.

574 Karl Cobbe, Vineet Kosaraju, Mohammad Bavar-
575 ian, Jacob Hilton, Reiichiro Nakano, Christopher
576 Hesse, and John Schulman. 2021. [Training veri-
577 fiers to solve math word problems](#). *arXiv preprint*
578 *arXiv:2110.14168*.

579 Kawin Ethayarajh, Yejin Choi, and Swabha
580 Swayamdipta. 2021. [Understanding dataset
581 difficulty with v-usable information](#). In *Proceedings*
582 *of the 38th International Conference on Machine*
583 *Learning*, volume 139, pages 3011–3021. PMLR.

584 Vitaly Feldman and Chiyuan Zhang. 2020. [What neural
585 networks memorize and why: Discovering the long
586 tail via influence estimation](#). In *Advances in Neural*
587 *Information Processing Systems*, volume 33, pages
588 2881–2891. Curran Associates, Inc.

589 Joel Hestness, Sharan Narang, Newsha Ardalani,
590 Gregory Diamos, Heewoo Jun, Hassan Kianine-
591 jad, Mostofa Patwary, Md. Mostofa Ali, Yang
592 Young, and Yanqi Zhou. 2017. [Deep learning
593 scaling is predictable, empirically](#). *arXiv preprint*
594 *arXiv:1712.00409*.

595 Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch,
596 Elena Buchatskaya, Trevor Cai, Eliza Rutherford,
597 Diego de Las Casas, Lisa Anne Hendricks, Johannes
598 Welbl, Aidan Clark, Tom Hennigan, Eric Noland,
599 Katie Millican, George Van Den Driessche, Bog-
600 dan Damoc, Aurelia Guy, Simon Osindero, Karen
601 Simonyan, Erich Elsen, and 3 others. 2022. [Train-
602 ing compute-optimal large language models](#). In *Ad-
603 vances in Neural Information Processing Systems*,
604 volume 35, pages 30016–30030. Curran Associates,
605 Inc.

Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B.
Brown, Benjamin Chess, Rewon Child, Scott Gray,
Alec Radford, Jeffrey Wu, and Dario Amodei. 2020.
[Scaling laws for neural language models](#). *arXiv*
preprint arXiv:2001.08361.

Pang Wei Koh and Percy Liang. 2017. [Understand-
ing black-box predictions via influence functions](#).
In *International Conference on Machine Learning*
(ICML).

M. Pawan Kumar, Benjamin Packer, and Daphne Koller.
2010. [Self-paced learning for latent variable mod-
els](#). In *Advances in Neural Information Processing*
Systems, volume 23, pages 1189–1197. Curran Asso-
ciates, Inc.

Qian Liu, Wenbo Li, and Yue Zhang. 2023. [Curricu-
lum learning for mathematical reasoning with large
language models](#). In *Findings of the Association for*
Computational Linguistics: ACL 2023, pages 918–
931.

Baharan Mirzasoleiman, Jeff Bilmes, and Jure Leskovec.
2020. [Coresets for data-efficient training of machine
learning models](#). In *Proceedings of the 37th Inter-
national Conference on Machine Learning*, volume
119 of *Proceedings of Machine Learning Research*,
pages 6950–6960. PMLR.

Mansheej Paul, Surya Ganguli, and Gintare Karolina
Dziugaite. 2021. [Deep Learning on a Data Diet:
Finding Important Examples Early in Training](#). In
Advances in Neural Information Processing Systems
34 (*NeurIPS*), pages 20596–20607.

Garima Pruthi, Frederick Liu, Satyen Kale, and Mukund
Sundararajan. 2020. [Estimating training data influ-
ence by tracing gradient descent](#). In *Advances in*
Neural Information Processing Systems, volume 33,
pages 19920–19930. Curran Associates, Inc.

Mariya Toneva, Alessandro Sordoni, Remi Tachet des
Combes, Adam Trischler, Yoshua Bengio, and Geof-
frey J. Gordon. 2019. [An empirical study of example
forgetting during deep neural network learning](#). In
*International Conference on Learning Representa-
tions*.

Kai Wei, Rishabh Iyer, and Jeff Bilmes. 2015. [Submod-
ularity in data subset selection and active learning](#).
In *Proceedings of the 32nd International Conference*
on Machine Learning, volume 37 of *Proceedings of*
Machine Learning Research, pages 1954–1963, Lille,
France. PMLR.

Mengzhou Xia, Sadhika Akula, Arindam Mitra, Shuo
Sun, Reza Bashizade, Wenhui Chen, Yizhong Wang,
Zhaojiang Chen, Yashar Mehdad, and Luke Zettle-
moyer. 2024. [Data valuation for efficient fine-tuning
of large language models](#). *Transactions of the Asso-
ciation for Computational Linguistics*, 12:895–911.

Chunting Zhou, Pengfei Liu, Puxin Xu, Srini Iyer,
Jiao Sun, Yuning Wang, Zheyuan Wang, Chu-Cheng
Chen, Benjamin Deshpande, Omer Levy, Luke Ma,

662	Marjan Ghazvininejad, Luke Zettlemoyer, and Badih	710
663	Gupta. 2023a. LIMA: Less Is More for Alignment.	711
664	In <i>Advances in Neural Information Processing Sys-</i>	712
665	<i>tems 36 (NeurIPS)</i> .	713
666	Xuhui Zhou, Zhihan Kan, Guanbin Li, Jieyu Zhang,	714
667	and Yiren Zhao. 2023b. Difficulty-aware curricu-	715
668	lum learning for mathematical reasoning. In <i>Find-</i>	716
669	<i>ings of the Association for Computational Linguistics:</i>	717
670	<i>EMNLP 2023</i> , pages 13508–13521.	
671	A Selection Signals	
672	This appendix documents the auxiliary signals used	
673	exclusively for subset construction in Hardness-	
674	Weighted Diversity (HWD) . These signals are	
675	used <i>only</i> during data selection and are not accessed	
676	during fine-tuning or evaluation.	
677	A.1 Hardness Score Computation	
678	We define a per-example hardness score from re-	
679	peated evaluation of each training problem un-	
680	der a strict numeric-answer rule. For each	
681	GSM8K training item i , we query a fixed refer-	
682	ence model (teacher) for M independent attempts	
683	using stochastic decoding (enabled when $M > 1$),	
684	and mark an attempt as correct only if the model’s	
685	output ends with a line matching the exact for-	
686	mat <code>#### <number></code> and the extracted number	
687	equals the ground-truth numeric answer.	
688	Let $\text{win}_i \in [0, 1]$ denote the winning rate (frac-	
689	tion of correct attempts) across the M trials. We	
690	convert this success rate into hardness	
691	$H_i = 1 - \text{win}_i,$	
692	so that larger values correspond to harder prob-	
693	lems. The resulting hardness file stores, for each	
694	example id, the winning rate <code>acc</code> and/or hardness	
695	<code>hardness</code> , together with the number of trials	
696	<code>n=M</code> .	
697	All values are clamped to $[0, 1]$. If a percentage	
698	scale is detected, values are rescaled by dividing	
699	by 100.	
700	A.2 Semantic Embedding Computation	
701	We precompute semantic embeddings for all	
702	GSM8K training problems using a frozen	
703	sentence encoder. Specifically, we use	
704	BAAI/bge-small-en-v1.5 to encode	
705	each problem statement independently. For each	
706	input, we extract the hidden state of the first token	
707	from the final encoder layer (CLS-style pooling)	
708	and apply ℓ_2 normalization to obtain a unit-length	
709	vector.	
	Embeddings are computed in batches and stored	
	as a dense matrix $\mathbf{E} \in \mathbb{R}^{N \times d}$ in NumPy format.	
	A separate JSON file records the mapping from	
	example ids to embedding rows. During subset	
	selection, these embeddings are used exclusively	
	to compute cosine-based novelty scores. The em-	
	bedding model remains frozen throughout and is	
	never updated or fine-tuned.	
	A.3 Skill Taxonomy and Labeling	
	We annotate each GSM8K problem with one or	
	more skills drawn from a fixed taxonomy of 17	
	mathematical reasoning categories. Each skill	
	captures a distinct type of reasoning or computa-	
	tion pattern, and problems may belong to multiple	
	skills.	
	Skill definitions. The 17 skills are defined as	
	follows (IDs shown in parentheses):	
	1. Basic Arithmetic (1) : Single-operation or short-	
	chain arithmetic.	
	2. Multi-step Arithmetic (2) : Sequential depen-	
	dent computations.	
	3. Fractions & Ratios (3) : Fractional and ratio	
	reasoning.	
	4. Algebraic Reasoning (4) : Equation-based rea-	
	soning with variables.	
	5. Unit Conversion (5) : Measurement or currency	
	conversion.	
	6. Counting / Combinatorics (6) : Counting and	
	simple combinatorics.	
	7. Proportional Reasoning (7) : Rates, scaling,	
	and proportionality.	
	8. Geometry / Spatial Reasoning (8) : Shapes,	
	areas, and volumes.	
	9. Logical Deduction / Constraints (9) :	
	Constraint-based deduction.	
	10. Temporal Reasoning (10) : Time intervals and	
	schedules.	
	11. Money / Cost Calculations (11) : Prices, dis-	
	counts, and costs.	
	12. Table Lookup / Structured Data (12) : Reason-	
	ing over tables.	

- 751 13. **Common Sense / World Knowledge (13):** Ev- 793
 752 eryday factual knowledge. 794
 753 14. **Multi-hop Chain-of-Thought Reasoning (14):** 795
 754 Multi-step inference. 796
 755 15. **Comparison / Difference Finding (15):** Rela- 797
 756 tive quantity reasoning. 798
 757 16. **Estimation / Rounding (16):** Approximate nu- 799
 758 meric reasoning. 800
 759 17. **Nested Reasoning (17):** Hierarchically struc-
 760 tured reasoning.

761 **Annotation source.** Skill labels are obtained via
 762 automatic annotation using a large language model
 763 (GPT), which assigns one or more skills from the
 764 predefined taxonomy to each problem based on
 765 its problem statement. These annotations are used
 766 solely as a weak supervisory signal for regulariza-
 767 tion during subset selection and are not treated as
 768 ground-truth labels.

769 **Empirical distribution.** The resulting skill dis-
 770 tribution on GSM8K is highly imbalanced: multi-
 771 step arithmetic dominates as the primary skill for
 772 the majority of problems, while other skills occur
 773 infrequently and often only as secondary labels.
 774 Consequently, HWD applies a soft, asymmetric
 775 skill regularization rather than enforcing uniform
 776 per-skill coverage.

777 B Experimental Protocol

778 B.1 Subset Construction Procedure

779 **Candidate universe.** We retain only items that
 780 have both (i) a valid hardness value H_i and (ii) a
 781 valid embedding row mapping. All other items are
 782 excluded from selection.

783 **Top- M restriction.** To reduce computation, se-
 784 lection is restricted to a Top- M pool formed by
 785 sorting eligible items in descending hardness and
 786 retaining

$$787 M = \text{clamp}(\text{mult} \cdot K, M_{\min}, M_{\max}),$$

788 with mult , M_{\min} , and M_{\max} controlled by
 789 `--v1-topM-mult`, `--v1-topM-min`, and
 790 `--v1-topM-max`. This restriction limits the
 791 search space only; within Top- M , greedy selec-
 792 tion uses the full scoring function.

Soft hardness-mix control. We encourage a tar-
 get hardness composition $\mathbf{t} = [t_1, t_2, t_3]$ (passed
 via `--v1-hard-targets`). At step $m+1$, the
 running target for bin b is $\tau_b^{(m+1)} = (m+1)t_b$.
 With slack δ (`--v1-hard-slack`), we apply an
 overshoot-only quadratic penalty if adding a candi-
 date would exceed $(1+\delta)\tau_b^{(m+1)}$:

$$\text{pen}_{\text{Mix}}(i) = \left(\frac{[C_{b(i)} + 1 - (1 + \delta)\tau_{b(i)}^{(m+1)}]_+}{\max(1, (1 + \delta)\tau_{b(i)}^{(m+1)})} \right)^2.$$

Greedy score. At each iteration, we select the
 candidate maximizing:

$$\text{Score}(i | S) = \lambda_H H_i + \lambda_D D_i(S) - \lambda_S \text{pen}_S(i) - \lambda_{\text{Mix}} \text{pen}_{\text{Mix}}(i).$$

where $(\lambda_H, \lambda_D, \lambda_S, \lambda_{\text{Mix}})$ correspond
 to `--v1-lambda-h`, `--v1-lambda-d`,
`--v1-lambda-s`, and `--v1-hard-weight`.
 The greedy procedure initializes with the single
 hardest item in Top- M .

Swap polishing. After greedy selection, we per-
 form up to N_{swaps} random one-for-one swaps
 (`--v1-swaps`). A proposed swap is accepted
 if it improves the global objective computed on the
 full set.

814 B.2 Fine-Tuning and Evaluation

Training. We fine-tune **Gemma-2-2B-it** with
 LoRA under a fixed configuration (batch size, se-
 quence length, epochs, and LoRA hyperparam-
 eters), held constant across all subset strategies.
 Each run is performed for a single epoch.

Learning-rate sweep and seeds. For each
 (method, K) pair, we sweep five learning rates and
 evaluate five random seeds. We select the learning
 rate that maximizes the *mean* accuracy across seeds
 and report the mean and standard deviation under
 that learning rate.

Metric. We report **strict pass@1 accuracy** on
 GSM8K, where a prediction is counted correct only
 if the final numeric answer exactly matches the
 ground truth under a strict formatting rule.

830 B.3 Reproducibility Details

Fine-tuning. We fine-tune **Gemma-2-2B-it** with
 LoRA for one epoch under a fixed setup across
 all runs. We use **batch size 1**, mixed precision
 (`bf16`), and evaluate **five random seeds** (41–45)
 for each configuration.

836 **Budgets and learning rates.**

837 We evaluate subset sizes $K \in$
838 $\{690, 960, 1200, 1600, 2400, 3200, 4800\}$ (and
839 $K = 7473$ for the full set). For each (method, K)
840 pair we sweep **five learning rates** and report the
841 learning rate that maximizes mean accuracy across
842 seeds. The per- K learning-rate grids are:

- 843 • $K=690$: $\{3.5, 4.5, 5.0, 6.0, 7.5\} \times 10^{-6}$
- 844 • $K=960, 1200$: $\{2.1, 2.7, 3.0, 3.6, 4.5\} \times 10^{-6}$
- 845 • $K=1600$: $\{1.4, 1.6, 1.8, 2.0, 2.2\} \times 10^{-6}$
- 846 • $K=2400$: $\{1.0, 1.2, 1.4, 1.6, 1.8\} \times 10^{-6}$
- 847 • $K=3200$: $\{0.8, 1.0, 1.2, 1.4, 1.6\} \times 10^{-6}$
- 848 • $K=4800$: $\{0.4, 0.5, 0.6, 0.8, 1.0\} \times 10^{-6}$

849 **HWD selection hyperparameters.** Unless
850 stated otherwise, HWD uses hardness bins
851 $(0.0, 0.5, 0.8, 1.0)$ and the target mix \mathbf{t} speci-
852 fied in each experiment. We set Top- M via
853 $M = \text{clamp}(4K, 2000, 10000)$, use slack
854 $\delta = 0.01$, and apply swap polishing with
855 $N_{\text{swaps}} = 300$. We use $\lambda_S = 0.10$ and skill
856 tolerance $\alpha = 1.5$, and sweep (λ_H, λ_D) as
857 reported in the main tables.

858 **C Additional Analyses**

859 **C.1 Stability Across Random Seeds**

860 **Confidence interval computation.** For all re-
861 ported confidence intervals, we compute a two-
862 sided 95% confidence interval over five independ-
863 ent runs as

$$864 \quad \bar{x} \pm t_{0.975,4} \cdot \frac{s}{\sqrt{n}},$$

865 where \bar{x} is the mean accuracy, s is the sample stan-
866 dard deviation, $n = 5$, and $t_{0.975,4} = 2.776$ is
867 the Student- t critical value with four degrees of
868 freedom. Narrower intervals indicate more stable
869 performance across random seeds.

Table 5: 95% confidence intervals (5 seeds) for baseline sampling strategies.

Uniform Random Sampling			
K	Mean	Std	95% CI
690	63.85	0.40	[63.50, 64.20]
960	63.87	0.11	[63.77, 63.97]
1200	63.97	0.55	[63.49, 64.45]
1600	63.79	0.33	[63.50, 64.08]
2400	63.87	0.33	[63.58, 64.16]
3200	64.03	0.31	[63.76, 64.30]
4800	64.05	0.11	[63.95, 64.15]
7473	63.90	0.23	[63.70, 64.10]

Skill-Based Sampling			
K	Mean	Std	95% CI
690	63.88	0.31	[63.61, 64.15]
960	63.73	0.38	[63.40, 64.06]
1200	64.11	0.34	[63.81, 64.41]
1600	63.78	0.28	[63.53, 64.03]
2400	64.41	0.31	[64.14, 64.68]
3200	64.03	0.39	[63.69, 64.37]
4800	64.35	0.51	[63.90, 64.80]
7473	64.03	0.27	[63.79, 64.27]

Table 6: 95% confidence intervals (5 seeds) for selected HWD configurations under three hardness mixes.

Mix (0.1, 0.6, 0.3)			
K	Mean	Std	95% CI
690	64.47	0.46	[64.06, 64.88]
960	63.79	0.37	[63.47, 64.12]
1200	64.31	0.68	[63.71, 64.91]
1600	64.29	0.30	[64.03, 64.55]
2400	64.02	0.32	[63.74, 64.30]
3200	63.78	0.39	[63.44, 64.12]
4800	63.65	0.23	[63.45, 63.85]

Mix (0.05, 0.45, 0.50)			
K	Mean	Std	95% CI
690	63.44	0.64	[62.65, 64.23]
960	63.84	0.48	[63.24, 64.43]
1200	64.34	0.39	[63.85, 64.82]
1600	64.44	0.47	[63.86, 65.03]
2400	64.15	0.26	[63.83, 64.48]
3200	63.76	0.17	[63.55, 63.98]
4800	63.76	0.34	[63.34, 64.19]

Mix (0.34, 0.33, 0.33)			
K	Mean	Std	95% CI
690	64.17	0.24	[63.87, 64.47]
960	63.87	0.20	[63.62, 64.12]
1200	64.22	0.43	[63.68, 64.76]
1600	64.29	0.26	[63.97, 64.61]
2400	64.29	0.35	[63.86, 64.72]
3200	63.81	0.45	[63.25, 64.37]
4800	63.67	0.19	[63.44, 63.90]