KNOWLEDGE-CENTRIC DATA SELECTION FOR EFFECTIVE DOMAIN ADAPTATION OF LARGE LANGUAGE MODELS

Anonymous authorsPaper under double-blind review

000

001

002

004

006

008 009 010

011 012

013

014

015

016

017

018

019

021

023

025

026

027

028

029

031

032

034

036

037

040

041

042

043

044

046

047

048

051

052

ABSTRACT

Domain adaptation of language models is critical for specialized applications in fields, but its success hinges on high-quality data selection rather than sheer volume. Current methods, such as heuristic filters, perplexity pruning, and embedding-based clustering, often fail to address domain-specific redundancy and interferences, leading to inefficient training and models that overfit to frequent linguistic forms rather than core knowledge. This inefficiency leads to limited generalization, as well as high curation and computational costs, particularly in high-stakes domains where errors carry significant consequences. We propose a knowledge-centric approach that redefines data quality around discrete knowledge procedures and theorems. Our framework introduces Knowledge Coverage Entropy (KCE), a metric quantifying knowledge diversity, and Entropy-Driven Selection (EDS), which optimizes data selection for compact, high-quality datasets. Experiments in supervised fine-tuning (SFT) and retrieval-augmented generation (RAG) demonstrate EDS's efficacy. In SFT on the MATH-500 benchmark, at matched data budgets, our method consistently yields the best post-training accuracy among data selection methods. In RAG on medical datasets, our method delivers the best retrieval quality with mean reciprocal rank (MRR) improvements of approximately 11% to 42% and substantial coverage gains while using significantly fewer samples. Enhanced performance in both SFT and RAG demonstrates that KCE offers a principled metric for data quality, and that EDS facilitates efficient in domain-specific tasks.

1 Introduction

Domain adaptation tailors general-purpose language models for specialized tasks, embedding domain-specific knowledge and reasoning (Howard & Ruder, 2018; Longpre et al., 2023; Seto et al., 2025; Parmar et al., 2024). Unlike broad fluency training, adaptation via supervised fine-tuning (SFT) or retrieval-augmented generation (RAG) prioritizes precise, context-relevant concepts (Shum et al., 2024; Muennighoff et al., 2025a). Effective adaptation hinges on data quality, not volume, requiring corpora that capture essential knowledge for model internalization (Pang et al., 2025; Xia et al., 2024; Liu et al., 2023). Uncurated datasets yield diminishing returns or degraded performance under minimum description length principles (Li & Vitányi, 2008), emphasizing the need for high-quality selection to optimize learning, efficiency, and robustness (Longpre et al., 2023; Seto et al., 2025).

Current methods, including heuristic filters (e.g., text length, readability) (Xia et al., 2024; Liu et al., 2023), perplexity-based pruning (Pang et al., 2025; Ankner et al., 2024), embedding-based clustering (Xie et al., 2024), and entropy-driven approaches (Song et al., 2012; Lairez, 2022), manage large corpora but struggle with domain-specific redundancy, such as rephrased definitions or overlapping evidence (Lee et al., 2022; Hei et al., 2024). These methods, relying on tokens or embeddings, fail to distinguish novel knowledge from stylistic variations, leading to reliance on unscalable manual curation (Liu et al., 2024; Wang et al., 2024). Consequently, SFT models memorize frequent patterns, faltering on edge cases, while RAG retrieves irrelevant or redundant data, increasing compute demands and reducing generalization (Amiraz et al., 2025; Hager et al., 2024; Fayyaz et al., 2025). In high-stakes domains like medicine or law, these inefficiencies heighten risks.

 Empirical studies highlight these challenges. Noisy SFT datasets with incorrect or mismatched pairs degrade accuracy and introduce biases (Liu et al., 2024; Wang et al., 2024). Unfiltered RAG datasets reduce retrieval precision (Amiraz et al., 2025). Fine-tuning models like BERT on uncurated data causes overfitting to superficial patterns, impairing out-of-distribution performance. Importance sampling amplifies errors with scarce in-domain data, and unsupervised curricula fail to address conceptual overlap without validation (Ankner et al., 2024; Pang et al., 2025). These gaps underscore the need for balanced data quality: correctness, relevance, coverage, minimal redundancy, and informativeness for efficient SFT and RAG.

A knowledge-centric approach addresses these issues by focusing on discrete, auditable knowledge units (e.g., definitions, theorems, guidelines), evaluated for quality dimensions (Fu et al., 2025; Deb et al., 2025; Agarwal et al., 2024; Muennighoff et al., 2025a). Correctness ensures factual accuracy; relevance aligns with domain queries; coverage includes diverse concepts; redundancy avoids overrepresented ideas; and informativeness targets challenging samples. Using information-theoretic principles like mutual information (Song et al., 2012; Lairez, 2022; Majenz, 2018), this framework prioritizes novel knowledge via a binary coverage matrix, enhancing interpretability and scalability for SFT and RAG across applications like mathematical reasoning and clinical support.

We propose a knowledge-centric framework for data selection in domain adaptation of large language models that combines Knowledge Coverage Entropy (KCE) and Entropy-Driven Filtering (EDF). KCE quantifies diversity and balance over discrete knowledge units, and EDF prioritizes novel, high-information samples to reduce redundancy. By leveraging entropy to emphasize informative coverage, the framework strengthens supervised fine-tuning learning signals and improves retrieval-augmented generation retrieval precision. On the MATH-500 benchmark, at matched data budgets, KCE-selected data yields the best post-training accuracy among data selection methods and reaches 456/500 with substantially fewer samples. In medical retrieval-augmented generation, the framework delivers the best retrieval quality with mean reciprocal rank improvements of approximately 11% to 42% alongside large coverage gains under significant data reduction. These results establish KCE and EDF as principled tools for efficient and high-performance domain adaptation.

2 METHODOLOGY

The Entropy-Driven Selection (EDS) methodology selects a diverse and informative subset of data samples by maximizing Knowledge Coverage Entropy (KCE) within a binary information-knowledge matrix. This approach constructs a matrix representing knowledge points across samples, computes entropy-based scores to quantify diversity, and employs a set-aware lazy-greedy algorithm to optimize subset selection under cardinality constraints.

2.1 BINARY INFORMATION-KNOWLEDGE MATRIX

We construct a knowledge set \mathcal{K} of domain-relevant concepts and map each data sample to a binary vector over \mathcal{K} , forming a matrix $\mathbf{B} \in \{0,1\}^{n \times m}$, where n is the number of samples, $m = |\mathcal{K}|$ is the number of knowledge points, and $\mathbf{B}_{i,j} = 1$ if sample i covers knowledge point j, and 0 otherwise. The matrix \mathbf{B} is built using Qwen-max-0125 (Team, 2025) with task-specific prompts to extract and tag concepts, as detailed in Appendix C. This matrix underpins the computation of Knowledge Coverage Entropy (KCE).

2.2 COVERAGE PROBABILITY DEFINITIONS

For the matrix $\mathbf{B} \in \{0,1\}^{n \times m}$, we define the smoothed coverage probability for sample a as $P_a = \frac{\sum_{j=1}^m \mathbf{B}_{a,j} + \alpha}{m + \alpha m}$, where $\alpha = 10^{-6}$ ensures numerical stability. The joint probability distribution is computed as $P_{i,j} = \frac{\mathbf{B}_{i,j} + \alpha/(nm)}{\sum_{i=1}^n \sum_{j=1}^m (\mathbf{B}_{i,j} + \alpha/(nm))}$. These probabilities support entropy calculations, with further details in Appendix B.2.

2.3 KNOWLEDGE COVERAGE ENTROPY (KCE)

For a subset $S \subseteq \{1, \ldots, n\}$ of size |S| = h, KCE is defined as $H(S) = -\sum_{j=1}^m p_j \log_2 p_j$, where $p_j = \frac{1}{h} \sum_{a \in S} \mathbf{B}_{a,j}$ is the average coverage of knowledge point j in S. The entropy is normalized by subset size, $H_n(S) = H(S)/h$, to compare diversity across subsets. Upper bounds and properties are discussed in Appendix B.2.

2.4 Entropy-Driven Selection Algorithm (EDS)

The EDS algorithm selects a subset S of size |S| = s that maximizes KCE, addressing a combinatorial optimization problem. Below, we describe the optimization goal and the set-aware lazy-greedy algorithm used to achieve it efficiently, with theoretical justifications provided in Appendices B.6 and B.5.

2.4.1 Optimization Objective

The goal is to identify a subset S that maximizes KCE:

$$S^* = \underset{S \subseteq \{1, \dots, n\}, |S| = s}{\operatorname{arg\,max}} H(S).$$

This problem is computationally intractable due to its combinatorial nature, necessitating approximate strategies. We employ a submodular optimization approach, leveraging the diminishing returns property of KCE (see Appendix B.4).

2.4.2 SET-AWARE LAZY-GREEDY SELECTION

To maximize KCE efficiently, we define a concave-over-coverage objective:

$$F(S) = \sum_{j=1}^{m} w_j f(c_j(S)), \quad c_j(S) = \sum_{a \in S} \mathbf{B}_{a,j},$$

where $w_j \in \mathbb{R}_+^m$ are weights reflecting the importance of knowledge point j (estimated from the dataset distribution), and f is a concave, nondecreasing function. This objective is nonnegative, monotone, and submodular, ensuring that a greedy algorithm achieves a (1-1/e) approximation to the optimal solution, as detailed in Appendix B.4. The lazy-greedy algorithm, summarized in Appendix B.3. The choice of f balances fidelity to KCE (using the entropy-derived h) and computational efficiency (using $\log(1+x)$). Each marginal gain evaluation has complexity $O(\operatorname{nnz}(\mathbf{B}_{a,\cdot}))$, and the lazy-greedy approach scales efficiently with sparse matrices. An optional early stopping criterion, based on a revenue boundary, is discussed in Appendix B.5.

2.4.3 WEIGHTED ENTROPY SCORING

To incorporate domain-specific priorities, we encode concept priorities with a weight vector $\mathbf{k} \in \mathbb{R}^m$ (e.g., from concept frequencies). For sample a, define $P_a = \frac{1}{m} \sum_{j=1}^m \mathbf{B}_{a,j}$ and $H(a) = -P_a \log_2 P_a$. The weighted score is

$$Score(a) = H(a) \left(1 + \gamma \sum_{i=1}^{m} k_i \mathbf{B}_{a,i} \right),$$

where γ trades off diversity and importance. This heuristic steers greedy selection toward diverse samples emphasizing high-priority concepts. Estimation of \mathbf{k} and single-pass variants are in Appendix B.1.

3 EXPERIMENTS AND EVALUATIONS

We evaluate our entropy-driven data selection framework in two paradigms: supervised fine-tuning (SFT) for mathematical chain-of-thought (CoT) and retrieval-augmented generation (RAG). Baselines include QuRating (Wettig et al., 2024), SuperFiltering (Li et al., 2024), Structure Entropy (Xie et al., 2024), random sampling, and the human-curated S1 subset (Muennighoff et al., 2025a). For RAG, we construct proprietary diabetes and general medical corpora and compare matched-size selections across methods. Ablations vary retrieval depth, corpus size, and top-k.

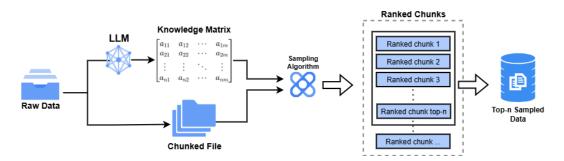


Figure 1: Overview of the Entropy-Driven Data Selection (EDS) workflow. EDS integrates normalized entropy, marginal gain, and knowledge-aware weighting to balance information efficiency with domain coverage.

3.1 SUPERVISED FINE-TUNING EVALUATION BENCHMARK

We perform CoT SFT on the S1 data-ablation-full59K pool (Muennighoff et al., 2025b), using the human-curated 1k subset (simplescaling/s1K-tokenized) (Muennighoff et al., 2025b;a) as a high-quality reference. Our method selects matched-size subsets (1k, 5k, 10k) under identical preprocessing and prompting, with baselines producing size- and token-matched counterparts. We fine-tune Qwen-32B-Instruct (Team, 2024) via standard next-token cross-entropy with consistent schedules across conditions. Selection is guided by knowledge coverage entropy (KCE), computed over a knowledge–sample matrix to balance per-sample uncertainty (row entropy) and global coverage, reducing redundancy and promoting diverse reasoning structures. The only difference across conditions is the upstream selection criterion.

We evaluate on the MATH-500 exam set using VLLM inference, reporting exact-answer accuracy. First, we compare KCE with non-negative knowledge-point weights to prioritize rare but critical units—against random sampling and the human-curated S1 subset. Next, we test alternative weighting schemes across baseline selectors by generating matched-size subsets and retraining under identical SFT protocols. We report overall accuracy, sample efficiency, and training stability.

3.2 CONSTRUCTION OF PROPRIETARY RAG CORPORA

We programmatically compile domain-relevant sources (diabetes: textbooks and clinical guidelines; general medical) (Holt & Flyvbjerg, 2024; Royal Government of Bhutan, Ministry of Health, Department of Medical Service, 2007; fun) and use LLMs to: (i) segment texts into atomic chunks, (ii) normalize to a controlled vocabulary of knowledge IDs, and (iii) finalize retrieval-ready passages with titles and structured metadata (knowledge IDs, source, language, timestamps). For each domain, we generate 1000 LLM-authored questions with automatic validation and light manual spot checks. We embed passages with BAAI/bge-large-zh and BAAI/bge-large-en (Chen et al., 2023; Xiao et al., 2023) and retrieve by cosine similarity (Salton et al., 1975) (top-k). Matched-size corpus variants are produced via our selection, QuRating, SuperFiltering, Structure Entropy, and the unselected full corpus.

To assess the selected corpora, we compute knowledge-point coverage rate $\mathrm{Hit}@k$ (the proportion of ground-truth knowledge points covered within the $\mathrm{top}\text{-}k$ retrieved passages) and conventional MRR, and analyze the accuracy–efficiency trade-off as a function of corpus size. We first evaluate at $\mathrm{top}\text{-}10$ retrieval, where each selection method operates at its theoretical data-efficiency point. We then vary (i) retrieval depth with $k \in \{5, 10, 20, 50\}$ and (ii) corpus size, always comparing under matched-size settings.

3.3 RAG EXPERIMENTS AND EVALUATION BENCHMARK

Let Q be the query set with |Q| = N. For each query $q \in Q$, let K(q) denote the required knowledge points and $R_k(q)$ the set of knowledge points covered by the top-k retrieved entries (from annotated knowledge IDs).

The per-query knowledge-point hit rate at depth k is:

$$HitRate_k(q) = \frac{|K(q) \cap R_k(q)|}{|K(q)|}.$$
 (1)

The average knowledge-point hit rate is:

$$AverageHitRate_k = \frac{1}{N} \sum_{q \in Q} HitRate_k(q). \tag{2}$$

Define $r_k(q)$ as the smallest $r \in \{1, 2, ..., k\}$ such that the union of knowledge points covered by the top-r retrieved entries contains all elements of K(q). If no such r exists within the top-k entries, set $r_k(q) = 0$ by convention.

The per-query reciprocal rank is

$$RR(q) = \begin{cases} \frac{1}{r_k(q)}, & \text{if } r_k(q) \ge 1, \\ 0, & \text{if } r_k(q) = 0. \end{cases}$$
 (3)

The average multi-point MRR (distinct from conventional MRR (Voorhees & Tice, 2000), as it requires setwise completion of K(q) is:

$$AvgMRR_k = \frac{1}{N} \sum_{q \in Q} RR(q).$$
 (4)

We compute knowledge-point coverage at depth k (Hit@k) and the multi-point MRR, and also report conventional MRR for comparison. All configurations use identical embedding models, cosine-similarity retrieval, indexing, and query/annotation sets; selection methods differ only in the upstream criterion (KCE vs. baselines).

We then conduct two classes of experiments. (i) Fixed-size corpora: for each domain, we construct a matched-size evaluation corpus (Diabetes: 3K; Medical: 8K) for each selector and vary retrieval depth with $k \in \{5, 10, 20, 50\}$. (ii) Variable-size corpora: for each selector, we subsample 1%, 5%, 10%, 20%, and 50% of the full corpus and evaluate at multiple k. To operationalize the "revenue boundary," we sweep corpus-size-performance curves and select the smallest subset within 1% relative performance of the maximum Hit@10, yielding the data-efficiency point.

4 RESULTS

4.1 Entropy-Driven SFT Performance on MATH-500

To evaluate our SFT data selection algorithm, we conducted experiments on the MATH-500 benchmark. Specifically, we compared 28 randomly sampled subsets with 28 entropy-selected subsets across different dataset sizes. All models were trained to convergence using an early stopping criterion (loss ≤ 0.05 , patience = 5), and inference was performed with the VLLM framework (Kwon et al., 2023), with the temperature fixed at 0 to eliminate stochastic variation. The model performance curves are shown in Fig. 2, and the complete performance results are summarized in Table 6. Across all dataset scales, entropy-based selection consistently outperforms random sampling, highlighting its ability to identify high-quality training data. Even relatively small entropy-selected subsets achieve performance comparable to much larger randomly sampled sets, demonstrating strong data efficiency. Notably, the entropy-selected subset reaches 450/500 at size 1000, closely matching the manually curated S1 dataset (452/500), and even exceeds it at size 500 (456/500). This consistent advantage across scales validates knowledge-point entropy as a principled and effective criterion for data selection.

The training loss trajectories are demonstrated in Figure 4 for models trained on 40K and 50K samples, selected via entropy-based selection or random sampling. Entropy-selected subsets consistently converge faster and more stably, with the 40K subset exhibiting lower loss and a steeper early

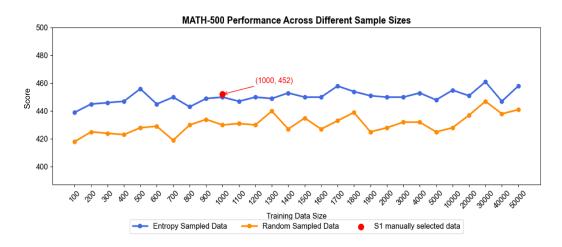


Figure 2: MATH-500 Performance with different trainind data size, The horizontal axis shows the dataset size, and the vertical axis shows the model's test scores. The red point at (1000, 452) indicates that the S1 team selected 1000 samples, achieving a score of 452 on the MATH-500 benchmark.

decline, indicating stronger gradient signals from high-quality data. Interestingly, the 40K entropy-selected subset achieves slightly lower final loss than the 50K subset, suggesting that smaller, carefully curated datasets can provide more informative gradients than larger, unfiltered ones. The Entropy-driven selection delivers strong performance with fewer samples, demonstrating that principled data selection is a more efficient and practical strategy for supervised fine-tuning (SFT) than indiscriminate dataset expansion.

4.1.1 SFT EVALUATION OF BASELINE DATA SELECTION ALGORITHMS

We compared our method against several baseline algorithms under LoRA fine-tuning across varying training data sizes, as shown in Figure 3. Overall, our KCE-based method consistently achieves higher exact answer accuracy than the baseline algorithms (Structure Entropy, QuRating, and Super-Filtering) at most dataset sizes, demonstrating its effectiveness in selecting high-quality, informative samples. Notably, KCE with knowledge-point weighting outperforms the unweighted variant in most cases (e.g., 455 vs. 444 at size 1000, 450 vs. 447 at size 2000), indicating that incorporating knowledge-point weights helps prioritize rare but critical knowledge units, further enhancing model performance. These results validate both the superiority of our entropy-driven selection method and the utility of weighted knowledge coverage for efficient and effective SFT.

4.2 RETRIEVAL EFFICIENCY ON MEDICAL KNOWLEDGE DATASETS

We evaluated our entropy-driven data selection framework, focusing on the proposed Knowledge-Centric Entropy (KCE) method, on two medical datasets: Diabetes and General Medical, comparing it with Qurating, Structural Entropy, and Superfiltering. The revenue boundary, illustrated in Figure 5, indicates where adding more samples yields diminishing returns, allowing redundant data to be discarded while preserving the most valuable knowledge. As summarized in Table 1, KCE consistently improves retrieval performance. On the General Medical dataset, which is high-dimensional with many sparse attributes and a less pronounced revenue boundary (as shown in Figure 5), KCE effectively prioritizes the most informative samples, leading to a notable improvement in MRR. Although the average coverage rate shows a slight decrease, it remains high, demonstrating that KCE enhances retrieval quality with minimal impact on overall coverage.

For the **Diabetes** dataset (251 Attributes), KCE achieves the highest coverage rate and MRR among all selection methods, increasing coverage from 75.5% to 79.3% and MRR from 0.431 to 0.480, while reducing the dataset size from 12K to 3K. In the **General Medical** dataset (1,122 Attributes), KCE maintains coverage, slightly decreasing from 73.9% to 72.9%, and further improves MRR from 0.451 to 0.468, despite a significant reduction in data size from 20K to 8K.

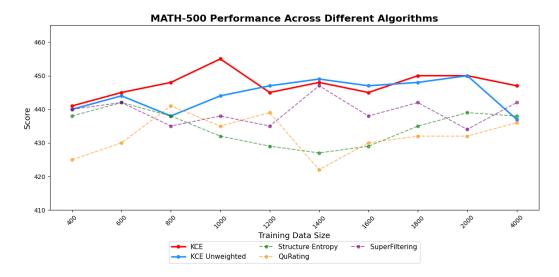


Figure 3: For LoRA fine-tuning, datasets of varying sizes were sampled using different algorithms. The red solid line represents KCE with knowledge-point weighting, the blue solid line represents KCE without weighting, and the remaining three dotted lines correspond to the other baseline data filtering algorithms.

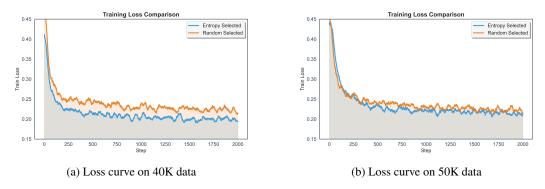


Figure 4: Training loss curves for models trained with entropy-based selection (blue) and random selection (brown) on 40K (a) and 50K (b) datasets. Entropy-based selection accelerates convergence and achieves lower, more stable training loss compared to random selection.

4.2.1 RETRIEVE WITH VARYING DATA SIZES

To evaluate the robustness of data selection algorithms under varying dataset sizes, we conducted experiments on the Diabetes and General Medical datasets using 1%, 5%, 10%, 20%, and 50% subsets. KCE consistently achieves the highest coverage and ranking quality across all scales. For instance in Table 2 shows, on the Diabetes dataset, KCE reaches 68.2%, 79.3%, 86.4%, and 90.1% coverage for the top 5, 10, 20, and 50 retrieved entries, respectively, outperforming Structure Entropy, QuRating, and Superfiltering by a significant margin. Similarly, on the Medical dataset, KCE attains 62.6%, 72.9%, 81.6%, and 89.0% coverage on the corresponding top retrieved entries, demonstrating a persistent advantage over all baselines. These results confirm that KCE robustly prioritizes high-value knowledge, even with limited or sparse data subsets.

4.2.2 Retrieve with different top-k

In this experiment, we evaluated retrieval performance on fixed-size datasets (Diabetes top 3K and Medical top 8K) by varying the top-k retrieved items from 5 to 50 to assess how well each algorithm ranks the most relevant knowledge. KCE consistently outperforms other methods across all top-k settings. For example in Table 3, on the Diabetes dataset, KCE achieves Top@10 coverage of

Table 1: Evaluation of data selection algorithms on RAG metrics. Reported metrics are average coverage rate and MRR. KCE achieves higher coverage and MRR with reduced dataset size compared to other methods.

Dataset	Algorithm	Avg. MRR	Avg. Coverage Rate	Data Size
	Full Dataset	0.4314	75.5%	12K
	KCE	0.4802	79.3%	3K
Diabetes	Structure Entropy	0.3372	61.9%	3K
	QuRating	0.3699	70.0%	3K
	Superfiltering	0.3695	68.7%	3K
	Full Dataset	0.4511	73.9%	20K
	KCE	0.4685	72.9%	8K
Medical	Structure Entropy	0.3952	67.6%	8K
	QuRating	0.3992	68.7%	8K
	Superfiltering	0.4227	69.2%	8K

Table 2: RAG evaluation of different data selection algorithms across varying dataset sizes (% of full dataset). Metrics reported are average coverage rate (%) and MRR. KCE consistently achieves higher coverage and MRR at all data sizes compared to other methods.

Dataset	Algorithm	Data Size (% of full dataset)					
Dataset	Aigoriumi	1%	5%	10%	20%	50%	
	KCE	59.4	62.8	66.9	77.10	79.8	
	KCE	0.3620	0.3562	0.3679	0.4452	0.4885	
	Structure Entropy	22.9	41.8	54.6	60.5	68.5	
Diabetes	Structure Entropy	0.1385	0.2516	0.3102	0.3229	0.3959	
Diabetes	QuRating	17.1	33.0	50.2	62.1	71.3	
	Quixating	0.08	0.17	0.2826	0.339	0.408	
	Cumanfiltanina	49.1	68.9	72.9	74.5	75.0	
	Superfiltering	0.2443	0.3731	0.4083	0.4416	0.4342	
	KCE	32.5	48.2	67.0	73.5	74.8	
	KCE	0.1927	0.2825	0.3855	0.4568	0.4805	
	Structure Entropy	22.1	46.7	58.0	63.7	69.5	
Medical	Structure Entropy	0.1393	0.2828	0.3343	0.3835	0.4143	
	QuRating	13.3	23.1	52.8	59.7	70.3	
	Quixating	0.0850	0.1300	0.2973	0.3641	0.4301	
	Superfiltering	41.4	50.9	65.9	70.0	73.9	
	Supermening	0.2324	0.2816	0.3843	0.4288	0.4511	

79.3% with MRR 0.4802, compared to Structure Entropy (61.9% / 0.3372), QuRating (70.0% / 0.3699), and Superfiltering (68.7% / 0.3695). Similarly, on the Medical dataset, KCE attains superior coverage and ranking quality across Top@5 to Top@50 (e.g., Top@50 coverage 89.0% with MRR 0.4762), demonstrating its persistent advantage in prioritizing high-value knowledge over competing algorithms.

These results demonstrate that KCE consistently outperforms other algorithms in retaining essential knowledge and improving retrieval quality. By effectively prioritizing high-value information and removing redundancy, KCE enables substantial dataset reduction without sacrificing performance, reducing computational cost and enhancing retrieval-augmented generation on both low- and high-dimensional medical datasets.

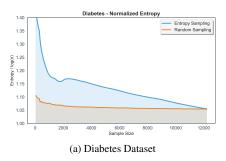
4.3 REVENUE BOUNDARIES AND INFORMATION GAIN ACROSS DOMAINS

Entropy-based sampling improves data utilization efficiency on both Diabetes and General Medical datasets. Normalized entropy curves show that entropy-selected subsets achieve higher information gain per sample than unfiltered data, with a clear revenue boundary beyond which additional samples provide diminishing returns. In low-sample regimes, steeper slopes indicate faster acquisition of high-value data, while flattening slopes mark diminishing marginal returns and a natural stopping

Table 3: RAG evaluation of different data selection algorithms across varying top retrieval sizes (Top@5, 10, 20, 50). Metrics reported are average coverage rate (%) and MRR. KCE consistently achieves higher coverage and MRR than other methods across all top-k settings.

Dataset	Algorithm	Retrieve Top Entries			
Dataset	Algorithm	Top@5	Top@10	Top@20	Top@50
	KCE	68.2	79.3	86.4	90.1
		0.4655	0.4802	0.4853	0.4867
	Structure Entropy	50.1	61.9	74.9	85.6
Diabetes top3K		0.3209	0.3372	0.3465	0.3510
Diabetes top3K	QuRating	54.4	70.0	81.6	89.5
		0.3491	0.3699	0.3778	0.3808
	Superfiltering	54.6	68.7	78.2	89.5
		0.3503	0.3695	0.3761	0.3798
	KCE Structure Entropy	62.6	72.9	81.6	89.0
		0.4539	0.4685	0.4737	0.4762
		57.0	67.6	76.0	84.9
Madical top9V		0.3810	0.3952	0.4018	0.4044
Medical top8K	QuRating	57.4	68.7	77.4	85.4
		0.3838	0.3992	0.4052	0.4079
	Superfiltering	58.0	69.2	77.1	84.7
	Supermering	0.4076	0.4227	0.4281	0.4308

criterion. Across domains, entropy-based sampling consistently attains higher coverage efficiency than random selection, enabling the construction of compact, high-quality datasets for LLM training and retrieval-augmented generation.



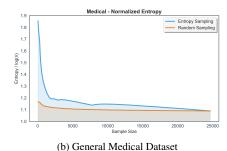


Figure 5: Comparison of normalized entropy vs. sample size for two datasets using entropy sampling (blue line with orange variability) and random sampling (brown line). (a) Diabetes Dataset (0–12,000 samples, entropy 1.0–1.4). (b) General Medical Dataset (0–25,000 samples, entropy 1.0–1.9). Entropy sampling consistently yields higher normalized entropy than random sampling.

5 DISCUSSION

In this work, we propose a knowledge-centric data selection framework for domain adaptation, formalized through Knowledge Coverage Entropy (KCE) and instantiated via an entropy-driven, submodular selection algorithm (EDS). The approach models discrete knowledge units and prioritizes coverage diversity under cardinality constraints, aiming to reduce redundancy and improve sample efficiency in both supervised fine-tuning and retrieval-augmented generation in domain adaptation of large language models. Empirical results on MATH-500 and medical RAG indicate consistent gains with smaller datasets and more stable training dynamics.

6 REPRODUCIBILITY STATEMENT

To ensure reproducibility of our results, all code used for data processing, model training, and evaluation will be provided in a zip file as part of the supplementary materials. Detailed descriptions of the datasets, preprocessing steps, and experimental settings are included in the main text and appendices. This will allow readers to reproduce the reported experiments and verify the findings.

REFERENCES

- Funpang medical dataset. https://huggingface.co/datasets/FunPang/medical_dataset. Accessed: 2025-08-26.
- Ishika Agarwal, Krishnateja Killamsetty, Lucian Popa, and Marina Danilevksy. Delift: Data efficient language model instruction fine tuning. *arXiv preprint arXiv:2411.04425*, 2024.
- Chen Amiraz, Florin Cuconasu, Simone Filice, and Zohar Karnin. The distracting effect: Understanding irrelevant passages in rag, 2025. URL https://arxiv.org/abs/2505.06914.
- Zachary Ankner, Cody Blakeney, Kartik Sreenivasan, Max Marion, Matthew L. Leavitt, and Mansheej Paul. Perplexed by perplexity: Perplexity-based data pruning with small reference models, 2024. URL https://arxiv.org/abs/2405.20541.
- Jianly Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu Lian, and Zheng Liu. Bge m3-embedding: Multi-lingual, multi-functionality, multi-granularity text embeddings through self-knowledge distillation, 2023.
- Rohan Deb, Kiran Thekumparampil, Kousha Kalantari, Gaurush Hiranandani, Shoham Sabach, and Branislav Kveton. Fishersft: Data-efficient supervised fine-tuning of language models using information gain. *arXiv preprint arXiv:2505.14826*, 2025.
- Mohsen Fayyaz, Ali Modarressi, Hinrich Schuetze, and Nanyun Peng. Collapse of dense retrievers: Short, early, and literal biases outranking factual evidence, 2025. URL https://arxiv.org/abs/2503.05037.
- Ronald A Fisher. On the mathematical foundations of theoretical statistics. *Philosophical transactions of the Royal Society of London. Series A, containing papers of a mathematical or physical character*, 222(594-604):309–368, 1922.
- Yanjun Fu, Faisal Hamman, and Sanghamitra Dutta. T-shirt: Token-selective hierarchical data selection for instruction tuning, 2025. URL https://arxiv.org/abs/2506.01317.
- Paul Hager, Friederike Jungmann, Robbie Holland, Kunal Bhagat, Inga Hubrecht, Manuel Knauer, Jakob Vielhauer, Marcus Makowski, Rickmer Braren, Georgios Kaissis, and Daniel Rueckert. Evaluation and mitigation of the limitations of large language models in clinical decision-making. *Nature Medicine*, 30(9):2613–2622, 2024. doi: 10.1038/s41591-024-03097-1.
- Zijian Hei, Weiling Liu, Wenjie Ou, Juyi Qiao, Junming Jiao, Guowen Song, Ting Tian, and Yi Lin. Dr-rag: Applying dynamic document relevance to retrieval-augmented generation for question-answering, 2024. URL https://arxiv.org/abs/2406.07348.
- Richard IG Holt and Allan Flyvbjerg. *Textbook of diabetes*. John Wiley & Sons, 2024.
- Jeremy Howard and Sebastian Ruder. Universal language model fine-tuning for text classification. *arXiv preprint arXiv:1801.06146*, 2018.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*, 2023.
- Didier Lairez. What entropy really is: the contribution of information theory, 2022. URL https://arxiv.org/abs/2204.05747.

- Katherine Lee, Daphne Ippolito, Andrew Nystrom, Chiyuan Zhang, Douglas Eck, Chris Callison-Burch, and Nicholas Carlini. Deduplicating training data makes language models better, 2022. URL https://arxiv.org/abs/2107.06499.
 - Ming Li and Paul Vitányi. *An introduction to Kolmogorov complexity and its applications*. Springer, 3 edition, 2008.
 - Ming Li, Yong Zhang, Shwai He, Zhitao Li, Hongyu Zhao, Jianzong Wang, Ning Cheng, and Tianyi Zhou. Superfiltering: Weak-to-strong data filtering for fast instruction-tuning. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 14255–14273, Bangkok, Thailand, August 2024. Association for Computational Linguistics. URL https://aclanthology.org/2024.acl-long.769.
 - Wei Liu, Weihao Zeng, Keqing He, Yong Jiang, and Junxian He. What makes good data for alignment? a comprehensive study of automatic data selection in instruction tuning. *arXiv* preprint *arXiv*:2312.15685, 2023.
 - Wei Liu, Weihao Zeng, Keqing He, Yong Jiang, and Junxian He. What makes good data for alignment? a comprehensive study of automatic data selection in instruction tuning, 2024. URL https://arxiv.org/abs/2312.15685.
 - Shayne Longpre, Gregory Yauney, Emily Reif, Katherine Lee, Adam Roberts, Barret Zoph, Denny Zhou, Jason Wei, Kevin Robinson, David Mimno, and Daphne Ippolito. A pretrainer's guide to training data: Measuring the effects of data age, domain coverage, quality, toxicity, 2023. URL https://arxiv.org/abs/2305.13169.
 - Christian Majenz. Entropy in quantum information theory communication and cryptography, 2018. URL https://arxiv.org/abs/1810.10436.
 - Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. s1: Simple test-time scaling. *arXiv preprint arXiv:2501.19393*, 2025a.
 - Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. s1: Simple test-time scaling, 2025b. URL https://arxiv.org/abs/2501.19393.
 - Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback, 2022. URL https://arxiv.org/abs/2203.02155.
 - Jinlong Pang, Na Di, Zhaowei Zhu, Jiaheng Wei, Hao Cheng, Chen Qian, and Yang Liu. To-ken cleaning: Fine-grained data selection for llm supervised fine-tuning. *arXiv preprint arXiv:2502.01968*, 2025.
 - Jupinder Parmar, Shrimai Prabhumoye, Joseph Jennings, Bo Liu, Aastha Jhunjhunwala, Zhilin Wang, Mostofa Patwary, Mohammad Shoeybi, and Bryan Catanzaro. Data, data everywhere: A guide for pretraining dataset construction, 2024. URL https://arxiv.org/abs/2407.06380.
 - Royal Government of Bhutan, Ministry of Health, Department of Medical Service. *Managing Diabetes Mellitus: Guide for Health Workers*. Royal Government of Bhutan, Ministry of Health, Department of Medical Service, Thimphu, Bhutan, August 2007.
 - Gerard Salton, Anita Wong, and Chung-Shu Yang. A vector space model for automatic indexing. *Communications of the ACM*, 18(11):613–620, 1975.
 - Skyler Seto, Maartje ter Hoeve, Maureen de Seyssel, and David Grangier. Assessing the role of data quality in training bilingual language models, 2025. URL https://arxiv.org/abs/2506.12966.

- KaShun Shum, Shizhe Diao, and Tong Zhang. Automatic prompt augmentation and selection with chain-of-thought from labeled data, 2024. URL https://arxiv.org/abs/2302.12822.
 - Yan Song, Prescott Klassen, Fei Xia, and Chunyu Kit. Entropy-based training data selection for domain adaptation. In *Proceedings of COLING 2012: Posters*, pp. 1191–1200, 2012.
 - Qwen Team. Qwen2.5: A party of foundation models, September 2024. URL https://qwenlm.github.io/blog/qwen2.5/.
 - Qwen Team. Qwen_max_0125. https://www.alibabacloud.com/help/en/model-studio/what-is-qwen-llm, 2025. Accessed: 2025-09-11.
 - Ellen M. Voorhees and Dawn M. Tice. The TREC-8 question answering track. In M. Gavrilidou, G. Carayannis, S. Markantonatou, S. Piperidis, and G. Stainhauer (eds.), *Proceedings of the Second International Conference on Language Resources and Evaluation (LREC'00)*, Athens, Greece, May 2000. European Language Resources Association (ELRA). URL https://aclanthology.org/L00-1018/.
 - Zige Wang, Wanjun Zhong, Yufei Wang, Qi Zhu, Fei Mi, Baojun Wang, Lifeng Shang, Xin Jiang, and Qun Liu. Data management for training large language models: A survey, 2024. URL https://arxiv.org/abs/2312.01700.
 - Alexander Wettig, Aatmik Gupta, Saumya Malik, and Danqi Chen. QuRating: Selecting high-quality data for training language models. In *International Conference on Machine Learning (ICML)*, 2024.
 - Tingyu Xia, Bowen Yu, Kai Dang, An Yang, Yuan Wu, Yuan Tian, Yi Chang, and Junyang Lin. Rethinking data selection at scale: Random selection is almost all you need. *arXiv preprint arXiv:2410.09335*, 2024.
 - Shitao Xiao, Zheng Liu, Peitian Zhang, and Niklas Muennighoff. C-pack: Packaged resources to advance general chinese embedding, 2023.
 - Tianchi Xie, Jiangning Zhu, Guozu Ma, Minzhi Lin, Wei Chen, Weikai Yang, and Shixia Liu. Structural-entropy-based sample selection for efficient and effective learning. *arXiv preprint arXiv:2410.02268*, 2024.

A USE OF LLM

The text of this article has been refined with the assistance of a large language model (LLM). All scholarly opinions, factual content, and final expressions remain the responsibility of the authors; the model was used solely to enhance the clarity, readability, and linguistic quality of the manuscript.

B LIMITATION

While the proposed entropy-driven Selection framework demonstrates promising results, it is not without limitations. First, the method assumes a certain degree of redundancy in the corpus, as the entropy computation relies on overlapping knowledge points across samples to establish informative distributions. Consequently, the approach may underperform on highly sparse datasets with minimal overlap. Second, the framework presumes that each information unit contains multiple knowledge points, providing sufficient variability to compute knowledge-point entropy. In cases where samples are extremely atomic—e.g., containing only a single knowledge point—the resulting knowledge matrix becomes nearly diagonal, rendering Knowledge Coverage Entropy computation ineffective.

Additionally, we currently represent the Information–Knowledge Matrix as binary, indicating whether a sample fully covers a knowledge point or not. While this simplification facilitates computation and aligns with the current entropy formulation, it neglects partial or graded coverage. We acknowledge this limitation and note that a probabilistic or weighted representation could better capture the degree of knowledge coverage in future work.

B.1 SINGLE-PASS SCORING APPROXIMATION

For datasets with large n, a computationally efficient approximation ranks samples using a single-pass scoring function, independent of set interactions. Let $\mathbf{B} \in \{0,1\}^{n \times m}$ be the binary information-knowledge matrix, where n is the number of samples, m is the number of knowledge points, and $\mathbf{B}_{a,j} = 1$ if sample a covers knowledge point j, and 0 otherwise. The row coverage probability for sample a is

$$P_a = \frac{1}{m} \sum_{j=1}^{m} \mathbf{B}_{a,j},$$

with entropy $H(a) = -P_a \log_2 P_a$. A weight vector $\mathbf{k} \in \mathbb{R}^m$ assigns importance k_i to knowledge point i. The score for sample a is

Score(a) =
$$H(a) \cdot \left(1 + \gamma \sum_{i=1}^{m} k_i \mathbf{B}_{a,i}\right)$$
,

where γ is a balancing parameter. The top-s samples are selected based on this score. This approach, while linear in time, lacks submodular guarantees due to its neglect of set interactions, unlike the lazy-greedy method described in the main text.

B.2 Knowledge Coverage Entropy Definition and Bounds

The Knowledge Coverage Entropy (KCE) measures the diversity of knowledge coverage in a subset $S \subseteq \{1,\ldots,n\}$ of size |S|=h from a dataset represented by a binary matrix $\mathbf{B} \in \{0,1\}^{n\times m}$, where $\mathbf{B}_{i,j}=1$ if sample i covers knowledge point j, and 0 otherwise. To ensure numerical stability, we apply additive smoothing:

$$\mathbf{B}' = \mathbf{B} + \frac{\alpha}{nm}, \qquad \alpha = 10^{-6},$$

and normalize to obtain a joint probability distribution:

$$P_{i,j} = \frac{\mathbf{B}'_{i,j}}{\sum_{i=1}^{n} \sum_{j=1}^{m} \mathbf{B}'_{i,j}}.$$

The KCE for subset S is defined as

$$H(S) = -\sum_{j=1}^{m} p_j \log_2 p_j, \qquad p_j = \frac{1}{h} \sum_{a \in S} \mathbf{B}_{a,j}.$$

The maximum entropy occurs when $p_i = 1/m$, yielding $H(S) \le \log_2 m$. For the joint distribution over nm outcomes, the upper bound is

$$H(S) \le \log_2(nm) = \log_2 n + \log_2 m.$$

The normalized entropy is

$$H_n(S) = \frac{H(S)}{\log_2 h},$$

with $H_n(S) \leq 1 + \frac{\log_2 m}{\log_2 n}$. As $n \to \infty$, $H_n(S) \to 1$ (or less with redundancy). Redundancy in **B** (e.g., samples covering identical points) reduces H(S) to $\approx \log_2 m$, exhibiting sublinear growth.

B.3 Entropy-Driven Selection with Lazy-Greedy Algorithm

Algorithm 1: Lazy-Greedy EDS (Set-Aware Selection)

Input: Binary matrix $\mathbf{B} \in [0, 1]^{n \times m}$; weights $w \in \mathbb{R} + [n]^m$; budget s; concave f; tolerance $\varepsilon \ge 0$ **Output:** Selected indices S

 $S \leftarrow \emptyset; c \leftarrow \mathbf{0}m;$

// Coverage counts

for $a \in 1, \ldots, n$ do

Compute initial upper bound U_a on $\Delta F(\emptyset; a)$; Push (a, U_a) into max-heap \mathcal{H} ;

```
(a, U_a) \leftarrow \text{PopMax}(\mathcal{H}); // \text{ Exact marginal gain using current } c
g_a \leftarrow \sum j : \mathbf{B}a, j = 1w_j [f(c_j + 1) - f(c_j)]; U_{\max} \leftarrow \mathsf{CurrentMaxKey}(\mathcal{H}) \text{ (or } -\infty \text{ if }
 empty); if g_a \ge U_{\max} - \varepsilon then
    else
```

Push (a, g_a) back into \mathcal{H} ;

return S

SUBMODULARITY OF KNOWLEDGE COVERAGE ENTROPY

The effectiveness of the greedy algorithm relies on the submodular properties of KCE. Let $\mathbf{B} \in$ $\{0,1\}^{n\times m}$ be the binary matrix, and $H(S)=-\sum_{j=1}^m p_j\log_2 p_j$ the KCE for subset S, where $p_j = \frac{1}{|S|} \sum_{a \in S} \mathbf{B}_{a,j}$. Although KCE is not strictly submodular, it exhibits diminishing marginal gains. For nested subsets $S_A \subseteq S_B$ and a sample $a \notin S_B$, the marginal gain satisfies

$$\Delta H(S_A; a) = H(S_A \cup \{a\}) - H(S_A) \ge \Delta H(S_B; a).$$

To derive this, consider the entropy function $H(p) = -\sum_{i} p_{i} \log_{2} p_{j}$, which is concave in the probability vector p. When adding sample a to subset S, define the coverage distribution induced by a as

$$\delta_j = \frac{\mathbf{B}_{a,j}}{\sum_j \mathbf{B}_{a,j}}, \qquad c_a = \sum_j \mathbf{B}_{a,j},$$

and let $K(S) = \sum_{a \in S} \sum_{j} \mathbf{B}_{a,j}$ be the total coverage of S. The mixing parameter is

$$\lambda = \frac{c_a}{K(S) + c_a}.$$

The updated probability vector p' is a convex combination:

$$p_i' = (1 - \lambda)p_i + \lambda \,\delta_i.$$

Since H(p) is concave, by Jensen's inequality applied to the convex combination,

$$H(p') \geq (1 - \lambda)H(p) + \lambda H(\delta),$$

which yields the marginal gain bound

$$\Delta H = H(p') - H(p) \ge \lambda (H(\delta) - H(p)).$$

Because the Hessian of H(p) is negative semi-definite, entropy changes are smaller when p is near uniform (as in larger sets). For $|S_B| > |S_A|$, $K(S_B) > K(S_A)$, so $\lambda_B < \lambda_A$, and the distribution p_B is closer to uniform, reducing $\Delta H(S_B;a)$. Alternatively, one can approximate $\Delta H \approx -D_{\mathrm{KL}}(p' \parallel p)$, where D_{KL} decreases with set size due to smaller λ , reinforcing the inequality for non-redundant a. This property supports the greedy algorithm's effectiveness, as detailed in the main text.

B.5 Information Gain and Revenue Boundary

The information gain (IG) monitors the marginal contribution of adding samples to a subset. For a binary matrix $\mathbf{B} \in \{0,1\}^{n \times m}$ and subset S_t of size t, the normalized entropy is $H_n(S_t) = H(S_t)/\log_2 t$, where $H(S_t) = -\sum_{j=1}^m p_j \log_2 p_j$ and $p_j = \frac{1}{t} \sum_{a \in S_t} \mathbf{B}_{a,j}$. The discrete information gain is

$$IG(t) = H_n(S_t) - H_n(S_{t-1}).$$

Due to diminishing returns (see Appendix B.4), IG(t) decays as t increases. The revenue boundary is defined as

$$t^* = \min\{t : IG(t) < \delta\},\$$

where $\delta > 0$ is a task-specific threshold. To derive the decay, note that entropy is subadditive: for a new sample a with row entropy $H(a) = -P_a \log_2 P_a$, where $P_a = \frac{1}{m} \sum_{j=1}^m \mathbf{B}_{a,j}$,

$$H(S \cup \{a\}) \le H(S) + H(a),$$

and $H(a) \leq \log_2 m$ for uniform coverage. The marginal gain is

$$\Delta H = H(S \cup \{a\}) - H(S) \le H(a).$$

Accounting for redundancy,

$$\Delta H = H(a \mid S) = H(a) - I(a; S),$$

where I(a; S) is the mutual information measuring overlap. For large t, the expected $\Delta H_t \frac{\log_2 m}{t}$, as new samples cover at most m/t new points on average (pigeonhole principle). Entropy concavity implies successive gains diminish:

$$\Delta H_t \le \frac{\Delta H_{t-1}}{1+\epsilon}, \quad \epsilon > 0,$$

in redundant regimes. Summing the series,

$$H(S_t) = H(S_1) + \sum_{k=2}^{t} \Delta H_k \le H(S_1) + \sum_{k=2}^{t} O\left(\frac{1}{k}\right) = H(S_1) + O(\log t).$$

Thus, $H_n(S_t) = O(1)$, and

$$IG(t) \approx \frac{\Delta H_t}{\log_2 t} = O\left(\frac{1}{t \log t}\right),$$

which asymptotically simplifies to O(1/t). This decay justifies the revenue boundary for efficient stopping.

B.6 MUTUAL INFORMATION APPROXIMATION

Maximizing KCE approximates maximizing mutual information I(R;C) between samples (rows R) and knowledge points (columns C). Let $\mathbf{B} \in \{0,1\}^{n \times m}$ be the binary matrix, and $S \subseteq \{1,\ldots,n\}$ a subset. Define R as a uniform random variable over S and C as a knowledge point conditioned or coverage. The joint entropy is H(R,C) = H(S), where $H(S) = -\sum_{j=1}^m p_j \log_2 p_j$, $p_j = \frac{1}{|S|} \sum_{a \in S} \mathbf{B}_{a,j}$. The mutual information is

$$I(R; C) = H(R) + H(C) - H(R, C) = \log_2 |S| + H(C) - H(S)$$

where $H(R) = \log_2 |S|$ (uniform over rows) and $H(C) = -\sum_{i=1}^m P(\cdot, j) \log_2 P(\cdot, j)$, with

$$P(\cdot, j) = \frac{1}{|S|} \sum_{a \in S} \mathbf{B}_{a,j}$$
 (column marginals).

Maximizing I(R;C) requires maximizing H(C) (broad coverage) while minimizing H(S) (low redundancy). Per-sample row entropy $H(a) = -P_a \log_2 P_a$, where $P_a = \frac{1}{m} \sum_{j=1}^m \mathbf{B}_{a,j}$, peaks at $P_a \approx 0.5$, favoring balanced samples that diversify C and reduce $H(C \mid R)$. Under row independence, H(S) = H(R) + H(C), so I(R;C) = 0; selection induces correlations, increasing I. The score

$$Score(a) = H(a) \cdot \left(1 + \gamma \sum_{i=1}^{m} k_i \mathbf{B}_{a,i}\right)$$

prioritizes task-relevant balance, approximating greedy I(R;C) maximization (similar to submodular set cover).

B.7 Data Distribution Effects in Supervised Fine-Tuning

In supervised fine-tuning (SFT), let z denote the logits, p_{θ} the predicted probability via softmax, q the target distribution, and L the cross-entropy loss (Ouyang et al., 2022):

$$L(\theta) = -\sum_{i=1}^{m} q_i \log p_{\theta,i}, \quad p_{\theta,i} = \frac{e^{z_i}}{\sum_{j} e^{z_j}}.$$
 (5)

The gradient with respect to logits is

$$\nabla_z L = p_\theta - q. \tag{6}$$

The Fisher information matrix with respect to logits (Fisher, 1922) is defined as

$$F_z(q) = \mathbb{E}Y \sim q \left[(\nabla_z L(Y))(\nabla_z L(Y))^\top \right], \tag{7}$$

where Y is a one-hot random variable drawn from q. Expanding this gives

$$F_z(q) = (p\theta - q)(p_\theta - q)^\top + \text{Cov}(Y).$$
(8)

Near convergence, $p_{\theta} \approx q$, so the rank-one term vanishes, and we have

$$F_z(q) \approx \text{Cov}(Y) = \text{diag}(q) - qq^{\top}.$$
 (9)

The expected squared gradient norm is

$$\mathbb{E}[||\nabla_z L||^2] = \text{Tr}(F_z(q)) = 1 - \sum_i i = 1^m q_i^2,$$
(10)

which is maximized for uniform q (high diversity) and minimized for skewed q (low diversity).

From this perspective, selecting datasets with high Knowledge Coverage Entropy (KCE) promotes a more uniform empirical knowledge distribution $p_j(S)$, ensuring that minibatches sampled from S maintain high average gradient norms. This leads to faster and more stable convergence during SFT by avoiding overly skewed label distributions that would produce weak learning signals. In other words, maximizing row entropy H(q) through KCE naturally aligns the data distribution to enhance both gradient strength and training efficiency.

B.7.1 EFFICIENCY IN MODEL TRAINING

To validate the Revenue Boundary Theory, we prepared two sets of sampled datasets: (1) 28 subsets randomly sampled from the original dataset, with sizes ranging from 100 to 50,000; and (2) 28 subsets selected using the Entropy-Driven Data Selection algorithm. We trained 56 models in total using these datasets and visualized their performance trends.

B.7.2 NORMALIZED ENTROPY AND INFORMATION GAIN

We conducted experiments on mathematical dataset by applying the proposed Entropy-Driven Data Selection algorithm to generate subsets with sizes ranging from 100 to 30,000. For each subset, we computed the normalized entropy and visualized its variation trend as the sample size increased. Furthermore, we plotted the information gain efficiency curves for both datasets to illustrate the points of maximum efficiency.

B.8 ALTERNATIVE JOINT ENTROPY FORMULATION

An alternative joint entropy formulation is

$$H_{\text{joint}}(S) = -\sum_{i \in S} \sum_{j=1}^{m} P_{i,j} \log_2 P_{i,j},$$

where

$$P_{i,j} = \frac{\mathbf{B}'_{i,j}}{\sum_{i=1}^{n} \sum_{j=1}^{m} \mathbf{B}'_{i,j}}, \quad \mathbf{B}' = \mathbf{B} + \frac{\alpha}{nm}, \quad \alpha = 10^{-6}.$$

This accounts for row and column dependencies but is computationally costly and sensitive to redundancy. The marginal KCE in the main text is more efficient for diversity-focused selection.

B.9 STOCHASTIC-GREEDY VARIANT

A stochastic-greedy variant samples a subset R of size $r \approx \frac{n}{s} \log \frac{1}{s}$ at each iteration, selecting

$$a^* = \arg\max_{a \in R} \Delta F(S; a).$$

This achieves a $(1 - 1/e - \varepsilon)$ guarantee with reduced computational cost.

B.10 Hybrid Objective

A hybrid objective combines coverage and similarity:

$$F_{\text{hybrid}}(S) = \lambda \sum_{j=1}^{m} w_j f(c_j(S)) + (1 - \lambda) \sum_{x=1}^{n} \max_{a \in S} \text{sim}(x, a),$$

where f is concave, and the second term is a facility-location function over a similarity graph. Both terms are submodular, preserving the (1-1/e) guarantee of the lazy-greedy algorithm.

B.11 PARAMETER SENSITIVITY ANALYSIS

We conducted a sensitivity analysis of KCE with respect to the smoothing parameter α and the weight balance γ on sample sizes 500 and 1000. The normalized KCE (H_n) remains nearly constant across $\alpha \in [0.1, 2.0]$ and $\gamma \in [0, 1]$. For instance, with sample size 500, H_n varies only from 1.2300 to 1.2320 (< 0.2%), and with sample size 1000, from 1.2037 to 1.2062 (< 0.3%). These small variations indicate that KCE is robust to both α and γ , and the algorithm reliably selects diverse knowledge subsets without significant sensitivity to hyperparameter choices.

B.12 EMPIRICAL VALIDATION VIA SIMULATIONS

To empirically validate the decay in information gain, experiments were conducted on four datasets, with knowledge points m ranging from 200 to 1000 and sample sizes n between 20,000 and 60,000, averaged over 5 runs. Across all datasets, the normalized entropy grows sublinearly:

$$H_n(S) = \frac{H(S)}{\log m}, \qquad H(S) = -\sum_{i=1}^m p_i \log p_i,$$

where p_i denotes the empirical frequency of knowledge point i in the subset S. The information gain (IG) at step t is defined as the marginal increase in entropy:

$$IG(t) = H_n(S_t) - H_n(S_{t-1}), S_t = S_{t-1} \cup \{x_t\}.$$

Empirically, IG(t) starts high (approximately 0.99 at t=1) and decays to near-zero (around 10^{-7} by t=1000), following an overall O(1/t) trend:

$$IG(t) \approx \frac{c}{t}, \qquad c > 0.$$

Moreover, we examine the slope of IG(t), i.e., its discrete derivative:

$$\Delta IG(t) = IG(t+1) - IG(t).$$

On the diabetes dataset, the slope decreases from approximately

$$\Delta IG \approx -0.1 \times 10^{-4}$$
 (at the best advantage point)

to

$$\Delta IG \approx -1 \times 10^{-6}$$
, after which it stabilizes.

This behavior confirms diminishing returns and validates the revenue boundary condition, where

$$IG(t) < \delta$$
 for $t \ge T_{\delta}$.

Simulations on random binary matrices (e.g., m=251, varying n) show $H_n(S)$ peaking early and IG(n) decaying from ~ 1.22 to near-zero, confirming theorems. For entropy-selected subsets, I(R;C) is 10-20% higher than random, tying theory to empirical wins.

B.13 COMPUTATIONAL COMPLEXITY ANALYSIS

The lazy-greedy algorithm has time complexity that scales with the sparsity of the matrix \mathbf{B} . Each exact marginal gain evaluation is $O(\operatorname{nnz}(\mathbf{B}_{a,\cdot}))$, where nnz denotes the number of non-zero entries in row a. The lazy variant reduces the number of full evaluations by using upper bounds in the heap, leading to near-linear time in the total number of non-zero entries in \mathbf{B} for sparse matrices. For dense matrices, the complexity is $O(nm\log n)$ in the worst case, but practical datasets are often sparse. The single-pass approximation is O(nm), linear in the matrix size. Memory requirements are O(n+m) for the heap and counts, making it scalable for large n and m.

Remarks and interpretation.

- The quantity $1 \sum_i q_i^2$ is closely related to the Gini impurity and measures the distributional uncertainty: it is zero for a one-hot (deterministic) q and maximized when q is uniform.
- The derivation above is performed in the *logit space*. For the Fisher information with respect to model parameters θ , one needs to apply the Jacobian chain rule $F_{\theta} = J_{z \to \theta}^{\top} F_z J_{z \to \theta}$; nevertheless, the qualitative conclusion—uncertainty in q increases the expected gradient magnitude—remains valid.
- The approximation $F_z(q) \approx \operatorname{Cov}(Y)$ relies on $p \approx q$. When the model is far from well calibrated, the additional term $(p-q)(p-q)^{\top}$ may be non-negligible and should be accounted for.

C KNOWLEDGE DISTIL PROMPT

Prompt Example:

- "You are a medical knowledge summarization assistant."
- "Task: Summarize the given text and extract only concise knowledge points directly related to diabetes."
- "Requirements:"
- "1. Focus only on diabetes and its directly related aspects (symptoms, complications, treatments, risk factors, diagnostic methods, pathophysiology)."
- "2. If there is **no diabetes-related content**, output exactly: NO"
- "3. The output must consist only of short words or phrases (concise terms)."
- "4. Do not output personal names, study names, or any content unrelated to diabetes."
- "5. Do not add external knowledge, only use the given content."
- "6. Output multiple knowledge points separated by commas, without extra text or explanations."

"Example: In diabetes management, α -glucosidase inhibitors may cause gastrointestinal side effects such as flatulence, abdominal discomfort, and diarrhea, particularly with high doses relative to carbohydrate intake, but these improve with gradual titration. Hypoglycemia is rare, and drug interactions are minimal, though concomitant use with motility agents or cholestyramine is not recommended. In the STOP-NIDDM trial, 31% of acarbose-treated patients discontinued early due to adverse effects compared to 19% with placebo."

"Response: α -glucosidase inhibitors, gastrointestinal side effects, hypoglycemia rare, acarbose"

D KNOWLEDGE POINT EXAMPLE

diabetes	cardiovascular disease	genetic factors	insulin resistance
insulin secretion	diabetes complications	diabetes management	diagnostic criteria
type 2 diabetes	GLP-1	cardiovascular benefits	complications
diabetes care	HbA1c	type 1 diabetes	insulin
diabetic complications	cardiovascular risk	chronic kidney disease	continuous glucose monitoring
coronary heart disease	congestive heart failure	diabetic ketoacidosis	diabetic retinopathy
GIP	impaired fasting glucose	impaired glucose tolerance	glucose metabolism
hyperglycemia	beta-cell dysfunction	liraglutide	oral glucose tolerance test
endothelial dysfunction	oxidative stress	proliferative diabetic retinopa-	glucose monitoring
		thy	
inflammation	neuropathy	retinopathy	type 1 diabetes mellitus
type 2 diabetes mellitus	thiazolidinediones	insulin sensitivity	urinary albumin excretion
Diabetes	polyuria	diabetic neuropathy	albuminuria
sulfonylureas	metformin	blood pressure management	microaneurysms
peripheral neuropathy	nephropathy	nausea	insulin deficiency
mortality	hepatic glucose production	blood glucose regulation	insulin production
glucose regulation	Diabetic retinopathy	diabetes treatment	insulin therapy
hypoglycemia risk	renal failure	diabetic nephropathy	Type 1 diabetes
MODY	pancreatic beta cells	hypertension	glucagon
lifestyle modifications	oral hypoglycemic agents	proteinuria	glycemic control
insulin pumps	rosiglitazone	pioglitazone	severe hypoglycemia
insulin use	cardiovascular risk reduction	microalbuminuria	blood pressure control
UKPDS	blood glucose control	insulin treatment	patient education
pregnancy	gestational diabetes mellitus	fasting plasma glucose	chronic hyperglycemia
microvascular complications	foot ulcers	macrovascular disease	weight loss
ketoacidosis	gestational diabetes	Type 2 diabetes T2DM	elevated blood glucose
increased diabetes risk	T1DM disease progression		Type 1 diabetes mellitus physical inactivity
Type 2 diabetes mellitus	1 0	obesity Gestational diabetes mellitus	macrosomia
glucocorticoids type 2 diabetes mellitus	infections	DCCT Gestational diabetes mellitus	islet autoantibodies
type 2 diabetes mellitus (T2DM)	fasting glucose	DCC1	isiet autoantibodies
children	environmental factors	increased risk	islet autoimmunity
prevention	weight gain	Diabetic ketoacidosis	hypoglycemia
cardiovascular mortality	risk factors	macrovascular complications	vascular complications
metabolic syndrome	dyslipidemia	diabetes risk	physical activity
family history	β -cell dysfunction	free fatty acids	impaired insulin secretion
type 2 diabetes risk	diabetes prevalence	Diabetes prevalence	undiagnosed diabetes
smoking	age	glucose intolerance	coronary artery disease
early detection	adolescents	risk factor	stroke
infection	cardiovascular risk factors	end-stage renal disease	myocardial infarction
diet	quality of life	lifestyle interventions	weight reduction
ACE inhibitors	diabetes prevention	genetic predisposition	exercise
smoking cessation	Metformin	low- and middle-income coun-	overweight
		tries	
alcohol consumption	hyperinsulinemia	insulin administration	depression
combination therapy	blood pressure	comorbidities	glucose homeostasis
	hypoglycemia prevention	blood glucose levels	Sulfonylureas
insulin release			
β -cell function	DPP-4 inhibitors	exenatide	glucose uptake
β-cell function lipolysis	DPP-4 inhibitors insulin signaling	exenatide impaired glucose metabolism	glucose uptake skeletal muscle
β -cell function lipolysis TNF- α	DPP-4 inhibitors insulin signaling adipose tissue	exenatide impaired glucose metabolism growth hormone	glucose uptake skeletal muscle liver
β -cell function lipolysis TNF- α gluconeogenesis	DPP-4 inhibitors insulin signaling adipose tissue counter-regulatory hormones	exenatide impaired glucose metabolism growth hormone monogenic diabetes	glucose uptake skeletal muscle liver insulin secretagogues
eta-cell function lipolysis TNF- $lpha$ gluconeogenesis NAFLD	DPP-4 inhibitors insulin signaling adipose tissue counter-regulatory hormones atherosclerosis	exenatide impaired glucose metabolism growth hormone monogenic diabetes neonatal diabetes	glucose uptake skeletal muscle liver insulin secretagogues increased mortality
eta-cell function lipolysis TNF- $lpha$ gluconeogenesis NAFLD HbA1c levels	DPP-4 inhibitors insulin signaling adipose tissue counter-regulatory hormones atherosclerosis metabolic control	exenatide impaired glucose metabolism growth hormone monogenic diabetes neonatal diabetes sulfonylurea	glucose uptake skeletal muscle liver insulin secretagogues increased mortality glucagon suppression
β-cell function lipolysis TNF-α ggluconeogenesis NAFLD HbA1c levels delayed gastric emptying	DPP-4 inhibitors insulin signaling adipose tissue counter-regulatory hormones atherosclerosis metabolic control osmotic diuresis	exenatide impaired glucose metabolism growth hormone monogenic diabetes neonatal diabetes sulfonylurea vomiting	glucose uptake skeletal muscle liver insulin secretagogues increased mortality glucagon suppression HbA1c reduction
β-cell function lipolysis TNF-α gluconeogenesis NAFLD HbA1c levels delayed gastric emptying glibenclamide	DPP-4 inhibitors insulin signaling adipose tissue counter-regulatory hormones atherosclerosis metabolic control osmotic diuresis α-glucosidase inhibitors	exenatide impaired glucose metabolism growth hormone monogenic diabetes neonatal diabetes sulfonylurea vomiting asymptomatic	glucose uptake skeletal muscle liver insulin secretagogues increased mortality glucagon suppression HbA1c reduction cardiovascular disease risk
β-cell function lipolysis TNF-α gluconeogenesis NAFLD HbA1c levels delayed gastric emptying glibenclamide dehydration	DPP-4 inhibitors insulin signaling adipose tissue counter-regulatory hormones atherosclerosis metabolic control osmotic diuresis α-glucosidase inhibitors heart failure	exenatide impaired glucose metabolism growth hormone monogenic diabetes neonatal diabetes sulfonylurea vomiting asymptomatic lifestyle changes	glucose uptake skeletal muscle liver insulin secretagogues increased mortality glucagon suppression HbA1c reduction cardiovascular disease risk drug interactions
eta-cell function lipolysis TNF- $lpha$ gluconeogenesis NAFLD HbA1c levels delayed gastric emptying glibenclamide dehydration blood glucose monitoring	DPP-4 inhibitors insulin signaling adipose tissue counter-regulatory hormones atherosclerosis metabolic control osmotic diuresis α-glucosidase inhibitors heart failure insulin dose adjustment	exenatide impaired glucose metabolism growth hormone monogenic diabetes neonatal diabetes sulfonylurea vomiting asymptomatic lifestyle changes hypoglycemia unawareness	glucose uptake skeletal muscle liver insulin secretagogues increased mortality glucagon suppression HbA1c reduction cardiovascular disease risk drug interactions polydipsia
β-cell function lipolysis TNF-α gluconeogenesis NAFLD HbA1c levels delayed gastric emptying glibenclamide dehydration blood glucose monitoring screening	DPP-4 inhibitors insulin signaling adipose tissue counter-regulatory hormones atherosclerosis metabolic control osmotic diuresis α-glucosidase inhibitors heart failure insulin dose adjustment HbA1c	exenatide impaired glucose metabolism growth hormone monogenic diabetes neonatal diabetes sulfonylurea vomiting asymptomatic lifestyle changes hypoglycemia unawareness self-management	glucose uptake skeletal muscle liver insulin secretagogues increased mortality glucagon suppression HbA1c reduction cardiovascular disease risk drug interactions polydipsia diabetes education
β-cell function lipolysis TNF-α gluconeogenesis NAFLD HbA Ic levels delayed gastric emptying glibenclamide dehydration blood glucose monitoring screening autonomic neuropathy	DPP-4 inhibitors insulin signaling adipose tissue counter-regulatory hormones atherosclerosis metabolic control osmotic diuresis α-glucosidase inhibitors heart failure insulin dose adjustment HbA1c gastroparesis	exenatide impaired glucose metabolism growth hormone monogenic diabetes neonatal diabetes sulfonylurea vomiting asymptomatic lifestyle changes hypoglycemia unawareness self-management eating disorders	glucose uptake skeletal muscle liver insulin secretagogues increased mortality glucagon suppression HbA1c reduction cardiovascular disease risk drug interactions polydipsia diabetes education erectile dysfunction
β-cell function lipolysis TNF-α gluconeogenesis NAFLD HbA1c levels delayed gastric emptying glibenclamide dehydration blood glucose monitoring screening autonomic neuropathy diabetes self-management	DPP-4 inhibitors insulin signaling adipose tissue counter-regulatory hormones atherosclerosis metabolic control osmotic diuresis α-glucosidase inhibitors heart failure insulin dose adjustment HbA1c gastroparesis blood glucose management	exenatide impaired glucose metabolism growth hormone monogenic diabetes neonatal diabetes sulfonylurea vomiting asymptomatic lifestyle changes hypoglycemia unawareness self-management eating disorders carbohydrate intake	glucose uptake skeletal muscle liver insulin secretagogues increased mortality glucagon suppression HbA1c reduction cardiovascular disease risk drug interactions polydipsia diabetes education erectile dysfunction insulin-treated diabetes
β-cell function lipolysis TNF-α gluconeogenesis NAFLD HbA1c levels delayed gastric emptying glibenclamide dehydration blood glucose monitoring screening autonomic neuropathy diabetes self-management glycaemic control	DPP-4 inhibitors insulin signaling adipose tissue counter-regulatory hormones atherosclerosis metabolic control osmotic diuresis α-glucosidase inhibitors heart failure insulin dose adjustment HbA1c gastroparesis blood glucose management cardiovascular events	exenatide impaired glucose metabolism growth hormone monogenic diabetes neonatal diabetes sulfonylurea vomiting asymptomatic lifestyle changes hypoglycemia unawareness self-management eating disorders carbohydrate intake cardiovascular outcomes	glucose uptake skeletal muscle liver insulin secretagogues increased mortality glucagon suppression HbA1c reduction cardiovascular disease risk drug interactions polydipsia diabetes education erectile dysfunction insulin-treated diabetes clinical trials
β-cell function lipolysis TNF-α gluconeogenesis NAFLD HbA1c levels delayed gastric emptying glibenclamide dehydration blood glucose monitoring screening autonomic neuropathy diabetes self-management glycaemic control renal impairment	DPP-4 inhibitors insulin signaling adipose tissue counter-regulatory hormones atherosclerosis metabolic control osmotic diuresis α-glucosidase inhibitors heart failure insulin dose adjustment HbA1c gastroparesis blood glucose management cardiovascular events cognitive impairment	exenatide impaired glucose metabolism growth hormone monogenic diabetes neonatal diabetes sulfonylurea vomiting asymptomatic lifestyle changes hypoglycemia unawareness self-management eating disorders carbohydrate intake cardiovascular outcomes anxiety	glucose uptake skeletal muscle liver insulin secretagogues increased mortality glucagon suppression HbA1c reduction cardiovascular disease risk drug interactions polydipsia diabetes education erectile dysfunction insulin-treated diabetes clinical trials Thiazolidinediones
β-cell function lipolysis TNF-α gluconeogenesis NAFLD HbA1c levels delayed gastric emptying glibenclamide dehydration blood glucose monitoring screening autonomic neuropathy diabetes self-management glycaemic control renal impairment GLP-1 receptor agonists	DPP-4 inhibitors insulin signaling adipose tissue counter-regulatory hormones atherosclerosis metabolic control osmotic diuresis α-glucosidase inhibitors heart failure insulin dose adjustment HbA1c gastroparesis blood glucose management cardiovascular events cognitive impairment basal insulin	exenatide impaired glucose metabolism growth hormone monogenic diabetes neonatal diabetes sulfonylurea vomiting asymptomatic lifestyle changes hypoglycemia unawareness self-management eating disorders carbohydrate intake cardiovascular outcomes anxiety gastrointestinal side effects	glucose uptake skeletal muscle liver insulin secretagogues increased mortality glucagon suppression HbA1c reduction cardiovascular disease risk drug interactions polydipsia diabetes education erectile dysfunction insulin-treated diabetes clinical trials Thiazolidinediones depressive symptoms
β-cell function lipolysis TNF-α gluconeogenesis NAFLD HbA1c levels delayed gastric emptying glibenclamide dehydration blood glucose monitoring screening autonomic neuropathy diabetes self-management glycaemic control renal impairment	DPP-4 inhibitors insulin signaling adipose tissue counter-regulatory hormones atherosclerosis metabolic control osmotic diuresis α-glucosidase inhibitors heart failure insulin dose adjustment HbA1c gastroparesis blood glucose management cardiovascular events cognitive impairment	exenatide impaired glucose metabolism growth hormone monogenic diabetes neonatal diabetes sulfonylurea vomiting asymptomatic lifestyle changes hypoglycemia unawareness self-management eating disorders carbohydrate intake cardiovascular outcomes anxiety	glucose uptake skeletal muscle liver insulin secretagogues increased mortality glucagon suppression HbA1c reduction cardiovascular disease risk drug interactions polydipsia diabetes education erectile dysfunction insulin-treated diabetes clinical trials Thiazolidinediones

E SUPPLEMENTARY TABLES

Table 4: Main Training Configurations and LoRA Parameter

Training Configurations

LoRA Parameter

Parameter	Value	Parameter	Value
Learning Rate	1×10^{-5}	Learning Rate	5×10^{-5}
Epochs	5	Epochs	5
Method	Full Parameter	Method	LoRA
Model	Qwen2.5_32B_Instruct	PEFT	q,k,v,o,down,gate,up
Evaluation Dataset	MATH-500	LoRA Rank	64
Early Stop Threshold	$Loss \leq 0.05$	LoRA Alpha	128
Early Stop Patience	5 steps	LoRA Dropout	0.1
Deepspeed Stage	Zero Stage 3	Early Stop Threshold	$Loss \leq 0.15$
-		Early Stop Patience	5 steps
		Deepspeed Stage	Zero Stage 3

Table 5: MATH-500 Performance Across Different Sample Sizes and different Algorithms on LORA fine tune

Training Data Size	KCE	KCE Unweighted	Struct Entropy	QuRating	Superfiltering
400	441	440	438	425	440
600	445	444	442	430	442
800	448	438	438	441	435
1000	455	444	432	435	438
1200	445	447	429	439	435
1400	448	449	427	422	447
1600	445	447	429	430	438
1800	450	448	435	432	442
2000	450	450	439	432	434
4000	447	437	438	436	442

Table 6: MATH-500 Performance Across Different Sample Sizes with full parameter fine tuning. The table reports scores for models trained on subsets selected by entropy sampling, random sampling, and the manually curated S1 dataset (1000 samples).

1150	Training Data Size	Entropy Sampled Data	Random Sampled Data	S1 Manually Selected
1151	100	439	418	
1152	200	445	425	_
1153	300	446	424	_
1154	400	447	423	_
1155	500	456	428	_
1156	600	445	429	_
	700	450	419	_
1157	800	443	430	_
1158	900	449	434	_
1159	1000	450	430	452
1160	1100	447	431	_
1161	1200	450	430	_
1162	1300	449	440	_
1163	1400	453	427	_
1164	1500	450	435	_
	1600	450	427	_
1165	1700	458	433	_
1166	1800	454	439	_
1167	1900	451	425	_
1168	2000	450	428	_
1169	3000	450	432	_
1170	4000	453	432	_
1171	5000	448	425	_
1172	10000	455	428	_
	20000	451	437	_
1173	30000	461	447	_
1174	40000	447	438	_
1175	50000	458	441	_