

WINNING THE PRUNING GAMBLE: A UNIFIED APPROACH TO JOINT SAMPLE AND TOKEN PRUNING FOR EFFICIENT SUPERVISED FINE-TUNING

Anonymous authors

Paper under double-blind review

ABSTRACT

As supervised fine-tuning (SFT) evolves from a lightweight post-training step into a compute-intensive phase rivaling mid-training in scale, data efficiency has become critical for aligning large language models (LLMs) under tight budgets. Existing data pruning methods suffer from a fragmented design: they operate either at the sample level or the token level in isolation, failing to jointly optimize both dimensions. This disconnect leads to significant inefficiencies—high-value samples may still contain redundant tokens, while token-level pruning often discards crucial instructional or corrective signals embedded in individual examples. To address this bottleneck, we introduce the *Error–Uncertainty (EU) Plane*, a diagnostic framework that jointly characterizes the heterogeneous utility of training data across samples and tokens. Guided by this insight, we propose *Quadrant-based Tuning (Q-Tuning)*, a unified framework that strategically coordinates sample pruning and token pruning. Q-Tuning employs a two-stage strategy: first, it performs sample-level triage to retain examples rich in informative misconceptions or calibration signals; second, it applies an asymmetric token-pruning policy, using a context-aware scoring mechanism to trim less salient tokens exclusively from misconception samples while preserving calibration samples in their entirety. Our method sets a new state of the art across five diverse benchmarks. Remarkably, on SmolLM2-1.7B, Q-Tuning achieves a +38% average improvement over the full-data SFT baseline using only 12.5% of the original training data. As the first dynamic pruning approach to consistently outperform full-data training, Q-Tuning provides a practical and scalable blueprint for maximizing data utilization in budget-constrained LLM SFT.

1 INTRODUCTION

The explosive growth of alignment datasets—now routinely spanning billions of tokens—has fundamentally transformed Supervised Fine-Tuning (SFT) from a lightweight post-training step into a compute-intensive phase rivaling mid-training in scale (Ouyang et al., 2022; Dong et al., 2024; Yang et al., 2025; Achiam et al., 2023; Team et al., 2023). In this new regime, the primary challenge is no longer simply reducing data volume, but maximizing the utility of every retained token and sample—a task that demands accurate, on-the-fly estimation of data value. Yet despite the emphasis on data efficiency, recent work reveals a troubling paradox: even sophisticated dynamic pruning heuristics often underperform *simple random sampling* (Xia et al., 2024b). This starkly exposes a fundamental disconnect between current strategies and the true utility of alignment data.

At the heart of this challenge lies the fragmented design of existing pruning strategies. *Sample-level pruning* methods (Qin et al., 2024; Zhou et al., 2023a; Wang et al., 2025a; Yang et al., 2024) identify high-potential examples but treat all tokens within them as equally valuable—retaining redundant or even harmful content that dilutes alignment signals. Conversely, *token-level pruning* approaches (Lin et al., 2024; Xia et al., 2025; Chen et al., 2024; Zhang et al., 2024b) apply context-agnostic heuristics uniformly across the dataset, blind to the semantic role of each sample. Such a one-size-fits-all policy fails to differentiate, for instance, between a sample containing a correctable misconception—where only specific tokens need refinement—and one serving as a calibration anchor,

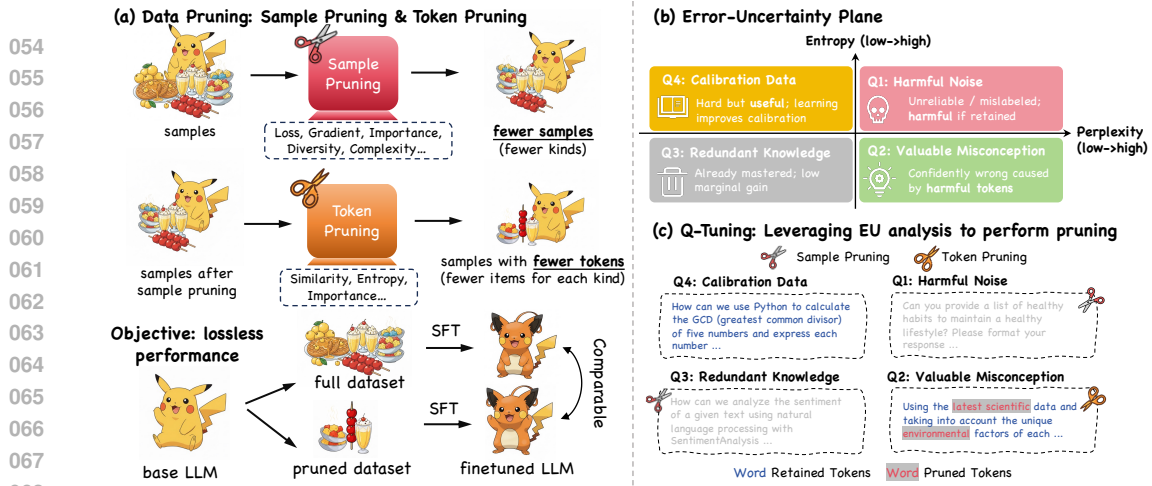


Figure 1: **(a) Data pruning.** *Sample pruning* selects training instances using signals such as loss, gradients, importance, diversity, or complexity; *token pruning* removes low-value tokens within retained samples (e.g., via similarity, entropy, or token importance). **(b) Error-Uncertainty (EU) plane.** We partition samples by perplexity and entropy into four regions: Q1 (harmful noise), Q2 (valuable misconceptions), Q3 (redundant knowledge), and Q4 (calibration data). **(c) Q-Tuning.** Q-Tuning performs joint pruning guided by the EU plane: it drops Q1 and Q3, selectively prunes tokens in Q2, and retains Q4 in full.

which must be preserved holistically to maintain model stability. By operating in isolation, neither paradigm captures the interdependent nature of sample and token utility. This raises a central question: *How can we dynamically coordinate sample selection and token pruning within a unified framework to maximize the true learning utility of limited data?*

To address this, we first formalize the problem as **Generalized Dynamic Data Pruning**, a bilevel optimization framework for jointly optimizing sample and token pruning. We then introduce a novel diagnostic lens, the **Error-Uncertainty (EU) Plane**, which categorizes training instances by mapping model error (perplexity) against model uncertainty (entropy). Specifically, as shown in Figure 1(b), training samples are categorized into four quadrants based on perplexity (higher indicating more wrong) and entropy (higher indicating more uncertain): (i) *Q1 (Harmful Noise)* — unreliable or mislabeled data that actively harms learning and should therefore be removed via sample-level pruning; (ii) *Q2 (Valuable Misconception)* — confidently wrong responses that, when pruned surgically at the token level, can be transformed into powerful teaching signals, making them ideal candidates for token-level pruning; (iii) *Q3 (Redundant Knowledge)* — mastered content offering diminishing returns, best eliminated through sample-level pruning to improve efficiency without sacrificing performance; and (iv) *Q4 (Calibration Data)* — hard but reliable samples essential for improving model confidence and robustness, which should be preserved in full — neither sample nor token pruning should be applied.

This insight motivates our solution: **Quadrant-based Tuning (Q-Tuning)**, the first principled, integrated method for dynamic data pruning. Guided by EU Plane analysis, Q-Tuning implements a two-stage, context-aware strategy, as shown in Figure 1(c). First, at the sample level, it acts as an intelligent triage system: retaining samples that offer clear signals for error correction or calibration (Q2 and Q4), while discarding those classified as harmful noise or redundant knowledge (Q1 and Q3). Crucially, Q-Tuning then applies an asymmetric token policy: for confidently wrong samples (valuable misconceptions), it performs token pruning to isolate the core misconception and amplify the learning signal; for uncertain but correct samples (calibration data), it preserves full token sequences to ensure robust uncertainty modeling. Our contributions are as follows:

1. We formulate the joint sample–token pruning problem as **Generalized Dynamic Data Pruning**, a bilevel optimization framework for hybrid pruning.
2. We introduce the **Error-Uncertainty (EU) Plane**, a tool that reveals and quantifies the heterogeneous value of data across error and uncertainty dimensions. Based on the EU analysis, we propose **Q-Tuning**, the first integrated, diagnosis-driven algorithm for dynamic pruning that coordinates sample and token decisions.
3. Experiments demonstrate that Q-Tuning exceeds full-data training and all existing pruning baselines across 4 models, 5 benchmarks, and 6 different kinds of data budgets. Particularly, with

LLaMA3-8B on GSM8K, Q-Tuning reaches 48.07 using only 35% of the data, outperforming full-data training by 6.0 points and the strongest baseline by 9.9 points.

2 GENERALIZED DYNAMIC DATA PRUNING: A UNIFIED FRAMEWORK

We first propose *Generalized Dynamic Data Pruning*, a framework for accelerating model training where samples and their constituent tokens are selectively and adaptively omitted at each step.

Specifically, the framework first considers the coarse-grained *sample level pruning*, which involves identifying and discarding examples from a mini-batch deemed less informative for the model’s state. Subsequently, *token-level pruning* operates on the retained samples, performing a finer-grained selection to keep a critical subset of tokens within each. This two-stage process is inherently *dynamic*: pruning criteria can be re-evaluated for each new mini-batch, allowing the data distribution used for gradient updates to evolve alongside the model’s parameters θ . The overarching objective is to focus computation on a “doubly-pruned” data subset to maximize training efficiency while preserving or enhancing the model’s generalization.

To formalize this framework, consider a model f_θ with parameters θ . At each training step t , a mini-batch \mathcal{B}_t is drawn from the training distribution \mathcal{D} . The hierarchical pruning process can be modeled by two abstract operators: a sample-level pruner Ψ and a token-level pruner Φ .

Stage 1: Sample-Level Pruning. The operator Ψ splits mini-batch \mathcal{B}_t into kept and discarded samples, governed by keep ratio $r_{\text{sample}} \in [0, 1]$. The retained set, \mathcal{B}'_t , is defined as:

$$\mathcal{B}'_t = \mathcal{B}_t \setminus \Psi(\mathcal{B}_t) \subseteq \mathcal{B}_t, \quad (1)$$

where $\Psi(\mathcal{B}_t)$ represents the set of discarded samples, and the size of the retained set is constrained such that $|\mathcal{B}'_t| = \lfloor r_{\text{sample}} \cdot |\mathcal{B}_t| \rfloor$.

Stage 2: Token-Level Pruning. For each sample in the retained set \mathcal{B}'_t , the operator Φ determines which tokens to keep, guided by a token keep ratio $r_{\text{token}} \in [0, 1]$. This is modeled by generating a binary mask $m(x) \in \{0, 1\}^{L(x)}$ for each sample $x \in \mathcal{B}'_t$, where $L(x)$ is its sequence length. The final, doubly-pruned mini-batch $\tilde{\mathcal{B}}_t$ is constructed by applying these masks:

$$\tilde{\mathcal{B}}_t = \Phi(\mathcal{B}'_t) = \{m(x) \odot x \mid x \in \mathcal{B}'_t\}, \quad (2)$$

where \odot denotes element-wise product, and each mask satisfies $\|m(x)\|_1 = \lfloor r_{\text{token}} L(x) \rfloor$.

Generalized Dynamic Data Pruning. We now put all things together into a unified framework. The central problem of this framework is to identify the optimal dynamic pruning operators, (Φ, Ψ) , that guide the training process to a final model θ^* with the best possible generalization performance. These operators are applied at each step t to transform a mini-batch \mathcal{B}_t into a computationally cheaper, pruned version $\tilde{\mathcal{B}}_t$, while adhering to predefined keep ratios.

This problem is naturally captured as a *bi-level optimization problem*. The outer loop seeks optimal pruners, while the inner loop represents the iterative training procedure that produces the final model parameters under the guidance of these pruners. Formally, the objective is as follows:

$$\begin{aligned} \min_{\Phi, \Psi} \quad & \mathbb{E}_{(x,y) \sim \mathcal{D}_{\text{test}}} [\mathcal{L}_{\text{test}}(f_{\theta^*}(x), y)] \\ \text{s.t.} \quad & \theta^* = \arg \min_{\theta} \sum_{t=1}^T \mathbb{E}_{\mathcal{B}_t \sim \mathcal{D}} \mathbb{E}_{(x,y) \sim \tilde{\mathcal{B}}_t} [\mathcal{L}_{\text{train}}(f_{\theta_t}(x), y)] \end{aligned} \quad (3)$$

where at each step t : $\tilde{\mathcal{B}}_t = \Phi(\mathcal{B}_t \setminus \Psi(\mathcal{B}_t))$, with keep ratios $(r_{\text{sample}}, r_{\text{token}})$.

Specifically, the *outer objective* defines the quality metric for any pair of pruners (Φ, Ψ) : the final test performance of the model they produce. The *inner objective* defines the training process itself, where the final parameters θ^* result from cumulatively minimizing the loss over a sequence of dynamically pruned mini-batches. All existing methods that instantiates these operators Φ and Ψ can be seen as a specific solution to this alignment problem.

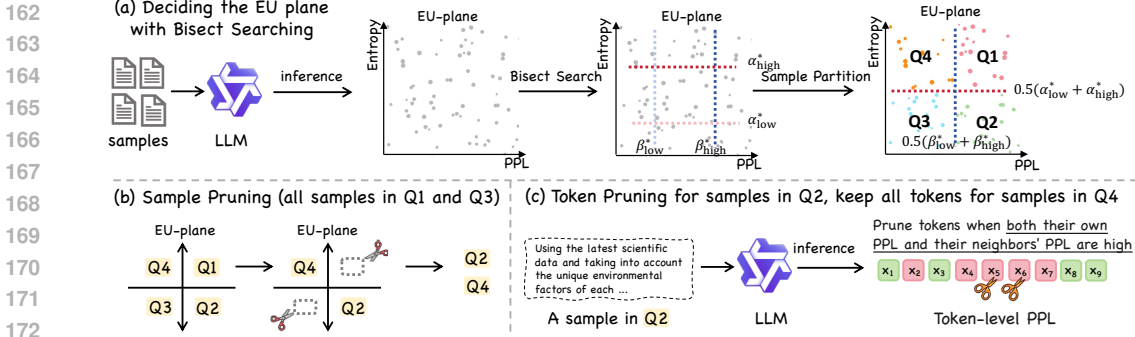


Figure 2: **(a) Constructing the Error-Uncertainty (EU) Plane via Bisect Search.** We run the base LLM to compute sample-level perplexity (PPL) and entropy, and use bisect search to set thresholds $(\alpha_{low}^*, \alpha_{high}^*, \beta_{low}^*, \beta_{high}^*)$ that partition the EU plane into Q1–Q4. **(b) Sample Pruning for Q1 and Q3.** Samples in Q1 and Q3 are pruned at the sample level, while those in Q2 and Q4 are retained. **(c) Token-Level Pruning for Q2 Samples.** For retained Q2 samples, token-level pruning is applied based on both the token’s own perplexity and the average perplexity of its neighboring tokens. Tokens with high local PPL are removed, preserving only the most informative ones.

3 WINNING THE PRUNING GAMBLE

Building on the previous analysis, we now introduce **Q-Tuning**, a dynamic pruning method guided by the Error–Uncertainty (EU) Plane. As shown in Figure 1, Q-Tuning proceeds in two coordinated stages: it first prunes harmful noise and redundant knowledge at the sample level, and then applies targeted token-level pruning to valuable misconceptions while preserving calibration data in full. This integrated strategy avoids the failure modes of one-dimensional heuristics and transforms pruning from a gamble into a principled, performance-enhancing process.

3.1 THE ERROR-UNCERTAINTY PLANE: A DIAGNOSTIC LENS

We now formalize the **Error–Uncertainty (EU) Plane** introduced in Figure 1(b). Each training sample is positioned by two orthogonal axes: *error*, quantified by perplexity (PPL), and *uncertainty*, quantified by predictive entropy (Ent).

Measure Error via Perplexity. For a training sample (x, y) with sequence length $L(x)$, let $T(x) \subseteq \{1, \dots, L(x)\}$ denote the set of trainable token positions. We define sample-level error as the perplexity (PPL), derived from the average token-level negative log-likelihood:

$$\text{PPL}(x, y; f_\theta) = \exp \left(\sum_{i \in T(x)} \frac{-\log p(y_i | x, y_{<i})}{|T(x)|} \right) \quad (4)$$

High perplexity indicates the model finds the ground-truth continuation highly surprising, which is a sign of either genuine difficulty or misconception.

Measure Uncertainty via Predictive Entropy. Independent of correctness, we quantify the model’s uncertainty using entropy (Ent), computed as the average token-level entropy over trainable positions, where $v \in \mathcal{V}$ denotes a vocabulary token:

$$\text{Ent}(x, y; f_\theta) = \frac{\sum_{i \in T(x)} H(p_\theta(\cdot | x, y_{<i}))}{|T(x)|}. \quad (5)$$

High entropy indicates that the model spreads probability mass broadly, reflecting persistent uncertainty even when the top prediction is correct. Taken together, PPL and Ent map each sample onto the EU Plane, providing a principled and interpretable basis for partitioning data into the four quadrants illustrated in Figure 1(b).

3.2 Q-TUNING: GENERALIZED DYNAMIC DATA PRUNING

Q-Tuning operationalizes the insights of the EU Plane into a two-stage, per-batch dynamic pruning strategy, as illustrated in Figure 2. In the first stage, bisect search determines quantile-based thresholds that partition samples into four quadrants, enabling sample-level pruning of uninformative data. In the second stage, token-level pruning is selectively applied within retained misconception samples to isolate informative signals, while calibration samples are preserved in full. The pseudocode of our method is illustrated in Algorithm 1 in Appendix D.

Stage 1: Constructing the EU Plane via Binary Search and Pruning Samples Accordingly. At each training step, we compute the perplexity (PPL) and entropy (Ent) for every sample $x \in \mathcal{B}_t$ using a gradient-free forward pass of the current model f_{θ_t} . These statistics map each sample to a point on the EU Plane (Figure 2(a)).

Our objective is to determine quantile thresholds (α^*, β^*) on the PPL and Entropy axes such that the retained proportion of samples in $Q_2 \cup Q_4$ exactly matches the target sample retention ratio r_{sample} . Here $\text{Quantile}_\gamma(X)$ denotes the γ -quantile of a variable X , *i.e.*, the smallest value q such that at least a fraction γ of samples satisfy $X \leq q$.

To locate (α^*, β^*) , we perform a **bisect search** on both axes. For each axis, the search interval is initialized as $[0, 0.5]$. At each iteration we set $\alpha = \frac{1}{2}(\alpha_{\text{low}} + \alpha_{\text{high}})$, $\beta = \frac{1}{2}(\beta_{\text{low}} + \beta_{\text{high}})$, and compute the proportion of samples that would fall into the tentative $Q_2 \cup Q_4$ region defined by the current thresholds (α, β) . The intervals are then updated as

$$(\alpha_{\text{low}}, \alpha_{\text{high}}, \beta_{\text{low}}, \beta_{\text{high}}) = \begin{cases} (\alpha_{\text{low}}, \alpha, \beta_{\text{low}}, \beta), & \text{if } \gamma > r_{\text{sample}}, \\ (\alpha, \alpha_{\text{high}}, \beta, \beta_{\text{high}}), & \text{otherwise.} \end{cases} \quad (6)$$

The search converges in $O(\log(1/\epsilon))$ iterations (typically < 10), incurring negligible overhead. The resulting thresholds (α, β) partition the EU Plane (Figure 2(b)), from which Q2 and Q4 samples are retained and Q1 and Q3 discarded.

Stage 2: Token-Level Pruning for Confident Errors. While Stage 1 removes entire samples deemed unhelpful, not all retained examples are equally homogeneous inside. In particular, samples in Q2 (Valuable Misconceptions) often contain a mix of informative context and locally harmful tokens that mislead the model. To extract the useful signal, we apply token-level pruning that discards only the most detrimental tokens while preserving the surrounding context. In contrast, samples in Q4 (Calibration Data) are challenging yet reliable, and every token contributes to improving model calibration. Therefore, Q4 sequences are kept intact without any token pruning.

For a retained sample $x \in Q_2$ with target sequence y , let $T(x)$ denote the set of trainable token positions. For each token $i \in T(x)$ (Figure 2(c)), we compute its token-level perplexity PPL_i , which measures how surprising the ground-truth token y_i is to the model.

To avoid pruning rare but meaningful tokens based on isolated spikes, we compute a **smoothed importance score** that incorporates local context:

$$s_i(x, y; f_\theta) = (1 - \lambda) \text{PPL}_i(x, y; f_\theta) + \lambda[\text{PPL}_{i-1}(x, y; f_\theta) + \text{PPL}_{i+1}(x, y; f_\theta)], \quad (7)$$

where $\lambda \in [0, 1]$ balances the contribution of neighboring tokens (default $\lambda = 0.5$). This smoothing reduces the risk of mistakenly removing single high-PPL tokens that may still be semantically critical, and we analyze the sensitivity of this choice in the ablation study (see Figure 5 in Section 4.2).

A token is deemed detrimental if both its own PPL_i and the average PPL of its immediate neighbors exceed a percentile-based threshold (*e.g.*, the median). All tokens in $T(x)$ are subsequently ranked by their smoothed scores s_i , and only the top- r_{token} fraction are retained to construct a binary mask $m(x)$. This mask selectively removes locally noisy tokens while preserving the informative context essential for learning. By contrast, no token-level pruning is applied to Q4 samples, as each token therein contributes to reliable calibration and must be preserved in full.

Table 1: Evaluation on WizardLM with a sample ratio of 12.5% and a token ratio of 50%, where \uparrow and \downarrow respectively denote improvements or degradations over the *Random-Random* baseline. Additional results under more sample and token pruning ratios are provided in the Table 8 and Table 9 in Appendix C.1 and results on Qwen3-8B trained on the OpenHermes are reported in Appendix C.2.

Sample Pruner	Token Pruner	ARC-E	ARC-C	LLaMA2-7B				Mistral-7B					
				GSM8K	SQuAD	TriviaQA	Avg.	ARC-E	ARC-C	GSM8K	SQuAD	TriviaQA	Avg.
Zero-Shot		53.44	38.98	5.31	12.18	43.00	30.58	66.67	46.10	18.35	10.01	43.77	36.98
12.5% Samples, 50% Tokens													
Random	Random	59.25	41.02	8.11	12.75	48.75	33.98	70.55	48.14	22.74	19.57	52.63	42.73
	PPL	60.49 ^{71.24}	43.39 ^{72.37}	7.20 ^{10.91}	12.20 ^{10.55}	48.04 ^{40.71}	34.26 ^{70.28}	70.72 ^{70.17}	48.47 ^{70.33}	25.78 ^{73.04}	21.36 ^{71.79}	53.92 ^{71.29}	44.05 ^{71.32}
	FastV	59.96 ^{70.71}	42.37 ^{71.35}	5.76 ^{12.35}	11.31 ^{11.44}	46.42 ^{42.33}	33.17 ^{40.81}	70.72 ^{70.17}	46.44 ^{41.70}	18.80 ^{43.94}	19.14 ^{40.43}	51.56 ^{41.07}	41.33 ^{41.40}
	SparseVLM	54.32 ^{44.93}	37.97 ^{43.05}	7.35 ^{10.76}	12.76 ^{70.01}	44.65 ^{44.10}	31.41 ^{42.57}	67.02 ^{43.53}	44.75 ^{43.39}	20.24 ^{42.50}	10.97 ^{48.60}	44.61 ^{48.02}	37.52 ^{45.21}
Longest	Random	59.96 ^{70.71}	44.41^{73.39}	7.51 ^{10.60}	15.34 ^{72.59}	48.91 ^{70.16}	35.22 ^{71.24}	74.25 ^{73.70}	48.81 ^{70.67}	28.73 ^{75.99}	17.66 ^{41.91}	55.73 ^{73.10}	45.04 ^{72.31}
	PPL	61.19 ^{71.84}	43.73 ^{72.71}	6.82 ^{11.29}	16.33 ^{73.58}	48.16 ^{40.59}	35.24 ^{71.26}	75.49^{74.94}	50.17 ^{72.03}	27.98 ^{73.24}	21.49 ^{71.92}	56.55 ^{73.92}	46.33 ^{73.60}
	FastV	59.25 ^{70.60}	43.05 ^{72.03}	5.69 ^{12.42}	13.64 ^{70.89}	46.98 ^{41.77}	33.72 ^{40.26}	74.43 ^{73.88}	49.15 ^{71.01}	25.70 ^{72.96}	22.89 ^{73.32}	54.15 ^{71.52}	45.26 ^{72.53}
	SparseVLM	54.32 ^{44.93}	38.31 ^{42.71}	7.13 ^{10.98}	10.92 ^{11.83}	43.77 ^{44.98}	30.89 ^{43.09}	69.49 ^{41.06}	46.10 ^{42.04}	28.89 ^{76.15}	8.62 ^{41.95}	50.30 ^{42.33}	40.68 ^{42.05}
InfoBatch	Random	60.31 ^{71.06}	41.36 ^{70.34}	5.38 ^{12.73}	15.71 ^{72.96}	47.74 ^{41.01}	34.10 ^{70.12}	69.31 ^{41.24}	45.76 ^{42.38}	18.95 ^{43.79}	21.23 ^{71.66}	50.39 ^{42.24}	41.13 ^{41.60}
	PPL	59.43 ^{70.18}	40.34 ^{40.68}	5.91 ^{12.20}	13.18 ^{70.43}	48.31 ^{40.44}	33.44 ^{40.54}	70.72 ^{70.17}	47.12 ^{41.02}	18.12 ^{44.62}	24.10 ^{74.53}	51.26 ^{41.37}	42.26 ^{40.47}
	FastV	58.90 ^{40.35}	43.39 ^{72.37}	3.34 ^{14.77}	12.37 ^{40.38}	46.88 ^{41.87}	32.98 ^{41.00}	69.14 ^{41.41}	45.42 ^{42.72}	14.86 ^{47.88}	23.19 ^{73.62}	50.58 ^{42.05}	40.64 ^{42.09}
	SparseVLM	54.67 ^{44.58}	40.00 ^{41.02}	7.73 ^{10.38}	12.41 ^{40.34}	45.07 ^{43.68}	31.98 ^{42.00}	68.25 ^{42.30}	45.08 ^{43.06}	23.43 ^{70.69}	10.17 ^{49.40}	45.34 ^{47.29}	38.46 ^{44.27}
Entropy	Random	60.31 ^{71.06}	42.37 ^{71.35}	6.44 ^{11.67}	14.10 ^{71.35}	48.09 ^{40.66}	34.27 ^{70.29}	72.13 ^{71.58}	48.81 ^{70.67}	20.09 ^{42.65}	17.55 ^{42.02}	54.69 ^{72.06}	42.66 ^{40.07}
	PPL	60.49 ^{71.24}	43.73 ^{72.71}	6.90 ^{11.21}	14.53 ^{71.78}	48.76 ^{70.01}	34.88 ^{70.90}	72.84 ^{72.29}	47.80 ^{40.34}	24.18 ^{71.44}	22.80 ^{73.23}	54.69 ^{72.06}	44.46 ^{71.73}
	FastV	58.91 ^{40.34}	43.05 ^{72.03}	6.37 ^{11.74}	13.03 ^{70.28}	47.05 ^{41.70}	33.68 ^{40.30}	73.90 ^{73.35}	47.12 ^{41.02}	24.56 ^{71.82}	23.96 ^{74.39}	54.67 ^{72.04}	44.84 ^{72.11}
	SparseVLM	55.20 ^{44.05}	38.98 ^{42.04}	7.51 ^{10.60}	12.65 ^{40.10}	46.14 ^{42.61}	32.10 ^{41.88}	68.08 ^{42.47}	44.07 ^{44.07}	24.87 ^{72.13}	10.72 ^{48.85}	47.00 ^{45.63}	38.95 ^{43.78}
Q-Tuning (Ours)		64.20^{74.95}	42.03 ^{71.01}	10.54^{72.43}	18.79^{76.04}	53.12^{74.37}	37.74^{73.76}	71.60 ^{71.05}	48.14 ^{70.00}	29.34^{66.60}	27.75^{78.18}	57.78^{78.15}	46.92^{74.19}
Full Dataset		61.55	42.37	8.64	13.80	50.45	35.36	71.25	45.76	26.68	31.81	53.67	45.84

4 EXPERIMENTS

We conducted experiments on language models from multiple families and scales, including LLaMA-series (LLaMA2-7B (Touvron et al., 2023), LLaMA3-8B (Grattafiori et al., 2024)), Qwen3-series (Qwen3-8B (Yang et al., 2025), Qwen3-14B, Qwen3-32B), Mistral-7B (Jiang et al., 2023a), and SmoLLM2-1.7B (Allal et al., 2025)). We fine-tuned models for two settings: (i) alignment on OpenHermes (Teknium, 2023) and WizardLM (Xu et al., 2024), and (ii) reasoning on MathInstruct (Yue et al., 2023). For evaluation, we used ARC-E, ARC-C (Clark et al., 2018), GSM8K (Cobbe et al., 2021), SQuAD (Rajpurkar et al., 2016), and TriviaQA (Joshi et al., 2017) for alignment task, and GSM8K and MATH (Hendrycks et al., 2021) for reasoning task. We compared against baselines that pair sample-level pruning methods (Random, Longest, Entropy, InfoBatch (Qin et al., 2024), Alpagasus (Chen et al., 2023), Deita (Liu et al., 2023), DS2 (Pang et al., 2024), and LESS (Xia et al., 2024a)) with token-level pruning methods (Random, PPL, FastV (Chen et al., 2024), SparseVLM (Zhang et al., 2024b), Rho-1 (Lin et al., 2024), and TokenCleaning (Pang et al., 2025)). More detailed experimental settings are in Appendix B.

4.1 MAIN RESULTS

Results on alignment datasets. Table 1 summarizes results in our alignment setting under a fixed low-budget regime (12.5% samples, 50% tokens) on WizardLM. Across both LLaMA2-7B and Mistral-7B, Q-Tuning achieves the best overall averages (37.74 and 46.92), outperforming all evaluated combinations of existing sample-pruning and token-pruning methods, and it also exceeds full-data fine-tuning (35.36 and 45.84). Many baseline pairings improve one or two benchmarks but exhibit noticeable trade-offs on others (e.g., degradation on GSM8K or SQuAD for several token-pruning choices), indicating that naively combining independently designed sample- and token-level pruners is insufficient under a tight compute budget. In contrast, Q-Tuning delivers consistent gains across all benchmarks, suggesting that jointly coordinating sample triage with targeted token reduction better preserves high-utility supervision in this alignment evaluation. We report additional results under more pruning ratios in Appendix C.1 (Table 8 and Table 9), showing similar trends, while Appendix C.2 presents detailed results and analysis on Qwen3-8B with larger performance gains of up to 8.23 points and improvements of up to 5.55 over the strongest baseline. Overall, these results highlight that Q-Tuning is robust to heterogeneous task mixes and remains effective even when both sample and token budgets are simultaneously constrained across diverse model families.

Results on reasoning dataset. Table 2 compares Q-Tuning against diverse baselines on math reasoning datasets under constrained sample and token budgets. Several methods, especially FastV-based token pruning, degrade sharply (e.g., driving LLaMA3-8B average accuracy below 6% in

Table 2: Evaluation of pruning strategies on GSM8K and MATH under 25% samples with 50% tokens. \uparrow and \downarrow respectively indicate improvements or degradations over the *Random-Random* baseline under the same sample and token keep ratio. Additional results under more sample and token pruning ratios are provided in the Table 12 in Appendix C.3.

Sample Pruner	Token Pruner	LLaMA3-8B			Mistral-7B			SmolLM2-1.7B		
		GSM8K	MATH	Avg.	GSM8K	MATH	Avg.	GSM8K	MATH	Avg.
Zero-Shot		27.82	2.26	15.04	19.86	3.30	11.58	15.47	2.20	8.83
25% Samples, 50% Tokens										
Random	Random	23.96	2.56	13.26	23.35	1.54	12.45	14.33	2.56	8.44
	PPL	24.18 ^{70.22}	2.58 ^{70.02}	13.38 ^{70.12}	24.94 ^{71.59}	2.02 ^{70.48}	13.48 ^{71.03}	14.18 ^{40.15}	2.08 ^{40.48}	8.13 ^{40.31}
	FastV	12.13 ^{411.83}	2.32 ^{40.24}	7.23 ^{46.03}	12.36 ^{410.99}	1.24 ^{40.30}	6.80 ^{45.65}	9.86 ^{44.47}	1.92 ^{40.64}	5.89 ^{42.55}
	SparseVLM	22.97 ^{40.99}	4.72 ^{72.16}	13.85 ^{70.59}	19.26 ^{44.09}	4.58 ^{73.04}	11.92 ^{40.53}	13.19 ^{41.14}	3.48 ^{70.92}	8.34 ^{40.10}
Longest	Random	22.14 ^{41.82}	3.18 ^{70.62}	12.66 ^{40.60}	21.91 ^{41.44}	2.18 ^{70.64}	12.05 ^{40.40}	12.89 ^{41.44}	2.06 ^{40.50}	7.47 ^{40.97}
	PPL	24.94 ^{70.98}	2.78 ^{70.22}	13.86 ^{70.60}	22.90 ^{40.45}	1.86 ^{70.32}	12.38 ^{40.07}	13.19 ^{41.14}	1.78 ^{40.78}	7.49 ^{40.95}
	FastV	9.48 ^{414.48}	2.26 ^{40.30}	5.87 ^{47.39}	7.13 ^{416.22}	1.46 ^{40.08}	4.29 ^{48.16}	12.36 ^{41.97}	1.82 ^{40.74}	7.09 ^{41.35}
	SparseVLM	26.91 ^{72.95}	4.68 ^{72.12}	15.80 ^{72.54}	24.34 ^{70.99}	4.84 ^{73.30}	14.59 ^{72.14}	12.43 ^{41.90}	3.60 ^{71.04}	8.02 ^{40.42}
InfoBatch	Random	26.23 ^{72.27}	2.42 ^{40.14}	14.33 ^{71.07}	27.14 ^{73.79}	2.24 ^{70.70}	14.69 ^{72.24}	14.33 ^{70.00}	1.66 ^{40.90}	7.99 ^{40.45}
	PPL	26.91 ^{72.95}	2.66 ^{70.10}	14.79 ^{71.53}	27.90 ^{74.55}	2.52 ^{70.98}	15.21 ^{72.76}	14.71 ^{70.38}	1.90 ^{40.66}	8.30 ^{40.14}
	FastV	7.58 ^{416.38}	1.88 ^{40.68}	4.73 ^{48.53}	6.44 ^{416.91}	1.34 ^{40.20}	3.89 ^{48.56}	7.51 ^{46.82}	1.62 ^{40.94}	4.56 ^{43.88}
	SparseVLM	14.63 ^{49.33}	3.26 ^{70.70}	8.95 ^{44.31}	11.90 ^{414.45}	1.94 ^{70.40}	6.92 ^{45.53}	11.90 ^{42.43}	4.36 ^{71.80}	8.13 ^{40.31}
Entropy	Random	30.02 ^{76.06}	3.66 ^{71.10}	16.84 ^{73.58}	26.61 ^{73.26}	2.08 ^{70.54}	14.35 ^{71.90}	14.18 ^{40.15}	2.24 ^{40.32}	8.21 ^{40.23}
	PPL	32.98 ^{79.02}	2.92 ^{70.36}	17.95 ^{74.69}	30.17 ^{76.82}	1.76 ^{70.22}	15.97 ^{73.52}	16.38 ^{72.05}	2.40 ^{40.16}	9.39 ^{70.95}
	FastV	17.29 ^{46.67}	2.66 ^{70.10}	9.97 ^{43.29}	14.56 ^{48.79}	1.12 ^{40.42}	7.84 ^{44.61}	12.59 ^{41.74}	2.36 ^{40.20}	7.47 ^{40.97}
	SparseVLM	20.85 ^{43.11}	5.12 ^{72.56}	12.98 ^{40.28}	19.56 ^{43.79}	4.20 ^{72.66}	11.88 ^{40.57}	14.18 ^{40.15}	2.90 ^{70.34}	8.54 ^{70.10}
Q-Tuning (Ours)		36.32 ^{712.36}	5.54 ^{72.98}	20.93 ^{71.67}	41.47 ^{718.12}	4.0 ^{72.46}	22.74 ^{710.29}	21.83 ^{77.50}	3.90 ^{71.34}	12.87 ^{74.43}
Full Dataset		32.90	3.02	17.96	42.08	3.08	22.58	16.53	2.10	9.31

multiple settings). In contrast, quality-aware sample selection (Entropy, InfoBatch) consistently improves over random sampling: Entropy–Random raises LLaMA3-8B from 13.26 to 16.84 and yields similar gains on Mistral-7B. Across all model scales, Q-Tuning is the strongest and most stable, improving the Random–Random baseline by 7.67 points on LLaMA3-8B and 10.29 on Mistral-7B, reaching 20.93 and 22.74, respectively. Notably, Q-Tuning exceeds full-data GSM8K performance for both models despite using only a quarter of the samples and half of the tokens, indicating strong data efficiency and better training-resource utilization; larger models benefit more from the joint strategy. Additional ratios are reported in Appendix C.3 (Table 12). This suggests that jointly optimized pruning can improve reasoning training signal quality, rather than merely reducing compute. In particular, it preferentially retains error-revealing and correction-rich spans while trimming redundant verbiage. As a result, the model receives cleaner gradient signals and learns more effectively.

Comparison with stronger baselines.

Figure 3 compares Q-Tuning with stronger, task-relevant baselines under low-budget settings. In the top figures with 12.5% samples and 50% tokens, Q-Tuning achieves the best average on both models, scoring 37.74 on LLaMA2-7B and 46.92 on Mistral-7B, which exceeds the strongest sample-selection baseline by +3.44 and +2.00, respectively. In the bottom figures with a fixed 50% token budget, Q-Tuning also outperforms token-level baselines at the same 12.5% sample ratio, improving over the best baseline from 33.82 to 37.74 on LLaMA2-7B (+3.92) and from 43.57 to 46.92 on Mistral-7B (+3.35). These results indicate that Q-Tuning remains competitive against stronger baselines under the same constrained budgets and delivers consistent gains across model families. Additional results under more pruning ratios are reported in Appendix C.4.

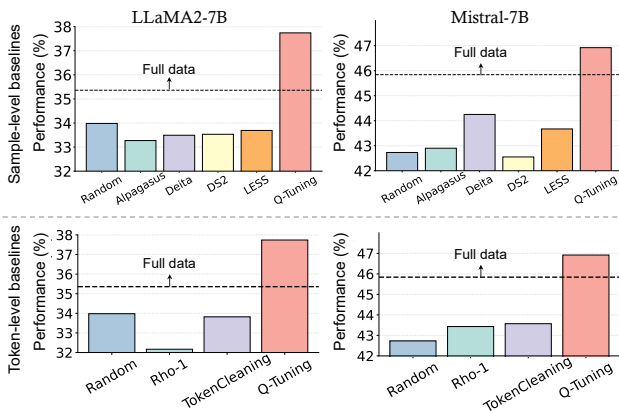


Figure 3: Comparison with stronger task-relevant baselines under matched low-budget settings. Top: Sample-pruner baselines with a sample ratio of 12.5% and a token ratio of 50%. Bottom: Token-pruner baselines with a sample ratio of 12.5% and a token ratio of 50%; the dashed line marks the full-data upper bound. Additional results under more pruning ratios are reported in Appendix C.4.

Robustness to extreme data budgets. Table 3 shows that Q-Tuning remains highly data-efficient even when retaining only 3–9% of the original training tokens on Mistral-7B. Across all four extreme settings, the average score ranges from 46.23 to 46.92, consistently outperforming both the zero-shot baseline (36.98) and the full-dataset SFT baseline (45.84), despite using at most 8.75% of the data and maintaining strong per-task accuracy. Notably, with just 6.25% samples and 50% tokens (*i.e.*, 3.125% data retention), Q-Tuning already achieves an average of 46.23, exceeding full-data fine-tuning. Increasing the token ratio to 70% at the same sample budget further improves the average to 46.62, while still operating in a regime where less than 5% of the original data is used in training epochs.

Table 3: Q-Tuning under extreme data-retention settings on the Wizard dataset using Mistral-7B.

Sample Ratio	Token Ratio	Data Retention	Mistral-7B					Avg.
			ARC-E	ARC-C	GSM8K	SQuAD	TriviaQA	
Zero-Shot		-	66.67	46.10	18.35	10.01	43.77	36.98
6.25%	50%	3.125%	70.72	46.10	40.41	17.25	56.68	46.23
6.25%	70%	4.375%	71.60	48.47	39.73	16.13	57.16	46.62
12.50%	50%	6.250%	71.60	48.14	29.34	27.75	57.78	46.92
12.50%	70%	8.750%	71.78	47.12	26.08	32.79	56.17	46.79
Full Dataset		100%	71.25	45.76	26.68	31.81	53.67	45.84

Generalization across model scales. Figure 4 studies Q-Tuning on Qwen3-series of increasing scale (8B, 14B, and 32B). Under the same constrained budgets (12.5% samples with 50% or 70% tokens), Q-Tuning consistently matches or exceeds full-dataset SFT across benchmarks, demonstrating strong scalability despite using only a small fraction of the training data. The gains are especially evident on reasoning- and QA-oriented tasks such as GSM8K, ARC-E, and SQuAD. Moreover, for the 32B model, increasing the token budget from 50% to 70% typically yields further improvements and delivers the strongest overall averages.

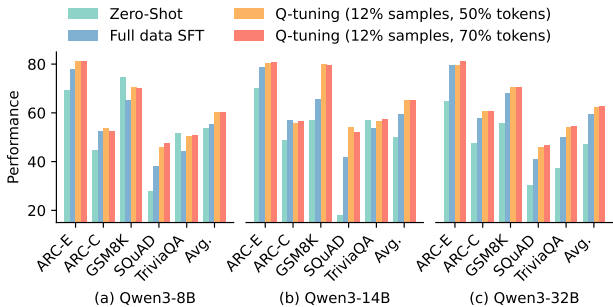


Figure 4: Comparison of Qwen3-series of varying scales (8B, 14B, 32B) across multiple benchmarks and their average. We report performance of Zero-shot, full dataset, and Q-Tuning with 12.5% samples under two token ratios (50% and 70%).

4.2 FURTHER ANALYSIS

Sensitivity to hyperparameters. An important question is how batch size and neighbor awareness influence the performance of Q-Tuning. We first varied the batch size across 8, 16, 32 and evaluated the method on GSM8K, SQuAD, and TriviaQA. As shown in the upper part of Figure 5, larger batch sizes tend to yield better performance, with more noticeable gains on GSM8K and SQuAD, while TriviaQA remains relatively stable, indicating limited sensitivity to batch size. We further investigated the effect of neighbor awareness controlled by the coefficient λ . As illustrated in the lower part of Figure 5, moderate values of λ lead to improvements on GSM8K and SQuAD, whereas extreme values result in diminishing or unstable gains; in contrast, performance on TriviaQA is largely unaffected across different λ settings. These observations indicate that Q-Tuning benefits from appropriate batch sizes and a moderate degree of neighbor awareness, though the sensitivity varies across tasks. Additional benchmark results are provided in the Figure 9 in Appendix C.5.1.

Effect of token pruning strategy. We compare different token pruning strategies under a fixed 25% × 50% budget on Mistral-7B, including Rho-1, PPL, and a reversed PPL variant, as shown in Table 4. The neighbor-aware PPL used in Q-Tuning achieves the best average performance, while alternative pruners are less consistent across tasks. This supports our design choice of locality-aware scoring that preserves salient error-relevant spans.

Table 4: Ablation study on token pruner in Q-tuning

Method	Avg.	ARC-E	ARC-C	GSM8K	SQuAD	TriviaQA
Rho-1	45.03	69.66	46.78	24.03	29.23	55.43
PPL ($\lambda = 0$)	45.92	71.60	46.44	25.32	29.71	56.54
Reversed PPL ($\lambda = 0.5$)	44.86	73.02	47.12	16.68	32.01	55.47
PPL ($\lambda = 0.5$)	46.79	71.78	47.12	26.08	32.79	56.17

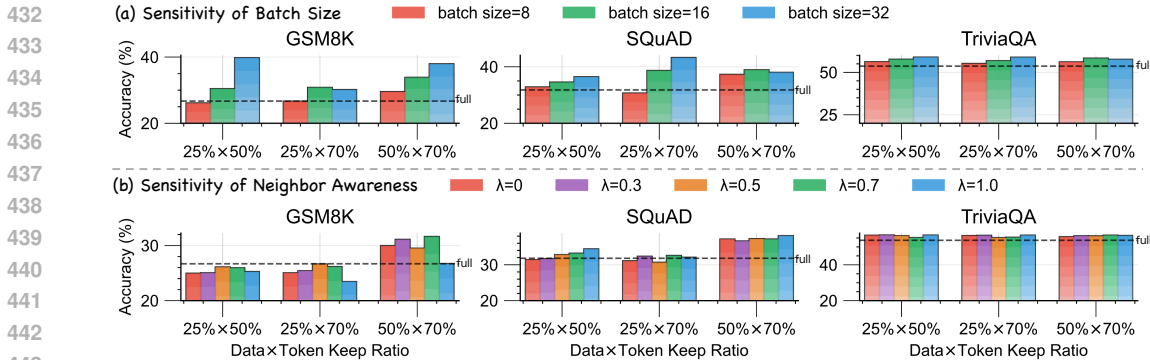


Figure 5: Effect of varying (a) batch size and (b) neighbor awareness for Mistral-7B under three keep ratio configurations. Additional benchmark results (Avg. of five benchmarks, ARC-E, ARC-C) are provided in the Figure 9 in Appendix C.5.1.

Effect of quadrant-wise pruning strategies. Table 5 reports a quadrant-wise ablation on Qwen3-8B under two budget settings. Retaining Q1 consistently yields the lowest averages, with pronounced drops on reasoning and QA tasks (e.g., GSM8K and SQuAD), whereas pruning Q1 leads to more stable performance. Token-level pruning of Q4 is generally harmful, indicating that these low-PPL yet diverse tokens carry useful signal and should be preserved. In contrast, pruning Q2 tokens is crucial, as retaining Q2’s high-PPL tokens degrades performance while removing them consistently improves results. Therefore, the best configurations prune Q1 and Q3 at the sample level and prune Q2 at the token level, supporting the design of Q-Tuning. Additional quadrant-wise ablations are provided in Appendix C.5.2.

Table 5: Ablation on four quadrants on OpenHerms and Qwen3-8B under different sample and token ratios, where \checkmark indicates that the corresponding quadrant is pruned. Additional quadrant-wise ablations are provided in Appendix C.5.2.

Sample pruning	Token pruning		Qwen3-8B					Avg.
	Q3	Q2 Q4	ARC-E	ARC-C	GSM8K	SQuAD	TriviaQA	
12.5% Samples, 50% Tokens								
	\checkmark	\checkmark	79.01	52.20	74.68	36.35	48.20	58.09
\checkmark	\checkmark	\checkmark \checkmark	79.37	52.20	71.80	35.55	47.23	57.23
\checkmark	\checkmark	\checkmark	78.84	53.90	68.99	44.32	49.15	59.04
\checkmark	\checkmark	\checkmark	80.95	53.56	70.51	45.76	50.16	60.19
25% Samples, 50% Tokens								
	\checkmark	\checkmark	75.31	50.17	68.08	29.47	45.90	53.79
\checkmark	\checkmark	\checkmark \checkmark	76.37	50.17	65.28	31.63	45.24	53.74
\checkmark	\checkmark	\checkmark	79.89	52.88	69.37	43.54	48.33	58.80

Case Study: Which Types of Tokens Are Preferentially Pruned in Q2? To better characterize the semantic effects of Q2 pruning, we perform a microscopic analysis within the Q2 quadrant. We randomly sample 300 Q2 instances and use GPT-5 (Singh et al., 2025) to annotate functional roles for 44K+ tokens, then compare category proportions between pruned and kept sets. Figure 6 shows that redundant content is over-represented among pruned tokens: boilerplate/templates and formatting increase, and non-answer content shifts toward pruning. In contrast, supervision-critical signals are enriched in kept tokens. Gold answer spans account for a higher proportion among retained tokens, and numbers/symbols are substantially more prevalent. We provide quadrant examples in Appendix E.1 and token-level pruning cases in Appendix E.2.

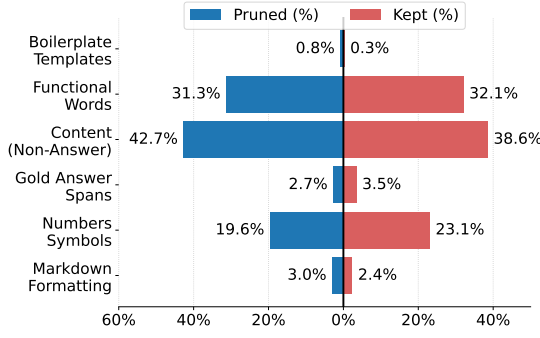


Figure 6: Semantic distribution of pruned vs. retained tokens in the Q2 quadrant, based on 300 instances annotated by GPT-5.

4.3 EFFICIENCY ANALYSIS: EMPIRICAL LATENCY

Table 6 quantifies the end-to-end efficiency of Q-Tuning on LLaMA2-7B under our standard three-epoch setup with 8 x A100 GPUs. Compared to full-data SFT (90 minutes, 1.0 x FLOPs, 78% peak memory), Q-Tuning with 12.5% samples and 50% tokens reduces training compute to 0.0625 x and completes in 65 minutes, yielding a ~28% wall-clock reduction while lowering peak memory to 27%.

The screening overhead introduced by dynamic sample pruning is modest relative to the overall training budget and is effectively amortized by the computational savings gained from skipping backward passes and optimizer updates on pruned data. Even as budget increases (25% and 50% samples at 50% tokens), total time remains below full-data SFT (70–80 minutes) with substantially reduced FLOPs (0.125–0.25 \times) and peak memory (35–50%), indicating a favorable practical cost–benefit trade-off.

Table 6: Efficiency and overhead breakdown on LLaMA2-7B under different pruning settings.

Pruner		Selection Time (mins)		Training Efficiency	
Sample Ratio (%)	Token Ratio (%)	Sample Selection	Token Selection	Total Time (mins)	Peak Mem (%)
12.5	50	35.88	6.96	65	27
25	50	38.64	7.49	70	35
50	50	44.16	8.56	80	50
100 (Full)	100	0.00	0.00	90	78

5 DISCUSSION

Can Q-Tuning outperform independent sample or token pruning? To isolate the effect of each pruning dimension, we conduct controlled ablations that apply *only* dynamic sample pruning (retaining all tokens) or *only* dynamic token pruning (retaining all samples). As shown in Figure 7, Q-Tuning consistently outperforms all baseline methods across keep ratios on both LLaMA2-7B and Mistral-7B, indicating that its advantage does not rely on a particular pruning axis but stems from more effective data utility modeling.

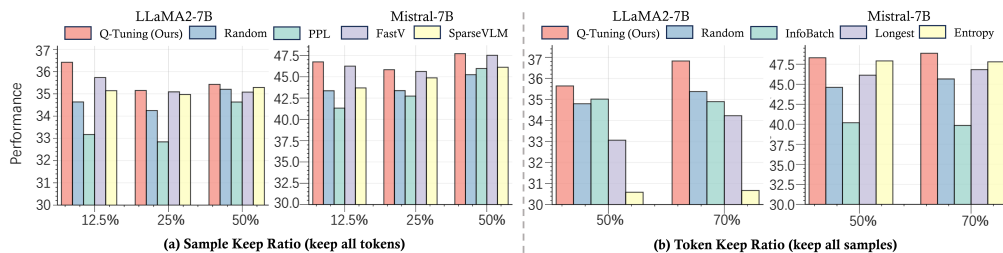


Figure 7: Comparison of independent (a) dynamic sample pruning and (b) dynamic token pruning across different keep ratios for LLaMA2-7B and Mistral-7B.

How does the sample distribution evolve during training?

To characterize training dynamics, we track the average perplexity and token entropy on 100 randomly sampled training instances throughout fine-tuning. As shown in Figure 8, Q-Tuning reduces both metrics more rapidly than alternative sample-pruning strategies, suggesting faster stabilization of the training signal and correlating with improved downstream performance.

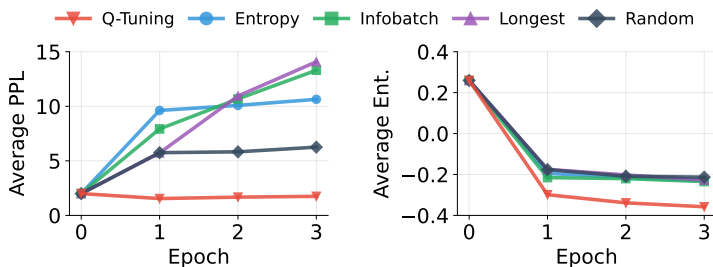


Figure 8: Training dynamics of different sample pruners. Compared to baseline strategies, Q-Tuning consistently reduces both metrics at a faster rate, indicating more efficient learning dynamics.

6 CONCLUSION

This work turns risky dynamic data pruning for LLM fine-tuning into a reliable, high-utility strategy. By analyzing sample modes via the Error–Uncertainty Plane, we expose data’s heterogeneous value and the need for nuance. Our Quadrant-based Tuning (Q-Tuning) uses a two-stage framework to coordinate sample- and token-level pruning, preserving valuable signals while removing noise, boosting efficiency without sacrificing overall stability, and often further improving performance across diverse tasks and models.

REFERENCES

- 540
541
542 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
543 Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report.
544 *arXiv preprint arXiv:2303.08774*, 2023.
- 545 Loubna Ben Allal, Anton Lozhkov, Elie Bakouch, Gabriel Martín Blázquez, Guilherme Penedo,
546 Lewis Tunstall, Andrés Marafioti, Hynek Kydlíček, Agustín Piqueres Lajarín, Vaibhav Srivastav,
547 et al. Smollm2: When smol goes big—data-centric training of a small language model. *arXiv*
548 *preprint arXiv:2502.02737*, 2025.
- 549 Abdul Hameed Azeemi, Ihsan Ayyub Qazi, and Agha Ali Raza. Language model-driven data pruning
550 enables efficient active learning. *arXiv preprint arXiv:2410.04275*, 2024.
- 551
552 Liang Chen, Haozhe Zhao, Tianyu Liu, Shuai Bai, Junyang Lin, Chang Zhou, and Baobao Chang.
553 An image is worth 1/2 tokens after layer 2: Plug-and-play inference acceleration for large vision-
554 language models. In *European Conference on Computer Vision*, pp. 19–35, 2024.
- 555 Lichang Chen, Shiyang Li, Jun Yan, Hai Wang, Kalpa Gunaratna, Vikas Yadav, Zheng Tang, Vijay
556 Srinivasan, Tianyi Zhou, Heng Huang, et al. Alpapasus: Training a better alpaca with fewer data.
557 *arXiv preprint arXiv:2307.08701*, 2023.
- 558
559 Tzu-Chun Chien, Chieh-Kai Lin, Shiang-Feng Tsai, Ruei-Chi Lai, Hung-Jen Chen, and Min Sun.
560 Grounding-aware token pruning: Recovering from drastic performance drops in visual grounding
561 caused by pruning. *arXiv preprint arXiv:2506.21873*, 2025.
- 562 Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and
563 Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge.
564 *arXiv preprint arXiv:1803.05457*, 2018.
- 565
566 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
567 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve
568 math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- 569 C Coleman, C Yeh, S Mussmann, B Mirzasoleiman, P Bailis, P Liang, J Leskovec, and M Zaharia.
570 Selection via proxy: Efficient data selection for deep learning. In *International Conference on*
571 *Learning Representations (ICLR)*, 2020.
- 572
573 OpenCompass Contributors. Opencompass: A universal evaluation platform for foundation models,
574 2023.
- 575 Xun Deng, Han Zhong, Rui Ai, Fuli Feng, Zheng Wang, and Xiangnan He. Less is more: Improving
576 llm alignment via preference data selection. *arXiv preprint arXiv:2502.14560*, 2025.
- 577
578 Guanting Dong, Hongyi Yuan, Keming Lu, Chengpeng Li, Mingfeng Xue, Dayiheng Liu, Wei Wang,
579 Zheng Yuan, Chang Zhou, and Jingren Zhou. How abilities in large language models are affected
580 by supervised fine-tuning data composition. In *Proceedings of the 62nd Annual Meeting of the*
581 *Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 177–198, 2024.
- 582 Qichen Fu, Minsik Cho, Thomas Merth, Sachin Mehta, Mohammad Rastegari, and Mahyar Na-
583 jibi. Lazyllm: Dynamic token pruning for efficient long context llm inference. *arXiv preprint*
584 *arXiv:2407.14057*, 2024.
- 585 Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad
586 Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of
587 models. *arXiv preprint arXiv:2407.21783*, 2024.
- 588
589 Yangyang Guo and Mohan Kankanhalli. Scan: Bootstrapping contrastive pre-training for data
590 efficiency. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp.
591 3662–3672, 2025.
- 592 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,
593 and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv*
preprint arXiv:2103.03874, 2021.

- 594 Xijie Huang, Li Lyna Zhang, Kwang-Ting Cheng, Fan Yang, and Mao Yang. Fewer is more: Boosting
595 llm reasoning with reinforced context pruning. *arXiv preprint arXiv:2312.08901*, 2023.
596
- 597 Hamish Ivison, Noah A Smith, Hannaneh Hajishirzi, and Pradeep Dasigi. Data-efficient finetuning
598 using cross-task nearest neighbors. In *Findings of the Association for Computational Linguistics:
599 ACL 2023*, pp. 9036–9061, 2023.
- 600 Hamish Ivison, Muru Zhang, Faeze Brahman, Pang Wei Koh, and Pradeep Dasigi. Large-scale data
601 selection for instruction tuning. *arXiv preprint arXiv:2503.01807*, 2025.
602
- 603 Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,
604 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier,
605 L  lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas
606 Wang, Timoth  e Lacroix, and William El Sayed. Mistral 7b, 2023a.
- 607 Huiqiang Jiang, Qianhui Wu, Chin-Yew Lin, Yuqing Yang, and Lili Qiu. Llmlingua: Compressing
608 prompts for accelerated inference of large language models. In *Proceedings of the 2023 Conference
609 on Empirical Methods in Natural Language Processing*, pp. 13358–13376, 2023b.
- 610 Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. Triviaqa: A large scale distantly
611 supervised challenge dataset for reading comprehension. *arXiv preprint arXiv:1705.03551*, 2017.
612
- 613 Christopher Keith, Michael Robinson, Francis Duncan, Allan Worthington, Joseph Wilson, and Sofia
614 Harris. Optimizing large language models: A novel approach through dynamic token pruning.
615 2024.
- 616 Yucheng Li, Bo Dong, Frank Guerin, and Chenghua Lin. Compressing context to enhance inference
617 efficiency of large language models. In *Proceedings of the 2023 conference on empirical methods
618 in natural language processing*, pp. 6342–6353, 2023.
619
- 620 Zhenghao Lin, Zhibin Gou, Yeyun Gong, Xiao Liu, Ruochen Xu, Chen Lin, Yujiu Yang, Jian Jiao,
621 Nan Duan, Weizhu Chen, et al. Not all tokens are what you need for pretraining. *Advances in
622 Neural Information Processing Systems*, 37:29029–29063, 2024.
- 623 Wei Liu, Weihao Zeng, Keqing He, Yong Jiang, and Junxian He. What makes good data for
624 alignment? a comprehensive study of automatic data selection in instruction tuning. *arXiv preprint
625 arXiv:2312.15685*, 2023.
626
- 627 Lingkun Long, Rubing Yang, Yushi Huang, Desheng Hui, Ao Zhou, and Jianlei Yang. Sliminfer: Ac-
628 celerating long-context llm inference via dynamic token pruning. *arXiv preprint arXiv:2508.06447*,
629 2025.
- 630 S  ren Mindermann, Jan M Brauner, Muhammed T Razzak, Mrinank Sharma, Andreas Kirsch, Winnie
631 Xu, Benedikt H  ltgen, Aidan N Gomez, Adrien Morisot, Sebastian Farquhar, et al. Prioritized
632 training on points that are learnable, worth learning, and not yet learnt. In *International Conference
633 on Machine Learning*, pp. 15630–15649, 2022.
634
- 635 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong
636 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow
637 instructions with human feedback. *Advances in neural information processing systems*, 35:27730–
638 27744, 2022.
- 639 Zhuoshi Pan, Qianhui Wu, Huiqiang Jiang, Menglin Xia, Xufang Luo, Jue Zhang, Qingwei Lin,
640 Victor R  hle, Yuqing Yang, Chin-Yew Lin, et al. Llmlingua-2: Data distillation for efficient and
641 faithful task-agnostic prompt compression. In *ACL (Findings)*, 2024.
- 642 Jinlong Pang, Jiaheng Wei, Ankit Parag Shah, Zhaowei Zhu, Yaxuan Wang, Chen Qian, Yang Liu,
643 Yujia Bao, and Wei Wei. Improving data efficiency via curating llm-driven rating systems. *arXiv
644 preprint arXiv:2410.10877*, 2024.
645
- 646 Jinlong Pang, Na Di, Zhaowei Zhu, Jiaheng Wei, Hao Cheng, Chen Qian, and Yang Liu. Token clean-
647 ing: Fine-grained data selection for llm supervised fine-tuning. *arXiv preprint arXiv:2502.01968*,
2025.

- 648 Ziheng Qin, Kai Wang, Zangwei Zheng, Jianyang Gu, Xiangyu Peng, Daquan Zhou, Lei Shang,
649 Baigui Sun, Xuansong Xie, Yang You, et al. Infobatch: Lossless training speed up by unbiased
650 dynamic data pruning. In *The Twelfth International Conference on Learning Representations*,
651 2024.
- 652 Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100,000+ questions for
653 machine comprehension of text. *arXiv preprint arXiv:1606.05250*, 2016.
- 654 Ravi S Raju, Kyle Daruwalla, and Mikko Lipasti. Accelerating deep learning with dynamic data
655 pruning. *arXiv preprint arXiv:2111.12621*, 2021.
- 656 Aaditya Singh, Adam Fry, Adam Perelman, Adam Tart, Adi Ganesh, Ahmed El-Kishky, Aidan
657 McLaughlin, Aiden Low, AJ Ostrow, Akhila Ananthram, et al. Openai gpt-5 system card. *arXiv*
658 *preprint arXiv:2601.03267*, 2025.
- 659 Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut,
660 Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly
661 capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- 662 Teknium. Openhermes 2.5: An open dataset of synthetic data for generalist llm assistants, 2023.
- 663 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay
664 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation
665 and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 666 Shaobo Wang, Xiangqi Jin, Ziming Wang, Jize Wang, Jiajun Zhang, Kaixin Li, Zichen Wen, Zhong
667 Li, Conghui He, Xuming Hu, and Linfeng Zhang. Data whisperer: Efficient data selection for
668 task-specific llm fine-tuning via few-shot in-context learning. *Annual Meeting of the Association*
669 *for Computational Linguistics*, 2025a.
- 670 Shenzhi Wang, Le Yu, Chang Gao, Chujiu Zheng, Shixuan Liu, Rui Lu, Kai Dang, Xionghui Chen,
671 Jianxin Yang, Zhenru Zhang, et al. Beyond the 80/20 rule: High-entropy minority tokens drive
672 effective reinforcement learning for llm reasoning. *arXiv preprint arXiv:2506.01939*, 2025b.
- 673 Zichen Wen, Yifeng Gao, Shaobo Wang, Junyuan Zhang, Qintong Zhang, Weijia Li, Conghui He, and
674 Linfeng Zhang. Stop looking for important tokens in multimodal language models: Duplication
675 matters more. *arXiv preprint arXiv:2502.11494*, 2025.
- 676 Heming Xia, Chak Tou Leong, Wenjie Wang, Yongqi Li, and Wenjie Li. Tokenskip: Controllable
677 chain-of-thought compression in llms. *arXiv preprint arXiv:2502.12067*, 2025.
- 678 Mengzhou Xia, Sadhika Malladi, Suchin Gururangan, Sanjeev Arora, and Danqi Chen. Less:
679 Selecting influential data for targeted instruction tuning. *arXiv preprint arXiv:2402.04333*, 2024a.
- 680 Tingyu Xia, Bowen Yu, Kai Dang, An Yang, Yuan Wu, Yuan Tian, Yi Chang, and Junyang Lin.
681 Rethinking data selection at scale: Random selection is almost all you need. *arXiv preprint*
682 *arXiv:2410.09335*, 2024b.
- 683 Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, Qingwei
684 Lin, and Daxin Jiang. Wizardlm: Empowering large pre-trained language models to follow complex
685 instructions. In *The Twelfth International Conference on Learning Representations*, 2024.
- 686 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang
687 Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*,
688 2025.
- 689 Yu Yang, Siddhartha Mishra, Jeffrey Chiang, and Baharan Mirzasoleiman. Smalltolarge (s2l):
690 Scalable data selection for fine-tuning large language models by summarizing training trajectories
691 of small models. *Advances in Neural Information Processing Systems*, 37:83465–83496, 2024.
- 692 Yixin Ye, Zhen Huang, Yang Xiao, Ethan Chern, Shijie Xia, and Pengfei Liu. Limo: Less is more for
693 reasoning. *arXiv preprint arXiv:2502.03387*, 2025.

- 702 Yang Yu, Kai Han, Hang Zhou, Yehui Tang, Kaiqi Huang, Yunhe Wang, and Dacheng Tao. Llm
703 data selection and utilization via dynamic bi-level optimization. *arXiv preprint arXiv:2507.16178*,
704 2025.
- 705
706 Zichun Yu, Spandan Das, and Chenyan Xiong. Mates: Model-aware data selection for efficient
707 pretraining with data influence models. *Advances in Neural Information Processing Systems*, 37:
708 108735–108759, 2024.
- 709 Xiang Yue, Xingwei Qu, Ge Zhang, Yao Fu, Wenhao Huang, Huan Sun, Yu Su, and Wenhao Chen.
710 Mammoth: Building math generalist models through hybrid instruction tuning. *arXiv preprint*
711 *arXiv:2309.05653*, 2023.
- 712
713 Guibin Zhang, Haonan Dong, Zhixun Li, Dingshuo Chen, Kai Wang, Tianlong Chen, Yuxuan Liang,
714 Dawei Cheng, Kun Wang, et al. Gder: Safeguarding efficiency, balancing, and robustness via
715 prototypical graph pruning. *Advances in Neural Information Processing Systems*, 37:50285–50312,
716 2024a.
- 717 Jia Zhang, Chen-Xi Zhang, Yao Liu, Yi-Xuan Jin, Xiao-Wen Yang, Bo Zheng, Yi Liu, and Lan-Zhe
718 Guo. D3: Diversity, difficulty, and dependability-aware data selection for sample-efficient llm
719 instruction tuning. *arXiv preprint arXiv:2503.11441*, 2025.
- 720
721 Yuan Zhang, Chun-Kai Fan, Junpeng Ma, Wenzhao Zheng, Tao Huang, Kuan Cheng, Denis A
722 Gudovskiy, Tomoyuki Okuno, Yohei Nakata, Kurt Keutzer, et al. Sparsevlm: Visual token sparsifi-
723 cation for efficient vision-language model inference. In *Forty-second International Conference on*
724 *Machine Learning*, 2024b.
- 725
726 Yaowei Zheng, Richong Zhang, Junhao Zhang, YeYanhan YeYanhan, and Zheyang Luo. Llamafactory:
727 Unified efficient fine-tuning of 100+ language models. In *Proceedings of the 62nd Annual Meeting*
728 *of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, pp. 400–410,
2024.
- 729
730 Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia
731 Efrat, Ping Yu, Lili Yu, et al. Lima: Less is more for alignment. *Advances in Neural Information*
Processing Systems, 36:55006–55021, 2023a.
- 732
733 Daquan Zhou, Kai Wang, Jianyang Gu, Xiangyu Peng, Dongze Lian, Yifan Zhang, Yang You, and
734 Jiashi Feng. Dataset quantization. In *Proceedings of the IEEE/CVF International Conference on*
735 *Computer Vision*, pp. 17205–17216, 2023b.
- 736
737 Fan Zhou, Zengzhi Wang, Qian Liu, Junlong Li, and Pengfei Liu. Programming every example:
738 Lifting pre-training data quality like experts at scale. *arXiv preprint arXiv:2409.17115*, 2024.
- 739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

A RELATED WORK

A.1 SAMPLE PRUNING

A growing body of work explores sample pruning as a means to reduce training cost by selecting smaller yet higher-quality subsets of data without compromising performance (Iverson et al., 2025; Yang et al., 2024; Xia et al., 2024a; Zhou et al., 2024). For example, LIMA (Zhou et al., 2023a) shows that aligning LLMs requires only a small collection of high-quality instruction–response pairs rather than large-scale corpora. Broadly, existing approaches fall into static and dynamic pruning. In the static setting, data subsets are determined in advance using fixed criteria: SVP (Coleman et al., 2020) leverages proxy models to estimate sample importance, D3 (Zhang et al., 2025) combines diversity, difficulty, and dependability into weighted coresets, and Less is More (Deng et al., 2025) shows that carefully curated preference subsets can outperform full datasets by filtering noisy or redundant examples. Other static approaches target in-context examples: LIMO (Ye et al., 2025) removes redundant demonstrations while distilling essential reasoning patterns, and methods such as DEFT (Iverson et al., 2023), Alpagasus (Chen et al., 2023), and Data Whisperer (Wang et al., 2025a) automatically select or reweight demonstrations based on influence estimation, contribution to performance, or few-shot evaluation.

In contrast, dynamic pruning adapts sample usage throughout training. Early work by (Raju et al., 2021) proposed two methods, ϵ -greedy and UCB, which retain uncertain examples while discarding easier ones. Subsequent approaches explore alternative criteria for adaptive pruning: RHO-LOSS (Mindermann et al., 2022) assigns importance using high-quality reference models, InfoBatch (Qin et al., 2024) removes low-loss examples on the fly, and more recent methods such as MATES (Yu et al., 2024) and DWM (Yu et al., 2025) reweight or select samples according to the evolving model state. Beyond single-example pruning, several methods consider structured or representation-based criteria: ActivePrune (Azeemi et al., 2024) selects demonstrations per input via similarity metrics, SCAN (Guo & Kankanhalli, 2025) iteratively removes ill-matched and redundant data during contrastive pretraining, GDeR (Zhang et al., 2024a) models data as a prototypical graph to prune noisy or redundant samples, and DQ (Zhou et al., 2023b) clusters data in representation space and replaces each cluster with representative subsets.

A.2 TOKEN PRUNING

Parallel to sample pruning, recent research (Wang et al., 2025b; Pan et al., 2024; Keith et al., 2024) has increasingly focused on token-level pruning to retain only the most informative portions of input sequences. At pretraining time, methods such as Rho-1 (Lin et al., 2024) estimate token importance using a reference model and remove low-utility tokens. During fine-tuning and inference, pruning is typically dynamic: Token Cleaning (Pang et al., 2025) treats harmful tokens as noisy labels and prunes them based on their influence on parameter updates, while LazyLLM (Fu et al., 2024) and SlimInfer (Long et al., 2025) accelerate long-context inference by dropping tokens or hidden-state blocks using attention signals and cache mechanisms. TokenSkip (Xia et al., 2025) further reduces redundant reasoning by skipping semantically less important tokens. A related line performs input compression before inference, as in LLMLingua (Jiang et al., 2023b) and Selective Context (Li et al., 2023), which apply coarse-to-fine filtering with budget control to prune low-information tokens, phrases, or sentences. Beyond text-only models, multimodal pruning exploits cross-modal attention patterns: FastV (Chen et al., 2024) and SparseVLM (Zhang et al., 2024b) drop redundant visual tokens guided by attention, while LMTL (Huang et al., 2023) prunes unnecessary visual features to adapt computation. Recent methods incorporate task objectives, GAP (Chien et al., 2025) maintains spatial grounding by correcting positional misalignment during pruning, and DART (Wen et al., 2025) reduces computation via confidence-based early stopping.

Despite their progress, prior work typically treats sample pruning and token pruning as independent lines of research. Sample pruning methods focus on reducing the number of training examples but leave token-level redundancy unaddressed, while token pruning techniques emphasize sequence-level efficiency without considering redundancy across training samples. Such separation limits the potential for jointly optimizing efficiency and effectiveness. In this work, we bridge these directions by proposing a unified framework that integrates sample-level and token-level pruning, enabling models to simultaneously filter uninformative data and compress redundant inputs.

B EXPERIMENTAL SETTINGS

B.1 MODELS AND DATASETS.

We conducted experiments on language models of different scales, including *Mistral-7B* (Jiang et al., 2023a), *LLaMA2-7B* (Touvron et al., 2023), *LLaMA3-8B* (Grattafiori et al., 2024), *SmolLM2-1.7B* (Allal et al., 2025), *Qwen3-8B*, *Qwen-14B* and *Qwen-32B* (Yang et al., 2025). For Qwen3-series, we fine-tune on the *OpenHermes* (Teknium, 2023), which contains diverse instruction–response pairs for general-purpose alignment. For other models, we use two complementary datasets: *WizardLM* (Xu et al., 2024) for alignment, which provides high-quality instruction–response pairs for supervised fine-tuning (SFT), and *MathInstruct* (Yue et al., 2023) for reasoning, which contains mathematically focused instructions designed to assess arithmetic and symbolic problem-solving.

B.2 EVALUATION.

For alignment task, we used five standard benchmarks: ARC-E, ARC-C (Clark et al., 2018), GSM8K (Cobbe et al., 2021), SQuAD (Rajpurkar et al., 2016), and TriviaQA (Joshi et al., 2017), covering knowledge-intensive question answering, commonsense reasoning, and reading comprehension. For reasoning task, we evaluated models fine-tuned on MathInstruct using GSM8K and MATH (Hendrycks et al., 2021), two widely used benchmarks for mathematical problem solving.

B.3 BASELINES.

We constructed baselines by pairing sample-level and token-level pruning strategies. For sample-level pruning, we considered eight methods: (i) Random, which dropped samples uniformly at random; (ii) Longest, which removed the longest sequences; (iii) Entropy, which retained high-entropy samples; (iv) InfoBatch (Qin et al., 2024), an information-theoretic data selection method; (v) Alpargus (Chen et al., 2023), which leveraged a LLM as an automatic quality rater to filter low-quality instruction–response pairs; (vi) Deita (Liu et al., 2023), which selected instruction–response pairs by modeling instruction complexity, quality, and diversity; (vii) DS2 (Pang et al., 2024), a diversity-aware data selection method that corrected LLM-generated quality scores via a transition matrix; and (viii) Less (Xia et al., 2024a), which selected instruction-tuning data by measuring gradient-based influence on few-shot target examples. For token-level pruning, we adopted six methods applied to the retained data: (i) Random, which masked tokens uniformly at random; (ii) PPL, which removed high-perplexity tokens; (iii) FastV (Chen et al., 2024), which pruned tokens receiving the least final-layer attention; (iv) SparseVLM (Zhang et al., 2024b), which combined attention and hidden-state features to score token importance; (v) Rho-1 (Lin et al., 2024), which scored tokens by their excess loss relative to a reference model; and (vi) TokenCleaning (Pang et al., 2025), which removed tokens deemed noisy or uninformative based on token-level statistics.

B.4 IMPLEMENTATION DETAILS

All experiments were conducted using the *LLaMA-Factory* framework (Zheng et al., 2024), which provided utilities for training and evaluation across diverse large language models. For benchmark evaluation, we used the *OpenCompass* framework (Contributors, 2023), which offered a standardized interface to a wide range of tasks. Unless otherwise specified, models were trained for 3 epochs with a learning rate of 1×10^{-4} and a batch size of 8 per device. In all pruning experiments, we applied pruning exclusively to the answers. All experiments were performed on NVIDIA A100 GPUs. For the *Random* baseline, we repeated each experiment five times and reported the averaged results.

B.5 HYPERPARAMETERS

We provide the full hyperparameter settings used in our experiments in Table 7. Unless otherwise noted, all other parameters follow the default settings of the *LLaMA-Factory* framework.

Table 7: Complete hyperparameter configurations employed in our experiments.

Hyperparameter	Value	Hyperparameter	Value
λ	0.5	Batch size	8
Epochs	3	LR	1×10^{-4}
LoRA rank	8	Cutoff length	2048
Grad. accum.	4	Scheduler	cosine

C ADDITIONAL EXPERIMENTAL RESULTS

C.1 ADDITIONAL RESULTS ON ALIGNMENT DATASETS

The detailed results of additional experiments on alignment dataset are presented in Table 8 and Table 9, providing a comprehensive comparison across different pruning strategies.

Table 8: Evaluation on Wizard under different sample ratios (12.5%, 25%) and token ratios (50%, 70%), where \uparrow and \downarrow respectively denote improvements or degradations over the *Random-Random* baseline.

Sample Pruner	Token Pruner	LLaMA2-7B						Mistral-7B					
		ARC-E	ARC-C	GSM8K	SQuAD	TriviaQA	Avg.	ARC-E	ARC-C	GSM8K	SQuAD	TriviaQA	Avg.
Zero-Shot		53.44	38.98	5.31	12.18	43.00	30.58	66.67	46.10	18.35	10.01	43.77	36.98
12.5% Samples, 70% Tokens													
Random	Random	59.43	41.02	6.97	13.64	47.97	33.81	71.08	47.46	24.34	21.64	53.15	43.53
	PPL	60.14 ^{↑0.71}	43.39 ^{↑2.37}	6.22 ^{↓0.75}	12.18 ^{↓1.46}	48.18 ^{↑0.21}	34.02 ^{↑0.21}	70.72 ^{↓0.36}	47.80 ^{↓0.34}	25.09 ^{↑0.75}	21.28 ^{↓0.36}	53.83 ^{↑0.68}	43.74 ^{↑0.21}
	FastV	58.20 ^{↓1.23}	41.02 ^{↑0.00}	6.29 ^{↓0.68}	13.42 ^{↓0.22}	45.32 ^{↓2.65}	32.85 ^{↓0.96}	70.72 ^{↓0.36}	46.44 ^{↓1.02}	19.56 ^{↓4.78}	21.38 ^{↓0.26}	53.34 ^{↑0.19}	42.29 ^{↓1.24}
	SparseVLM	54.67 ^{↓4.76}	37.97 ^{↓3.05}	8.04 ^{↑1.07}	13.06 ^{↓0.58}	44.87 ^{↓3.10}	31.72 ^{↓2.09}	67.72 ^{↓3.36}	44.75 ^{↓2.71}	23.65 ^{↓0.69}	11.76 ^{↓9.88}	44.90 ^{↓8.25}	38.58 ^{↓4.95}
Longest	Random	59.44 ^{↑0.01}	43.39 ^{↑2.37}	7.35 ^{↑0.38}	15.59 ^{↑1.95}	50.02 ^{↑2.05}	35.15 ^{↑1.34}	73.37 ^{↑2.29}	48.81 ^{↑1.35}	27.82 ^{↑3.48}	21.31 ^{↓0.33}	55.77 ^{↑2.62}	45.42 ^{↑1.89}
	PPL	60.85 ^{↑1.42}	43.39 ^{↑2.37}	7.73 ^{↑0.76}	16.21 ^{↑2.57}	48.57 ^{↑0.60}	35.35 ^{↑1.54}	74.96 ^{↑3.88}	49.83 ^{↑2.37}	28.73 ^{↑4.39}	21.62 ^{↓0.02}	56.59 ^{↑3.44}	46.35 ^{↑2.82}
	FastV	59.44 ^{↑0.01}	42.71 ^{↑1.69}	6.29 ^{↓0.68}	14.53 ^{↑0.89}	47.46 ^{↓0.51}	34.09 ^{↑0.28}	74.07 ^{↑2.99}	49.83 ^{↑2.37}	24.18 ^{↓0.16}	25.74 ^{↑4.10}	55.86 ^{↑2.71}	45.94 ^{↑2.41}
	SparseVLM	54.85 ^{↓4.58}	37.97 ^{↓3.05}	7.05 ^{↑0.08}	11.20 ^{↓2.44}	44.16 ^{↓3.81}	31.04 ^{↓2.77}	69.14 ^{↓4.94}	44.75 ^{↓2.71}	31.01 ^{↑6.67}	6.25 ^{↓15.39}	52.94 ^{↓0.21}	40.82 ^{↓2.71}
InfoBatch	Random	59.26 ^{↓0.17}	42.37 ^{↑1.35}	6.22 ^{↓0.75}	16.10 ^{↑0.46}	47.72 ^{↓0.25}	34.33 ^{↑0.52}	70.19 ^{↓0.89}	47.80 ^{↓0.34}	20.77 ^{↓3.57}	19.03 ^{↓2.61}	52.13 ^{↓1.02}	41.98 ^{↓1.55}
	PPL	60.49 ^{↑1.66}	39.32 ^{↓1.70}	5.76 ^{↓1.21}	14.47 ^{↑0.83}	48.06 ^{↑0.09}	33.62 ^{↓0.19}	70.72 ^{↓0.36}	46.44 ^{↓1.02}	19.03 ^{↓3.31}	23.20 ^{↑1.56}	51.75 ^{↓1.40}	42.23 ^{↓1.80}
	FastV	58.55 ^{↓1.08}	43.39 ^{↑2.37}	5.53 ^{↓1.44}	13.13 ^{↓0.51}	47.64 ^{↓0.33}	33.65 ^{↓0.16}	69.49 ^{↓1.59}	43.39 ^{↓4.07}	16.68 ^{↓7.66}	25.27 ^{↑3.63}	51.47 ^{↓1.68}	41.26 ^{↓2.27}
	SparseVLM	56.61 ^{↓2.82}	38.31 ^{↓2.71}	5.76 ^{↓1.21}	12.47 ^{↓1.17}	44.49 ^{↓3.48}	31.53 ^{↓2.28}	68.25 ^{↓2.83}	44.41 ^{↓3.05}	23.73 ^{↓0.61}	9.07 ^{↓12.57}	45.73 ^{↓2.42}	38.24 ^{↓5.29}
Entropy	Random	61.02 ^{↑1.59}	43.05 ^{↑2.03}	7.66 ^{↑0.69}	14.11 ^{↑0.47}	48.44 ^{↑0.47}	34.86 ^{↑1.05}	73.37 ^{↑2.29}	49.83 ^{↑2.37}	23.05 ^{↓1.29}	16.52 ^{↓5.12}	55.18 ^{↑2.03}	43.59 ^{↑0.06}
	PPL	61.02 ^{↑1.59}	43.39 ^{↑2.37}	6.97 ^{↑0.00}	14.94 ^{↑1.30}	48.94 ^{↑0.97}	35.05 ^{↑1.24}	73.02 ^{↑1.94}	47.46 ^{↑0.00}	24.03 ^{↑0.68}	22.85 ^{↑1.21}	54.89 ^{↑1.74}	44.45 ^{↑0.92}
	FastV	58.73 ^{↓0.70}	43.39 ^{↑2.37}	6.14 ^{↓0.83}	14.23 ^{↑0.59}	47.03 ^{↓0.94}	33.90 ^{↑0.09}	74.07 ^{↑2.99}	50.85 ^{↑3.39}	24.94 ^{↑0.60}	23.79 ^{↑2.15}	55.94 ^{↑2.79}	45.92 ^{↑2.39}
	SparseVLM	54.85 ^{↓4.58}	37.29 ^{↓3.73}	6.52 ^{↓0.45}	12.73 ^{↓0.91}	46.24 ^{↓1.73}	31.53 ^{↓2.28}	68.08 ^{↓3.00}	44.41 ^{↓3.05}	26.38 ^{↑2.04}	11.06 ^{↓10.58}	46.68 ^{↓6.47}	39.32 ^{↓4.21}
Q-Tuning (Ours)		64.37 ^{↑4.94}	42.37 ^{↑1.35}	10.84 ^{↑3.87}	17.63 ^{↑3.99}	52.17 ^{↑4.20}	37.48 ^{↑3.67}	71.78 ^{↑0.70}	48.14 ^{↑0.68}	30.33 ^{↑6.00}	28.59 ^{↑6.95}	57.93 ^{↑4.78}	47.35 ^{↑3.82}
Full Dataset		61.55	42.37	8.64	13.80	50.45	35.36	71.25	45.76	26.68	31.81	53.67	45.84
25% Samples, 50% Tokens													
Random	Random	60.32	41.69	5.76	13.43	48.41	33.92	70.19	46.10	20.62	24.07	53.74	42.95
	PPL	60.32 ^{↑0.00}	42.03 ^{↑0.34}	7.51 ^{↑1.75}	15.94 ^{↑2.51}	48.58 ^{↑0.17}	34.87 ^{↑0.95}	69.66 ^{↓0.53}	47.46 ^{↑1.36}	19.86 ^{↓0.76}	19.51 ^{↓4.56}	53.74 ^{↑0.00}	42.05 ^{↓0.90}
	FastV	59.08 ^{↓1.24}	41.69 ^{↑0.00}	3.56 ^{↓2.20}	12.78 ^{↓0.65}	45.60 ^{↓2.81}	32.54 ^{↓1.38}	71.78 ^{↑1.59}	47.12 ^{↑1.02}	15.77 ^{↓4.85}	26.97 ^{↑2.90}	50.84 ^{↓2.90}	42.50 ^{↓0.45}
	SparseVLM	54.50 ^{↓5.82}	38.64 ^{↓3.05}	6.44 ^{↑0.68}	12.04 ^{↓1.39}	44.79 ^{↓3.62}	31.28 ^{↓2.64}	67.55 ^{↓3.64}	46.44 ^{↑0.34}	24.41 ^{↑3.79}	11.80 ^{↓12.27}	48.14 ^{↓5.60}	39.67 ^{↓3.28}
Longest	Random	61.20 ^{↑0.88}	42.03 ^{↑0.34}	7.88 ^{↑2.12}	15.40 ^{↑1.97}	48.29 ^{↓0.12}	34.96 ^{↑1.04}	73.54 ^{↑3.38}	48.14 ^{↑2.04}	23.73 ^{↑3.11}	26.34 ^{↑2.27}	54.06 ^{↑0.32}	45.16 ^{↑2.21}
	PPL	60.85 ^{↑0.53}	43.39 ^{↑1.70}	7.20 ^{↑1.44}	13.88 ^{↑0.45}	48.48 ^{↑0.07}	34.76 ^{↑0.84}	72.31 ^{↑2.12}	48.14 ^{↑2.04}	24.34 ^{↑3.72}	23.84 ^{↑0.23}	55.22 ^{↑1.48}	44.77 ^{↑1.82}
	FastV	59.08 ^{↓1.24}	42.71 ^{↑1.02}	5.16 ^{↓0.60}	14.00 ^{↑0.57}	47.47 ^{↓0.94}	33.68 ^{↓0.24}	72.66 ^{↑2.47}	46.10 ^{↑0.00}	18.88 ^{↓1.74}	31.52 ^{↑7.45}	52.13 ^{↓1.61}	44.26 ^{↑1.31}
	SparseVLM	56.61 ^{↓3.71}	37.29 ^{↓4.40}	7.58 ^{↑1.82}	12.09 ^{↓1.34}	44.76 ^{↓3.65}	31.66 ^{↓2.26}	66.84 ^{↓3.35}	44.41 ^{↓1.69}	29.42 ^{↑8.80}	11.22 ^{↓12.85}	48.47 ^{↓5.27}	40.07 ^{↓2.88}
InfoBatch	Random	58.73 ^{↓1.59}	40.68 ^{↓1.01}	6.67 ^{↑0.91}	9.95 ^{↓3.48}	48.98 ^{↑0.57}	33.00 ^{↓0.92}	70.55 ^{↑0.36}	46.44 ^{↑0.34}	21.53 ^{↑0.91}	23.93 ^{↑0.14}	52.14 ^{↓1.60}	42.92 ^{↓0.03}
	PPL	59.96 ^{↑0.47}	42.71 ^{↑1.02}	6.52 ^{↑0.76}	14.58 ^{↑1.15}	48.57 ^{↑0.16}	34.47 ^{↑0.55}	71.08 ^{↑0.89}	47.80 ^{↑1.70}	20.62 ^{↑0.00}	24.88 ^{↑0.81}	51.61 ^{↓2.13}	43.20 ^{↑0.25}
	FastV	59.08 ^{↓1.24}	42.37 ^{↑0.68}	3.03 ^{↓2.73}	11.13 ^{↓2.30}	47.50 ^{↓0.91}	32.63 ^{↓1.29}	69.31 ^{↓0.88}	44.41 ^{↓1.69}	14.48 ^{↓6.14}	23.63 ^{↓0.44}	49.16 ^{↓4.58}	40.00 ^{↓2.75}
	SparseVLM	55.73 ^{↓4.59}	39.66 ^{↓2.03}	5.31 ^{↓0.45}	11.66 ^{↓1.77}	43.25 ^{↓5.16}	31.12 ^{↓2.80}	67.20 ^{↓3.00}	45.76 ^{↓0.34}	23.58 ^{↑2.96}	9.63 ^{↓14.44}	46.09 ^{↓7.65}	38.45 ^{↓4.50}
Entropy	Random	60.49 ^{↑0.17}	41.69 ^{↑0.00}	7.51 ^{↑1.75}	15.94 ^{↑2.51}	48.76 ^{↑0.35}	34.88 ^{↑0.96}	70.19 ^{↑0.00}	47.12 ^{↑1.02}	22.44 ^{↑1.82}	27.35 ^{↑3.28}	54.78 ^{↑1.04}	44.38 ^{↑1.43}
	PPL	60.49 ^{↑0.17}	41.02 ^{↓0.67}	6.60 ^{↑0.84}	14.92 ^{↑1.49}	49.33 ^{↑0.92}	34.47 ^{↑0.55}	71.43 ^{↑1.24}	48.14 ^{↑2.04}	21.30 ^{↑0.68}	25.52 ^{↑1.45}	55.60 ^{↑1.86}	44.40 ^{↑1.45}
	FastV	58.91 ^{↓1.41}	41.69 ^{↑0.00}	6.07 ^{↑0.31}	12.79 ^{↓0.64}	46.11 ^{↓2.30}	33.11 ^{↓0.81}	72.31 ^{↑2.12}	47.46 ^{↑1.36}	18.04 ^{↓2.58}	25.96 ^{↑1.89}	52.46 ^{↓1.28}	43.25 ^{↑0.30}
	SparseVLM	54.67 ^{↓5.82}	38.64 ^{↓3.05}	6.90 ^{↑1.14}	11.64 ^{↓1.79}	45.03 ^{↓3.38}	31.38 ^{↓2.54}	68.25 ^{↓4.94}	44.07 ^{↓2.03}	26.69 ^{↑6.07}	9.69 ^{↓14.38}	47.24 ^{↓6.50}	39.19 ^{↓3.76}
Q-Tuning (Ours)		63.14 ^{↑2.82}	42.03 ^{↑0.34}	8.87 ^{↑3.11}	16.76 ^{↑3.33}	51.52 ^{↑3.11}	36.47 ^{↑2.55}	71.78 ^{↑1.59}	47.12 ^{↑1.02}	26.08 ^{↑5.46}	32.79 ^{↑8.72}	56.17 ^{↑2.43}	46.79 ^{↑3.84}
Full Dataset		61.55	42.37	8.64	13.80	50.45	35.36	71.25	45.76	26.68	31.81	53.67	45.84

Table 9: Evaluation on Wizard under different sample ratios (25%, 50%) and token ratios (50%, 70%), where \uparrow and \downarrow respectively denote improvements or degradations over the *Random-Random* baseline.

Sample Pruner	Token Pruner	LLaMA2-7B						Mistral-7B					
		ARC-E	ARC-C	GSM8K	SQuAD	TriviaQA	Avg.	ARC-E	ARC-C	GSM8K	SQuAD	TriviaQA	Avg.
Zero-Shot		53.44	38.98	5.31	12.18	43.00	30.58	66.67	46.10	18.35	10.01	43.77	36.98
25% Samples, 70% Tokens													
Random	Random	60.67	41.69	6.22	13.91	48.77	34.25	70.02	46.78	19.71	23.50	52.93	42.59
	PPL	60.32 [↓] _{40.35}	42.03 [↑] _{70.34}	5.91 [↓] _{40.31}	15.88 [↑] _{11.97}	48.56 [↓] _{40.21}	34.54 [↑] _{70.29}	69.49 [↓] _{40.53}	48.14 [↑] _{71.36}	21.08 [↑] _{71.37}	20.63 [↓] _{42.87}	53.73 [↑] _{70.80}	42.61 [↑] _{70.02}
	SparseVLM	55.73 [↓] _{44.94}	37.97 [↑] _{43.72}	7.20 [↑] _{70.98}	11.83 [↓] _{42.08}	44.96 [↓] _{43.81}	31.54 [↓] _{42.71}	67.72 [↓] _{42.30}	44.75 [↓] _{42.03}	25.32 [↓] _{45.61}	12.69 [↓] _{41.03}	47.46 [↓] _{45.47}	39.59 [↓] _{43.00}
Longest	Random	61.02 [↑] _{70.35}	41.69 [↑] _{70.00}	7.20 [↑] _{70.98}	15.61 [↑] _{71.70}	49.53 [↑] _{70.76}	35.01 [↑] _{70.76}	74.43 [↑] _{74.41}	47.80 [↑] _{71.02}	24.56 [↑] _{74.85}	26.79 [↑] _{73.29}	55.24 [↑] _{72.31}	45.76 [↑] _{73.17}
	PPL	61.20 [↑] _{70.53}	43.05 [↑] _{71.36}	7.28 [↑] _{71.06}	14.06 [↑] _{70.15}	48.87 [↑] _{70.10}	34.89 [↑] _{70.64}	72.49 [↑] _{72.47}	47.46 [↑] _{70.68}	25.55 [↑] _{75.84}	23.92 [↑] _{70.04}	55.26 [↑] _{72.33}	44.93 [↑] _{72.34}
	SparseVLM	59.26 [↓] _{41.41}	43.05 [↑] _{71.36}	5.53 [↓] _{40.69}	13.98 [↑] _{70.07}	48.03 [↑] _{40.74}	33.97 [↑] _{70.28}	72.84 [↑] _{72.82}	47.80 [↑] _{71.02}	19.79 [↑] _{70.08}	31.54 [↑] _{70.04}	53.24 [↑] _{70.31}	45.04 [↑] _{72.45}
InfoBatch	Random	59.79 [↑] _{40.88}	42.71 [↑] _{71.02}	5.84 [↑] _{40.38}	10.42 [↑] _{43.89}	48.99 [↑] _{70.22}	33.55 [↑] _{40.70}	70.37 [↑] _{70.35}	46.44 [↑] _{40.34}	21.00 [↑] _{71.29}	24.09 [↑] _{70.59}	53.75 [↑] _{70.59}	43.13 [↑] _{70.54}
	PPL	60.32 [↑] _{70.35}	43.39 [↑] _{71.70}	6.67 [↑] _{70.45}	13.88 [↑] _{70.55}	49.03 [↑] _{70.26}	34.66 [↑] _{70.41}	71.08 [↑] _{71.06}	46.10 [↑] _{40.68}	22.44 [↑] _{72.73}	22.07 [↑] _{41.43}	52.17 [↑] _{40.76}	42.77 [↑] _{70.18}
	SparseVLM	59.61 [↓] _{41.06}	42.71 [↑] _{71.02}	3.34 [↓] _{42.88}	12.18 [↑] _{41.73}	47.98 [↑] _{40.79}	33.16 [↑] _{41.09}	69.66 [↑] _{40.36}	45.76 [↑] _{41.02}	15.39 [↑] _{45.21}	24.75 [↑] _{71.25}	49.51 [↑] _{43.42}	41.01 [↑] _{41.58}
Entropy	Random	60.67 [↑] _{70.00}	41.02 [↑] _{40.67}	5.99 [↑] _{40.23}	14.92 [↑] _{11.01}	49.34 [↑] _{70.57}	34.39 [↑] _{70.14}	70.37 [↑] _{70.35}	48.47 [↑] _{71.69}	22.21 [↑] _{72.80}	27.57 [↑] _{74.07}	56.05 [↑] _{73.42}	44.93 [↑] _{72.34}
	PPL	61.55 [↑] _{70.88}	41.02 [↑] _{40.67}	5.38 [↑] _{40.84}	14.90 [↑] _{70.99}	49.72 [↑] _{70.95}	34.51 [↑] _{70.26}	70.90 [↑] _{70.88}	48.14 [↑] _{71.36}	21.15 [↑] _{71.44}	26.22 [↑] _{72.72}	55.83 [↑] _{72.90}	44.45 [↑] _{71.86}
	SparseVLM	59.79 [↑] _{40.88}	42.03 [↑] _{70.34}	6.44 [↑] _{70.22}	13.21 [↑] _{40.70}	47.58 [↑] _{41.19}	33.81 [↑] _{40.44}	71.43 [↑] _{71.41}	47.12 [↑] _{70.34}	20.24 [↑] _{70.53}	26.06 [↑] _{72.56}	54.40 [↑] _{71.47}	43.85 [↑] _{71.26}
Q-Tuning (Ours)		62.43 [↑] _{71.76}	42.37 [↑] _{70.68}	9.25 [↑] _{73.03}	19.53 [↑] _{75.62}	50.78 [↑] _{72.01}	36.87 [↑] _{72.62}	71.60 [↑] _{71.58}	47.12 [↑] _{70.34}	26.61 [↑] _{76.90}	30.64 [↑] _{77.14}	55.13 [↑] _{72.20}	46.22 [↑] _{73.63}
Full Dataset		61.55	42.37	8.64	13.80	50.45	35.36	71.25	45.76	26.68	31.81	53.67	45.84
50% Samples, 50% Tokens													
Random	Random	62.08	41.36	6.75	12.14	48.86	34.24	71.25	46.44	21.53	24.91	54.16	43.66
	PPL	61.02 [↓] _{41.06}	43.05 [↑] _{71.69}	6.82 [↑] _{70.07}	15.08 [↑] _{72.94}	49.10 [↑] _{70.24}	35.01 [↑] _{70.77}	72.49 [↑] _{71.24}	46.78 [↑] _{70.34}	22.21 [↑] _{70.68}	33.28 [↑] _{72.37}	53.83 [↑] _{40.33}	45.72 [↑] _{72.06}
	SparseVLM	60.32 [↓] _{41.76}	42.03 [↑] _{70.67}	4.40 [↓] _{42.35}	11.80 [↓] _{40.34}	48.35 [↓] _{40.51}	33.38 [↓] _{40.86}	72.13 [↓] _{70.88}	45.42 [↓] _{41.02}	16.22 [↓] _{45.31}	27.69 [↑] _{72.78}	50.58 [↑] _{43.58}	42.41 [↓] _{41.25}
Longest	Random	60.85 [↑] _{41.23}	42.03 [↑] _{70.67}	7.05 [↑] _{70.30}	14.90 [↑] _{72.76}	49.46 [↑] _{70.60}	34.86 [↑] _{70.62}	72.49 [↑] _{71.24}	47.46 [↑] _{71.02}	21.91 [↑] _{70.38}	29.41 [↑] _{74.50}	55.99 [↑] _{71.83}	45.45 [↑] _{71.79}
	PPL	62.08 [↑] _{70.00}	41.69 [↑] _{70.33}	7.51 [↑] _{70.76}	14.89 [↑] _{72.75}	48.99 [↑] _{70.13}	35.03 [↑] _{70.79}	71.60 [↑] _{70.35}	46.78 [↑] _{70.34}	24.11 [↑] _{72.58}	29.17 [↑] _{74.26}	55.83 [↑] _{71.67}	45.50 [↑] _{71.84}
	SparseVLM	60.32 [↓] _{41.76}	40.68 [↓] _{40.68}	4.17 [↓] _{42.58}	13.26 [↓] _{41.12}	48.24 [↓] _{40.62}	33.33 [↓] _{40.91}	72.31 [↓] _{71.06}	45.76 [↓] _{40.68}	16.83 [↓] _{44.70}	30.37 [↑] _{75.46}	53.99 [↑] _{40.17}	43.85 [↑] _{71.89}
InfoBatch	Random	58.38 [↑] _{43.70}	42.71 [↑] _{71.35}	5.76 [↑] _{40.99}	13.52 [↑] _{71.38}	48.95 [↑] _{70.09}	33.86 [↑] _{40.38}	71.08 [↑] _{70.17}	46.44 [↑] _{70.00}	21.38 [↑] _{40.15}	26.73 [↑] _{71.82}	53.67 [↑] _{40.49}	43.86 [↑] _{70.20}
	PPL	60.67 [↑] _{41.41}	42.71 [↑] _{71.35}	5.76 [↑] _{40.99}	14.29 [↑] _{72.15}	49.18 [↑] _{70.32}	34.52 [↑] _{70.28}	71.43 [↑] _{70.18}	47.80 [↑] _{71.36}	24.11 [↑] _{72.58}	30.22 [↑] _{75.31}	53.97 [↑] _{40.19}	45.50 [↑] _{71.84}
	SparseVLM	59.44 [↓] _{42.84}	41.36 [↑] _{70.00}	3.87 [↓] _{42.88}	12.39 [↓] _{70.25}	48.25 [↓] _{40.61}	33.06 [↓] _{41.18}	70.72 [↓] _{40.53}	45.42 [↓] _{41.02}	13.87 [↓] _{47.66}	27.83 [↑] _{72.92}	51.53 [↑] _{42.63}	41.88 [↑] _{41.78}
Entropy	Random	59.96 [↑] _{42.12}	42.71 [↑] _{71.35}	7.66 [↑] _{70.91}	15.47 [↑] _{73.33}	49.07 [↑] _{70.21}	34.97 [↑] _{70.73}	72.13 [↑] _{70.88}	47.80 [↑] _{71.36}	22.67 [↑] _{71.14}	25.19 [↑] _{70.28}	55.70 [↑] _{71.54}	44.70 [↑] _{71.04}
	PPL	60.49 [↑] _{41.59}	41.02 [↑] _{40.34}	6.67 [↑] _{40.08}	13.51 [↑] _{71.37}	49.75 [↑] _{70.89}	34.29 [↑] _{70.05}	72.13 [↑] _{70.88}	47.80 [↑] _{71.36}	25.17 [↑] _{73.64}	30.25 [↑] _{75.34}	56.58 [↑] _{72.42}	46.39 [↑] _{72.73}
	SparseVLM	59.26 [↓] _{42.82}	42.37 [↑] _{71.01}	4.02 [↓] _{42.73}	13.81 [↑] _{71.67}	48.93 [↑] _{70.07}	33.68 [↑] _{40.56}	71.60 [↑] _{70.35}	47.46 [↑] _{71.02}	17.97 [↑] _{43.56}	29.82 [↑] _{74.91}	54.40 [↑] _{40.24}	44.25 [↑] _{70.59}
Q-Tuning (Ours)		62.79 [↑] _{70.71}	42.03 [↑] _{70.67}	10.46 [↑] _{73.71}	14.53 [↑] _{72.39}	51.05 [↑] _{72.19}	36.17 [↑] _{71.93}	73.37 [↑] _{72.12}	48.14 [↑] _{71.70}	28.81 [↑] _{77.28}	36.35 [↑] _{71.44}	56.30 [↑] _{72.14}	48.59 [↑] _{74.93}
Full Dataset		61.55	42.37	8.64	13.80	50.45	35.36	71.25	45.76	26.68	31.81	53.67	45.84
50% Samples, 70% Tokens													
Random	Random	61.02	41.36	7.43	15.56	48.91	34.85	71.60	47.12	22.59	27.86	53.68	44.57
	PPL	60.67 [↑] _{40.35}	42.71 [↑] _{71.35}	6.75 [↑] _{70.68}	14.93 [↑] _{40.63}	49.13 [↑] _{70.22}	34.84 [↑] _{40.01}	72.31 [↑] _{70.71}	47.12 [↑] _{70.00}	20.70 [↑] _{41.89}	32.02 [↑] _{74.16}	53.56 [↑] _{40.12}	45.14 [↑] _{40.57}
	SparseVLM	59.61 [↓] _{41.41}	42.37 [↑] _{71.01}	5.53 [↓] _{41.90}	13.36 [↓] _{42.30}	48.48 [↓] _{40.43}	33.87 [↓] _{40.98}	72.31 [↓] _{70.71}	45.08 [↓] _{42.04}	18.95 [↓] _{43.64}	28.60 [↑] _{70.74}	51.40 [↑] _{42.28}	43.27 [↑] _{41.30}
Longest	Random	61.20 [↑] _{70.18}	41.02 [↑] _{40.34}	6.37 [↑] _{41.06}	14.90 [↑] _{40.66}	49.79 [↑] _{70.88}	34.65 [↑] _{40.20}	72.13 [↑] _{70.53}	48.47 [↑] _{71.38}	23.05 [↑] _{70.46}	26.60 [↑] _{71.74}	55.94 [↑] _{72.26}	45.84 [↑] _{71.27}
	PPL	62.26 [↑] _{71.24}	41.36 [↑] _{70.00}	6.82 [↑] _{40.61}	14.83 [↑] _{40.73}	49.04 [↑] _{70.13}	34.86 [↑] _{70.70}	70.90 [↑] _{40.70}	47.80 [↑] _{70.68}	24.87 [↑] _{72.28}	27.33 [↑] _{40.53}	55.62 [↑] _{71.94}	45.30 [↑] _{70.73}
	SparseVLM	59.61 [↓] _{41.41}	42.71 [↑] _{71.38}	5.53 [↓] _{41.90}	14.12 [↓] _{41.44}	48.77 [↑] _{40.14}	34.15 [↑] _{40.70}	72.66 [↑] _{71.06}	45.08 [↑] _{42.04}	19.94 [↑] _{42.68}	31.53 [↑] _{73.67}	54.66 [↑] _{70.98}	44.77 [↑] _{70.20}
InfoBatch	Random	61.55 [↑] _{70.83}	40.34 [↓] _{41.02}	5.84 [↓] _{41.59}	14.99 [↑] _{40.57}	49.45 [↑] _{70.54}	34.43 [↑] _{40.42}	70.90 [↑] _{40.70}	45.76 [↑] _{41.36}	22.97 [↑] _{70.38}	28.85 [↑] _{70.99}	54.11 [↑] _{70.43}	44.52 [↑] _{40.05}
	PPL	60.85 [↑] _{40.17}	42.03 [↑] _{70.67}	6.75 [↑] _{40.68}	14.22 [↑] _{41.34}	48.41 [↑] _{40.50}	34.45 [↑] _{40.40}	71.08 [↑] _{40.52}	45.76 [↑] _{41.36}	22.82 [↑] _{70.23}	32.80 [↑] _{74.94}	54.63 [↑] _{70.95}	45.42 [↑] _{70.85}
	SparseVLM	58.91 [↓] _{42.11}	41.36 [↑] _{70.00}	4.02 [↓] _{43.41}	14.98 [↑] _{40.58}	48.76 [↑] _{40.15}	33.60 [↑] _{41.25}	70.19 [↑] _{41.41}	45.76 [↑] _{41.36}	19.86 [↑] _{42.73}	28.95 [↑] _{71.09}	50.92 [↑] _{42.76}	43.14 [↑] _{41.43</}

972 C.2 RESULTS ON QWEN3-8B ON ALIGNMENT DATASETS
973

974 The detailed results on the OpenHermes dataset using Qwen3-8B (Yang et al., 2025) are summarized
975 in Tables 10 and 11, which systematically evaluate different sample ratios and token ratios.
976

977 **Strong performance with limited data budgets.** Across all evaluated budgets, Q-Tuning consistently
978 outperforms traditional pruning-based fine-tuning methods and surpasses the full-data baseline. For
979 instance, under the extremely low-budget setting of 12.5% samples and 50% tokens, Q-Tuning
980 reaches an average of 60.19, outperforming the full-dataset result (55.45) by more than +4.7 points.
981 Even at 12.5% samples, 70% tokens, Q-Tuning maintains the same strong average (60.29), showing
982 remarkable robustness under aggressive pruning.

983 **Consistent superiority under larger budgets.** When the available data increases, Q-Tuning contin-
984 ues to dominate other baselines. At 25% samples, 50% tokens, Q-Tuning attains an average score of
985 58.80, exceeding the strongest baseline (InfoBatch–Random) by nearly +3.9 points. Similarly, at 25%
986 samples, 70% tokens, it reaches 59.20, still outperforming all competing methods such as InfoBatch,
987 Entropy, and SparseVLM by a clear margin, and exceeding the full-data fine-tuning by +3.75 points.
988

989 **Performance saturation and stable generalization.** At even larger budgets, such as 50% samples
990 with 50% or 70% tokens, Q-Tuning continues to deliver strong and stable improvements. In the 50%
991 / 50% setting, it achieves an average of 58.39, improving over all pruning-based baselines by 3–6
992 points. Under 50% / 70%, it reaches 57.86, maintaining competitive accuracy and demonstrating
993 consistent generalization despite heavy token-level sparsity.
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

Table 10: Evaluation on OpenHermes using Qwen3-8B under different sample ratios (12.5%, 25%) and token ratios (50%, 70%), where \uparrow and \downarrow respectively denote improvements or degradations over the *Random-Random* baseline.

Sample Pruner	Token Pruner	Qwen3-8B					
		ARC-E	ARC-C	GSM8K	SQuAD	TriviaQA	Avg.
Zero-Shot		69.31	44.41	74.45	27.63	51.76	53.51
12.5% Samples, 50% Tokens							
Random	Random	76.54	51.53	60.88	35.74	40.13	52.96
	PPL	77.60 \uparrow _{1.06}	52.20 \uparrow _{0.67}	58.00 \downarrow _{2.88}	35.55 \downarrow _{0.19}	40.73 \uparrow _{0.60}	52.82 \downarrow _{0.14}
	FastV	76.19 \downarrow _{0.35}	52.54 \uparrow _{1.01}	56.48 \downarrow _{4.40}	35.04 \downarrow _{0.70}	41.07 \uparrow _{0.94}	52.26 \downarrow _{0.70}
	SparseVLM	70.90 \downarrow _{5.64}	45.42 \downarrow _{6.11}	73.69 \uparrow _{12.81}	24.39 \downarrow _{11.35}	50.58 \uparrow _{10.45}	53.00 \uparrow _{0.04}
Longest	Random	78.48 \uparrow _{1.94}	53.22 \uparrow _{1.69}	62.70 \uparrow _{1.82}	33.97 \downarrow _{1.77}	42.13 \uparrow _{2.00}	54.10 \uparrow _{1.14}
	PPL	78.48 \uparrow _{1.94}	53.22 \uparrow _{1.69}	56.18 \downarrow _{4.70}	32.05 \downarrow _{3.69}	41.91 \uparrow _{1.78}	52.37 \downarrow _{0.59}
	FastV	77.60 \uparrow _{1.06}	51.86 \uparrow _{0.33}	60.12 \downarrow _{0.76}	32.57 \downarrow _{3.17}	43.01 \uparrow _{2.88}	53.03 \uparrow _{0.07}
	SparseVLM	69.14 \downarrow _{7.40}	44.75 \downarrow _{6.78}	75.59 \uparrow _{14.71}	22.24 \downarrow _{13.50}	50.02 \uparrow _{9.89}	52.35 \downarrow _{0.61}
InfoBatch	Random	80.25 \uparrow _{3.71}	53.56 \uparrow _{2.03}	60.58 \downarrow _{0.30}	38.51 \uparrow _{2.77}	44.36 \uparrow _{4.23}	55.45 \uparrow _{2.49}
	PPL	79.01 \uparrow _{2.47}	53.56 \uparrow _{2.03}	62.47 \uparrow _{1.59}	35.69 \downarrow _{0.05}	42.06 \uparrow _{1.93}	54.56 \uparrow _{1.60}
	FastV	79.01 \uparrow _{2.47}	51.19 \downarrow _{0.34}	45.49 \downarrow _{15.39}	38.71 \uparrow _{2.97}	44.72 \uparrow _{4.59}	51.82 \downarrow _{1.14}
	SparseVLM	68.96 \downarrow _{7.58}	44.41 \downarrow _{7.12}	77.71 \uparrow _{16.83}	19.88 \downarrow _{15.86}	49.35 \uparrow _{9.22}	52.06 \downarrow _{0.90}
Entropy	Random	76.37 \downarrow _{0.17}	51.53 \uparrow _{0.00}	62.77 \uparrow _{1.89}	33.90 \downarrow _{1.84}	43.07 \uparrow _{2.94}	53.53 \uparrow _{0.57}
	PPL	78.66 \uparrow _{2.12}	52.54 \uparrow _{1.01}	61.94 \uparrow _{1.06}	35.62 \downarrow _{0.12}	42.73 \uparrow _{2.60}	54.30 \uparrow _{1.34}
	FastV	76.90 \uparrow _{0.36}	53.56 \uparrow _{2.03}	65.05 \uparrow _{4.17}	36.14 \uparrow _{0.40}	43.99 \uparrow _{3.86}	55.13 \uparrow _{2.17}
	SparseVLM	71.08 \downarrow _{4.56}	44.07 \downarrow _{7.46}	75.89 \uparrow _{15.01}	22.77 \downarrow _{13.00}	50.47 \uparrow _{10.34}	52.85 \downarrow _{0.11}
Q-Tuning (Ours)		80.95 \uparrow _{4.41}	53.56 \uparrow _{2.03}	70.51 \uparrow _{9.63}	45.76 \uparrow _{10.02}	50.16 \uparrow _{10.03}	60.19 \uparrow _{7.23}
12.5% Samples, 70% Tokens							
Random	Random	76.54	52.20	61.64	34.88	41.02	53.26
	PPL	77.07 \uparrow _{0.53}	52.20 \uparrow _{0.00}	59.21 \downarrow _{2.43}	35.20 \uparrow _{0.32}	40.87 \downarrow _{0.15}	52.91 \downarrow _{0.35}
	FastV	76.37 \downarrow _{0.17}	50.51 \downarrow _{1.69}	59.36 \downarrow _{2.28}	35.25 \uparrow _{0.37}	41.01 \downarrow _{0.01}	52.50 \downarrow _{0.76}
	SparseVLM	70.90 \downarrow _{5.64}	45.42 \downarrow _{6.78}	73.69 \uparrow _{12.05}	24.39 \downarrow _{10.49}	50.58 \uparrow _{9.56}	53.00 \downarrow _{0.26}
Longest	Random	79.19 \uparrow _{2.65}	53.90 \uparrow _{1.70}	61.64 \uparrow _{0.00}	34.14 \downarrow _{0.74}	42.23 \uparrow _{1.21}	54.22 \uparrow _{0.96}
	PPL	78.48 \uparrow _{1.94}	51.53 \downarrow _{0.67}	55.65 \downarrow _{5.99}	32.95 \downarrow _{1.93}	42.45 \uparrow _{1.43}	52.21 \downarrow _{1.05}
	FastV	78.13 \uparrow _{1.59}	51.53 \downarrow _{0.67}	59.21 \downarrow _{2.43}	33.16 \downarrow _{1.72}	42.19 \uparrow _{1.17}	52.84 \downarrow _{0.42}
	SparseVLM	69.14 \downarrow _{7.40}	44.75 \downarrow _{7.45}	75.59 \uparrow _{13.95}	22.24 \downarrow _{12.64}	50.02 \uparrow _{9.00}	52.35 \downarrow _{0.91}
InfoBatch	Random	78.84 \uparrow _{2.30}	51.86 \downarrow _{0.34}	58.38 \downarrow _{3.26}	38.12 \uparrow _{3.24}	39.98 \downarrow _{1.04}	53.44 \uparrow _{0.18}
	PPL	78.48 \uparrow _{1.94}	52.88 \uparrow _{0.68}	62.17 \uparrow _{0.53}	36.28 \uparrow _{1.40}	42.96 \uparrow _{1.94}	54.55 \uparrow _{1.29}
	FastV	79.89 \uparrow _{3.35}	53.22 \uparrow _{1.02}	55.19 \downarrow _{6.45}	38.63 \uparrow _{3.75}	44.23 \uparrow _{3.21}	54.23 \uparrow _{0.97}
	SparseVLM	68.96 \downarrow _{7.58}	44.41 \downarrow _{7.79}	77.71 \uparrow _{16.07}	19.88 \downarrow _{15.00}	49.35 \uparrow _{7.33}	52.06 \downarrow _{1.20}
Entropy	Random	77.07 \uparrow _{0.53}	51.86 \downarrow _{0.34}	60.73 \downarrow _{0.91}	34.30 \downarrow _{0.58}	42.34 \uparrow _{1.32}	53.26 \uparrow _{0.00}
	PPL	78.84 \uparrow _{2.30}	52.88 \uparrow _{0.68}	61.26 \downarrow _{0.38}	36.44 \uparrow _{1.56}	43.53 \uparrow _{2.51}	54.59 \uparrow _{1.33}
	FastV	77.07 \uparrow _{0.53}	51.86 \downarrow _{0.34}	64.22 \uparrow _{2.58}	36.63 \uparrow _{1.75}	43.94 \uparrow _{2.92}	54.74 \uparrow _{1.48}
	SparseVLM	71.08 \downarrow _{4.56}	44.07 \downarrow _{8.13}	75.89 \uparrow _{14.25}	22.77 \downarrow _{12.11}	50.47 \uparrow _{9.45}	52.85 \downarrow _{0.41}
Q-Tuning (Ours)		80.95 \uparrow _{4.41}	52.20 \uparrow _{0.00}	70.13 \uparrow _{8.49}	47.60 \uparrow _{12.72}	50.58 \uparrow _{9.56}	60.29 \uparrow _{7.03}
25% Samples, 50% Tokens							
Random	Random	77.07	51.86	61.03	35.11	40.55	53.13
	PPL	77.78 \uparrow _{0.71}	54.58 \uparrow _{2.72}	58.53 \downarrow _{2.50}	35.43 \uparrow _{0.32}	40.68 \uparrow _{0.13}	53.40 \uparrow _{0.27}
	FastV	79.54 \uparrow _{2.47}	52.54 \uparrow _{0.68}	53.45 \downarrow _{7.58}	36.89 \uparrow _{1.78}	43.13 \uparrow _{2.58}	53.11 \downarrow _{0.02}
	SparseVLM	69.14 \downarrow _{7.93}	45.76 \downarrow _{6.10}	77.71 \uparrow _{16.68}	20.80 \downarrow _{14.31}	50.20 \uparrow _{9.65}	52.72 \downarrow _{0.41}
Longest	Random	79.01 \uparrow _{1.94}	54.58 \uparrow _{2.72}	51.10 \downarrow _{9.93}	36.30 \uparrow _{1.19}	41.93 \uparrow _{1.38}	52.58 \downarrow _{0.55}
	PPL	79.01 \uparrow _{1.94}	54.92 \uparrow _{3.06}	55.42 \downarrow _{5.61}	35.10 \downarrow _{0.01}	41.09 \uparrow _{0.54}	53.11 \downarrow _{0.02}
	FastV	80.95 \uparrow _{3.88}	51.53 \downarrow _{0.33}	53.30 \downarrow _{7.73}	39.73 \uparrow _{4.62}	44.34 \uparrow _{3.79}	53.97 \uparrow _{0.84}
	SparseVLM	72.49 \downarrow _{4.58}	44.41 \downarrow _{7.45}	74.83 \uparrow _{13.80}	22.08 \downarrow _{13.03}	48.93 \uparrow _{7.38}	52.55 \downarrow _{0.58}
InfoBatch	Random	77.25 \uparrow _{0.18}	51.53 \downarrow _{0.33}	63.61 \uparrow _{2.58}	36.99 \uparrow _{1.88}	45.15 \uparrow _{4.60}	54.91 \uparrow _{1.78}
	PPL	78.84 \uparrow _{1.77}	54.58 \uparrow _{2.72}	60.20 \downarrow _{0.83}	35.81 \uparrow _{0.70}	43.53 \uparrow _{3.00}	54.59 \uparrow _{1.46}
	FastV	80.25 \uparrow _{3.18}	52.54 \uparrow _{0.68}	48.52 \downarrow _{12.51}	39.62 \uparrow _{4.51}	44.87 \uparrow _{4.32}	53.16 \uparrow _{0.03}
	SparseVLM	69.49 \downarrow _{7.58}	44.07 \downarrow _{7.79}	76.88 \uparrow _{15.85}	18.29 \downarrow _{16.82}	50.16 \uparrow _{9.61}	51.78 \downarrow _{1.35}
Entropy	Random	77.07 \uparrow _{0.00}	53.22 \uparrow _{1.36}	62.02 \uparrow _{1.00}	35.40 \uparrow _{0.29}	42.91 \uparrow _{2.36}	54.12 \uparrow _{0.99}
	PPL	79.37 \uparrow _{2.30}	54.24 \uparrow _{2.38}	52.24 \downarrow _{8.79}	33.13 \downarrow _{1.98}	40.82 \uparrow _{0.27}	51.96 \downarrow _{1.17}
	FastV	79.54 \uparrow _{2.47}	51.86 \uparrow _{0.00}	52.92 \downarrow _{8.11}	35.68 \uparrow _{0.57}	43.79 \uparrow _{3.24}	52.76 \downarrow _{0.37}
	SparseVLM	70.19 \downarrow _{6.88}	43.39 \downarrow _{8.47}	75.21 \uparrow _{14.18}	19.97 \downarrow _{15.14}	50.46 \uparrow _{9.91}	51.84 \downarrow _{1.29}
Q-Tuning (Ours)		79.89 \uparrow _{2.82}	52.88 \uparrow _{1.02}	69.37 \uparrow _{78.34}	43.54 \uparrow _{78.43}	48.33 \uparrow _{77.78}	58.80 \uparrow _{75.67}
Full Dataset		77.78	52.20	65.13	38.05	44.08	55.45

Table 11: Evaluation on OpenHermes using Qwen3-8B under different sample ratios (25%, 50%) and token ratios (50%, 70%), where \uparrow and \downarrow respectively denote improvements or degradations over the *Random-Random* baseline.

Sample Pruner	Token Pruner	Qwen3-8B					
		ARC-E	ARC-C	GSM8K	SQuAD	TriviaQA	Avg.
Zero-Shot		69.31	44.41	74.45	27.63	51.76	53.51
25% Samples, 70% Tokens							
Random	Random	77.43	52.20	57.85	35.55	41.47	52.90
	PPL	76.90 $\downarrow_{0.53}$	53.90 $\uparrow_{1.70}$	60.27 $\uparrow_{2.42}$	35.70 $\uparrow_{0.15}$	40.73 $\downarrow_{0.74}$	53.50 $\uparrow_{0.60}$
	FastV	79.19 $\uparrow_{1.76}$	51.86 $\downarrow_{0.34}$	56.03 $\downarrow_{1.82}$	36.01 $\uparrow_{0.46}$	43.86 $\uparrow_{2.39}$	53.39 $\uparrow_{0.49}$
	SparseVLM	69.14 $\downarrow_{8.29}$	45.76 $\downarrow_{6.44}$	77.71 $\uparrow_{9.86}$	20.80 $\downarrow_{14.75}$	50.20 $\uparrow_{8.73}$	52.72 $\downarrow_{0.18}$
Longest	Random	79.01 $\uparrow_{1.58}$	54.58 $\uparrow_{2.38}$	53.83 $\downarrow_{4.02}$	37.56 $\uparrow_{2.01}$	42.40 $\uparrow_{0.93}$	53.48 $\uparrow_{0.58}$
	PPL	78.48 $\uparrow_{1.05}$	55.59 $\uparrow_{3.39}$	55.72 $\downarrow_{2.13}$	35.16 $\downarrow_{0.39}$	41.31 $\downarrow_{0.16}$	53.25 $\uparrow_{0.35}$
	FastV	79.89 $\uparrow_{2.46}$	52.54 $\uparrow_{0.34}$	53.22 $\downarrow_{4.63}$	37.90 $\uparrow_{2.35}$	43.79 $\uparrow_{2.32}$	53.47 $\uparrow_{0.57}$
	SparseVLM	72.49 $\downarrow_{4.94}$	44.41 $\downarrow_{7.79}$	74.83 $\uparrow_{6.98}$	22.08 $\downarrow_{13.47}$	48.93 $\uparrow_{7.46}$	52.55 $\downarrow_{0.35}$
InfoBatch	Random	78.66 $\uparrow_{1.23}$	52.54 $\uparrow_{0.34}$	63.68 $\uparrow_{5.83}$	37.56 $\uparrow_{2.01}$	44.51 $\uparrow_{3.04}$	55.39 $\uparrow_{2.49}$
	PPL	78.84 $\uparrow_{1.41}$	53.22 $\uparrow_{1.02}$	62.40 $\uparrow_{4.55}$	36.39 $\uparrow_{0.84}$	44.87 $\uparrow_{3.40}$	55.14 $\uparrow_{2.24}$
	FastV	79.89 $\uparrow_{2.46}$	53.90 $\uparrow_{1.70}$	55.04 $\downarrow_{2.82}$	38.79 $\uparrow_{3.24}$	44.28 $\uparrow_{2.81}$	54.38 $\uparrow_{1.48}$
	SparseVLM	69.49 $\downarrow_{7.94}$	44.07 $\downarrow_{8.13}$	76.88 $\uparrow_{9.03}$	18.29 $\downarrow_{17.26}$	50.16 $\uparrow_{8.69}$	51.78 $\downarrow_{1.12}$
Entropy	Random	76.90 $\downarrow_{0.53}$	52.54 $\uparrow_{0.34}$	57.47 $\downarrow_{0.38}$	34.58 $\downarrow_{0.97}$	42.50 $\uparrow_{1.03}$	52.80 $\downarrow_{0.10}$
	PPL	78.66 $\uparrow_{1.23}$	53.90 $\uparrow_{1.70}$	52.46 $\downarrow_{5.39}$	33.72 $\downarrow_{1.83}$	41.16 $\downarrow_{0.31}$	51.98 $\downarrow_{0.92}$
	FastV	79.01 $\uparrow_{1.58}$	52.20 $\uparrow_{0.00}$	57.54 $\downarrow_{0.31}$	35.72 $\uparrow_{0.17}$	42.94 $\uparrow_{1.47}$	53.48 $\uparrow_{0.58}$
	SparseVLM	70.19 $\downarrow_{7.24}$	43.39 $\downarrow_{8.81}$	75.21 $\uparrow_{17.36}$	19.97 $\downarrow_{15.58}$	50.46 $\uparrow_{9.00}$	51.84 $\downarrow_{0.66}$
Q-Tuning (Ours)		79.72 $\uparrow_{2.29}$	52.88 $\uparrow_{0.68}$	71.11 $\uparrow_{13.26}$	43.70 $\uparrow_{8.15}$	48.58 $\uparrow_{7.11}$	59.20 $\uparrow_{6.30}$
50% Samples, 50% Tokens							
Random	Random	77.95	51.86	63.61	36.64	42.22	54.46
	PPL	78.48 $\uparrow_{0.53}$	53.22 $\uparrow_{1.36}$	62.02 $\downarrow_{1.59}$	35.96 $\downarrow_{0.68}$	42.56 $\uparrow_{0.34}$	54.45 $\downarrow_{0.01}$
	FastV	79.19 $\uparrow_{1.24}$	51.53 $\downarrow_{0.33}$	49.73 $\downarrow_{13.88}$	40.64 $\uparrow_{4.00}$	44.45 $\uparrow_{2.23}$	53.11 $\downarrow_{1.35}$
	SparseVLM	72.13 $\downarrow_{5.82}$	46.10 $\downarrow_{5.76}$	73.69 $\uparrow_{10.08}$	18.74 $\downarrow_{17.90}$	50.39 $\uparrow_{8.17}$	52.21 $\downarrow_{2.25}$
Longest	Random	78.84 $\uparrow_{0.89}$	54.24 $\uparrow_{2.38}$	59.67 $\downarrow_{3.94}$	36.66 $\uparrow_{0.02}$	43.05 $\uparrow_{0.83}$	54.49 $\uparrow_{0.03}$
	PPL	79.72 $\uparrow_{1.77}$	54.24 $\uparrow_{2.38}$	58.15 $\downarrow_{5.46}$	39.88 $\uparrow_{3.24}$	44.01 $\uparrow_{1.79}$	55.20 $\uparrow_{0.74}$
	FastV	79.72 $\uparrow_{1.77}$	51.86 $\uparrow_{0.00}$	47.99 $\downarrow_{15.62}$	40.57 $\uparrow_{3.93}$	44.89 $\uparrow_{2.67}$	53.01 $\downarrow_{1.45}$
	SparseVLM	70.19 $\downarrow_{7.76}$	43.39 $\downarrow_{8.47}$	74.22 $\uparrow_{10.61}$	17.90 $\downarrow_{18.74}$	49.34 $\uparrow_{7.12}$	51.01 $\downarrow_{3.45}$
InfoBatch	Random	78.84 $\uparrow_{0.89}$	52.88 $\uparrow_{1.02}$	64.37 $\uparrow_{0.76}$	36.98 $\uparrow_{0.34}$	45.49 $\uparrow_{3.27}$	55.71 $\uparrow_{1.25}$
	PPL	77.95 $\uparrow_{0.00}$	53.56 $\uparrow_{1.70}$	60.73 $\downarrow_{2.88}$	38.10 $\uparrow_{1.47}$	44.39 $\uparrow_{2.17}$	54.95 $\uparrow_{0.49}$
	FastV	80.42 $\uparrow_{2.47}$	52.88 $\uparrow_{1.02}$	50.34 $\downarrow_{13.27}$	40.79 $\uparrow_{4.15}$	45.14 $\uparrow_{2.92}$	53.92 $\downarrow_{0.54}$
	SparseVLM	66.49 $\downarrow_{11.46}$	43.39 $\downarrow_{8.47}$	74.30 $\uparrow_{10.69}$	15.16 $\downarrow_{21.48}$	49.64 $\uparrow_{7.42}$	49.80 $\downarrow_{4.66}$
Entropy	Random	78.13 $\uparrow_{0.18}$	51.53 $\downarrow_{0.33}$	61.94 $\downarrow_{1.67}$	34.59 $\downarrow_{2.05}$	42.23 $\uparrow_{0.01}$	53.68 $\downarrow_{0.78}$
	PPL	78.66 $\uparrow_{0.71}$	52.88 $\uparrow_{1.02}$	54.97 $\downarrow_{8.61}$	36.11 $\downarrow_{0.53}$	42.06 $\downarrow_{0.16}$	52.94 $\downarrow_{1.52}$
	FastV	79.54 $\uparrow_{1.59}$	53.22 $\uparrow_{1.36}$	47.08 $\downarrow_{16.53}$	40.24 $\uparrow_{3.58}$	43.33 $\uparrow_{1.11}$	52.68 $\downarrow_{1.78}$
	SparseVLM	70.90 $\downarrow_{7.05}$	43.39 $\downarrow_{8.47}$	78.70 $\uparrow_{15.09}$	16.25 $\downarrow_{20.39}$	50.10 $\uparrow_{7.87}$	51.87 $\downarrow_{2.59}$
Q-Tuning (Ours)		80.25 $\uparrow_{2.30}$	53.90 $\uparrow_{2.04}$	70.05 $\uparrow_{6.44}$	39.03 $\uparrow_{2.39}$	48.70 $\uparrow_{6.48}$	58.39 $\uparrow_{3.93}$
50% Samples, 70% Tokens							
Random	Random	78.13	52.20	61.49	36.91	42.75	54.30
	PPL	77.95 $\downarrow_{0.18}$	52.20 $\uparrow_{0.00}$	59.59 $\downarrow_{1.90}$	36.14 $\downarrow_{0.77}$	42.51 $\downarrow_{0.24}$	53.68 $\downarrow_{0.62}$
	FastV	79.01 $\uparrow_{0.88}$	53.90 $\uparrow_{1.70}$	58.30 $\downarrow_{3.19}$	40.33 $\uparrow_{3.42}$	45.20 $\uparrow_{2.45}$	55.35 $\uparrow_{1.05}$
	SparseVLM	72.13 $\downarrow_{6.00}$	46.10 $\downarrow_{6.10}$	73.69 $\uparrow_{12.20}$	18.74 $\downarrow_{18.17}$	50.39 $\uparrow_{7.64}$	52.21 $\downarrow_{2.09}$
Longest	Random	78.48 $\uparrow_{0.35}$	52.54 $\uparrow_{0.34}$	59.14 $\downarrow_{2.35}$	37.27 $\uparrow_{0.36}$	42.38 $\downarrow_{0.37}$	53.96 $\downarrow_{0.34}$
	PPL	79.19 $\uparrow_{1.06}$	53.56 $\uparrow_{1.36}$	60.96 $\downarrow_{0.53}$	39.78 $\uparrow_{2.87}$	44.17 $\uparrow_{1.42}$	55.53 $\uparrow_{1.23}$
	FastV	80.25 $\uparrow_{2.12}$	52.20 $\uparrow_{0.00}$	56.79 $\downarrow_{4.70}$	40.53 $\uparrow_{3.62}$	45.62 $\uparrow_{2.87}$	55.08 $\uparrow_{0.78}$
	SparseVLM	70.19 $\downarrow_{7.94}$	43.39 $\downarrow_{8.81}$	74.22 $\uparrow_{12.73}$	17.90 $\downarrow_{19.01}$	49.34 $\uparrow_{6.59}$	51.01 $\downarrow_{3.29}$
InfoBatch	Random	78.31 $\uparrow_{0.18}$	52.88 $\uparrow_{0.68}$	62.32 $\uparrow_{0.83}$	35.55 $\downarrow_{1.36}$	45.15 $\uparrow_{2.40}$	54.84 $\uparrow_{0.54}$
	PPL	78.31 $\uparrow_{0.18}$	54.58 $\uparrow_{2.38}$	63.23 $\uparrow_{1.74}$	38.89 $\uparrow_{1.98}$	44.16 $\uparrow_{1.41}$	55.83 $\uparrow_{1.53}$
	FastV	79.01 $\uparrow_{0.88}$	52.88 $\uparrow_{0.68}$	58.30 $\downarrow_{3.19}$	39.05 $\uparrow_{2.14}$	46.05 $\uparrow_{3.30}$	55.06 $\uparrow_{0.76}$
	SparseVLM	66.49 $\downarrow_{11.64}$	43.39 $\downarrow_{8.81}$	74.30 $\uparrow_{12.81}$	15.16 $\downarrow_{21.75}$	49.64 $\uparrow_{6.89}$	49.80 $\downarrow_{4.50}$
Entropy	Random	78.48 $\uparrow_{0.35}$	53.22 $\uparrow_{1.02}$	60.65 $\downarrow_{0.84}$	35.64 $\downarrow_{1.27}$	42.31 $\downarrow_{0.44}$	54.06 $\downarrow_{0.24}$
	PPL	78.13 $\uparrow_{0.00}$	53.22 $\uparrow_{1.02}$	57.85 $\downarrow_{3.64}$	36.85 $\downarrow_{0.06}$	41.70 $\downarrow_{1.05}$	53.55 $\downarrow_{0.75}$
	FastV	79.37 $\uparrow_{1.24}$	52.88 $\uparrow_{0.68}$	52.62 $\downarrow_{8.87}$	39.76 $\uparrow_{2.85}$	44.18 $\uparrow_{1.43}$	53.76 $\downarrow_{0.54}$
	SparseVLM	70.90 $\downarrow_{7.23}$	43.39 $\downarrow_{8.81}$	78.70 $\uparrow_{17.21}$	16.25 $\downarrow_{20.66}$	50.10 $\uparrow_{7.79}$	51.87 $\downarrow_{2.43}$
Q-Tuning (Ours)		78.84 $\uparrow_{0.71}$	52.88 $\uparrow_{0.68}$	69.90 $\uparrow_{8.41}$	39.54 $\uparrow_{2.63}$	48.15 $\uparrow_{5.40}$	57.86 $\uparrow_{3.56}$
Full Dataset		77.78	52.20	65.13	38.05	44.08	55.45

C.3 ADDITIONAL RESULTS ON REASONING DATASETS

The detailed results of additional experiments on math reasoning are presented in Table 12, providing a comprehensive comparison across different pruning strategies.

Table 12: Evaluation of pruning strategies on GSM8K and MATH under 25% samples with 70% tokens, and 50% samples with 70% tokens settings. \uparrow and \downarrow respectively indicate improvements or degradations over the *Random-Random* baseline under the same sample and token keep ratio.

Sample Pruner	Token Pruner	LLaMA3-8B			Mistral-7B			SmolLM2-1.7B		
		GSM8K	MATH	Avg.	GSM8K	MATH	Avg.	GSM8K	MATH	Avg.
Zero-Shot		27.82	2.26	15.04	19.86	3.30	11.58	15.47	2.20	8.83
<i>25% Samples, 70% Tokens</i>										
Random	Random	25.09	2.20	13.65	24.11	1.68	12.89	13.80	2.22	8.01
	PPL	23.65 ^{↓1.44}	2.62 ^{↑0.42}	13.14 ^{↓0.51}	25.32 ^{↑1.21}	1.54 ^{↓0.14}	13.43 ^{↑0.54}	13.04 ^{↓0.76}	2.28 ^{↑0.06}	7.66 ^{↓0.35}
	FastV	16.91 ^{↓8.18}	2.16 ^{↓0.04}	9.53 ^{↓4.12}	16.07 ^{↓8.04}	1.60 ^{↓0.08}	8.84 ^{↓4.05}	12.89 ^{↓0.91}	1.94 ^{↓0.28}	7.41 ^{↓0.60}
	SparseVLM	22.97 ^{↓2.12}	4.72 ^{↑2.52}	13.85 ^{↑0.20}	19.26 ^{↓4.85}	4.58 ^{↑2.90}	11.92 ^{↓0.97}	13.19 ^{↓0.61}	3.48 ^{↑1.26}	8.34 ^{↑0.33}
Longest	Random	25.47 ^{↑0.38}	3.34 ^{↑1.14}	14.41 ^{↑0.76}	21.83 ^{↓2.28}	1.76 ^{↑0.08}	11.80 ^{↓1.09}	12.96 ^{↓0.84}	1.96 ^{↓0.26}	7.46 ^{↓0.55}
	PPL	24.87 ^{↓0.22}	3.22 ^{↑1.02}	14.04 ^{↑0.39}	23.65 ^{↓0.46}	2.22 ^{↑0.54}	12.94 ^{↑0.05}	14.33 ^{↑0.53}	1.60 ^{↓0.62}	7.96 ^{↓0.05}
	FastV	18.95 ^{↓6.14}	2.84 ^{↑0.64}	10.90 ^{↓2.75}	12.96 ^{↓1.15}	1.74 ^{↑0.06}	7.35 ^{↓5.54}	12.59 ^{↓1.21}	1.58 ^{↓0.64}	7.08 ^{↓0.93}
	SparseVLM	26.91 ^{↑1.82}	4.68 ^{↑2.48}	15.80 ^{↑2.15}	24.34 ^{↑0.23}	4.84 ^{↑3.16}	14.59 ^{↑1.70}	12.43 ^{↓1.37}	3.60 ^{↑1.38}	8.02 ^{↑0.01}
InfoBatch	Random	26.91 ^{↑1.82}	2.60 ^{↑0.40}	14.76 ^{↑1.11}	28.96 ^{↑4.85}	2.44 ^{↑0.76}	15.70 ^{↑2.81}	13.72 ^{↓0.08}	1.82 ^{↓0.40}	7.77 ^{↓0.24}
	PPL	25.93 ^{↑0.84}	2.48 ^{↑0.28}	14.20 ^{↑0.55}	31.46 ^{↑7.35}	2.18 ^{↑0.50}	16.82 ^{↑3.93}	14.86 ^{↑1.06}	1.90 ^{↓0.32}	8.38 ^{↑0.37}
	FastV	16.83 ^{↓8.26}	2.30 ^{↑0.10}	9.57 ^{↓4.08}	13.87 ^{↓10.24}	2.06 ^{↑0.38}	7.97 ^{↓4.92}	10.84 ^{↓2.96}	1.72 ^{↓0.50}	6.28 ^{↓1.73}
	SparseVLM	14.63 ^{↓10.46}	3.26 ^{↑1.06}	8.95 ^{↓4.70}	11.90 ^{↓12.21}	1.94 ^{↑0.26}	6.92 ^{↓5.97}	11.90 ^{↓1.90}	4.36^{↑2.14}	8.13 ^{↑0.12}
Entropy	Random	31.92 ^{↑6.83}	2.50 ^{↑0.30}	17.21 ^{↑3.56}	32.37 ^{↑8.26}	1.92 ^{↑0.24}	17.15 ^{↑4.26}	14.94 ^{↑1.14}	1.92 ^{↓0.30}	8.43 ^{↑0.42}
	PPL	33.13 ^{↑8.04}	2.86 ^{↑0.66}	18.00 ^{↑4.35}	30.17 ^{↑6.06}	1.96 ^{↑0.28}	16.07 ^{↑3.18}	14.94 ^{↑1.14}	1.96 ^{↓0.26}	8.45 ^{↑0.44}
	FastV	25.25 ^{↑0.16}	2.40 ^{↑0.20}	13.82 ^{↑0.17}	21.00 ^{↓3.11}	1.32 ^{↓0.36}	11.16 ^{↓1.73}	14.18 ^{↑0.38}	1.98 ^{↓0.24}	8.08 ^{↑0.07}
	SparseVLM	20.85 ^{↓4.24}	5.12^{↑2.92}	12.98 ^{↓0.67}	19.56 ^{↓4.55}	4.20 ^{↑2.52}	11.88 ^{↓1.01}	14.18 ^{↑0.38}	2.90 ^{↑0.68}	8.54 ^{↑0.53}
Q-Tuning (Ours)		37.23^{↑12.14}	4.86 ^{↑2.66}	21.04^{↑7.39}	42.99^{↑18.88}	5.08^{↑3.40}	24.56^{↑11.67}	22.90^{↑9.10}	3.64 ^{↑1.42}	13.27^{↑5.26}
<i>50% Samples, 70% Tokens</i>										
Random	Random	27.90	2.50	15.20	32.30	2.46	17.38	14.94	1.76	8.35
	PPL	27.45 ^{↓0.45}	2.50 ^{↑0.00}	14.97 ^{↓0.23}	31.99 ^{↓0.31}	2.04 ^{↓0.42}	17.02 ^{↓0.36}	16.15 ^{↑1.21}	1.80 ^{↑0.04}	8.97 ^{↑0.62}
	FastV	18.20 ^{↓9.70}	1.92 ^{↓0.58}	10.06 ^{↓5.14}	17.44 ^{↓14.86}	1.72 ^{↓0.74}	9.58 ^{↓7.80}	12.13 ^{↓2.81}	1.60 ^{↓0.16}	6.87 ^{↓1.48}
	SparseVLM	10.31 ^{↓17.59}	3.58 ^{↑1.08}	6.95 ^{↓8.25}	12.21 ^{↓20.09}	2.86 ^{↑0.40}	7.53 ^{↓9.85}	11.90 ^{↓3.04}	3.34 ^{↑1.58}	7.62 ^{↓0.73}
Longest	Random	26.23 ^{↓1.67}	2.76 ^{↑0.26}	14.50 ^{↓0.70}	27.82 ^{↓4.48}	2.20 ^{↓0.26}	15.01 ^{↓2.37}	15.69 ^{↑0.75}	1.92 ^{↑0.16}	8.81 ^{↑0.46}
	PPL	30.25 ^{↑2.35}	2.64 ^{↑0.14}	16.45 ^{↑1.25}	31.46 ^{↓0.84}	2.00 ^{↓0.46}	16.73 ^{↓0.65}	16.45 ^{↑1.51}	1.94 ^{↑0.18}	9.20 ^{↑0.85}
	FastV	18.42 ^{↓9.48}	2.28 ^{↓0.22}	10.35 ^{↓4.85}	18.65 ^{↓13.65}	1.82 ^{↓0.64}	10.24 ^{↓7.14}	12.43 ^{↓2.51}	1.94 ^{↑0.18}	7.19 ^{↓1.16}
	SparseVLM	18.95 ^{↓8.95}	4.84^{↑2.34}	11.90 ^{↓3.30}	19.48 ^{↓12.82}	4.74 ^{↑2.28}	12.11 ^{↓5.27}	11.07 ^{↓3.87}	3.90 ^{↑2.14}	7.48 ^{↓0.87}
InfoBatch	Random	29.42 ^{↑1.52}	2.92 ^{↑0.42}	16.17 ^{↑0.97}	35.03 ^{↑2.73}	2.50 ^{↑0.04}	18.76 ^{↑1.38}	16.00 ^{↑1.06}	1.78 ^{↑0.02}	8.89 ^{↑0.54}
	PPL	28.73 ^{↑0.83}	2.78 ^{↑0.28}	15.76 ^{↑0.56}	38.82 ^{↑6.52}	2.98 ^{↑0.52}	20.90 ^{↑3.52}	15.39 ^{↑0.45}	2.20 ^{↑0.44}	8.80 ^{↑0.45}
	FastV	19.33 ^{↓8.57}	1.76 ^{↓0.74}	10.55 ^{↓4.65}	18.04 ^{↓14.26}	1.94 ^{↓0.52}	9.99 ^{↓7.39}	12.05 ^{↓2.89}	1.86 ^{↑0.10}	6.96 ^{↓1.39}
	SparseVLM	9.10 ^{↓18.80}	3.12 ^{↑0.62}	6.11 ^{↓9.09}	12.05 ^{↓20.25}	2.32 ^{↓0.14}	7.19 ^{↓10.19}	10.84 ^{↓4.10}	3.96^{↑2.20}	7.40 ^{↓0.95}
Entropy	Random	32.60 ^{↑4.70}	2.14 ^{↓0.36}	17.37 ^{↑2.17}	38.67 ^{↑6.37}	2.10 ^{↓0.36}	20.38 ^{↑3.00}	18.88 ^{↑3.94}	2.18 ^{↑0.42}	10.53 ^{↑2.18}
	PPL	33.97 ^{↑6.07}	3.10 ^{↑0.60}	18.53 ^{↑3.33}	40.18 ^{↑7.88}	2.30 ^{↓0.16}	21.24 ^{↑3.86}	17.21 ^{↑2.27}	2.06 ^{↑0.30}	9.64 ^{↑1.29}
	FastV	24.94 ^{↓2.96}	2.26 ^{↓0.24}	13.60 ^{↓1.60}	21.53 ^{↓10.77}	1.56 ^{↓0.90}	11.55 ^{↓5.83}	16.22 ^{↑1.28}	1.78 ^{↑0.02}	9.00 ^{↑0.65}
	SparseVLM	10.31 ^{↓17.59}	4.56 ^{↑2.06}	7.44 ^{↑7.76}	9.40 ^{↓22.90}	4.40 ^{↑1.94}	6.90 ^{↓10.48}	13.95 ^{↓0.99}	3.68 ^{↑1.92}	8.82 ^{↑0.47}
Q-Tuning (Ours)		38.21^{↑10.31}	4.30 ^{↑1.80}	21.26^{↑6.06}	48.07^{↑15.77}	6.14^{↑3.68}	26.57^{↑9.19}	20.47^{↑5.53}	3.20 ^{↑1.44}	11.84^{↑3.49}
Full Dataset		32.90	3.02	17.96	42.08	3.08	22.58	16.53	2.10	9.31

C.4 ADDITIONAL RESULTS ON ADDITIONAL BASELINES

Tables 13 and 14 extend our evaluation to a broader set of pruning strategies under matched sample–token budgets. Table 13 varies sample ratios (12.5%, 25%) and token ratios (50%, 100%) under a fixed random token pruner, while Table 14 fixes the token ratio at 50% and varies the sample ratio (12.5%, 25%, 100%) under a random sample pruner.

Table 13: Evaluation of stronger task-relevant baselines under different sample ratios (12.5%, 25%) and token ratios (50%, 100%) with the same token pruner (*Random*), where \uparrow and \downarrow respectively denote improvements or degradations over the *Random* baseline.

Sample Pruner	Token Ratio	ARC-E	ARC-C	LLaMA2-7B				Mistral-7B					
		ARC-E	ARC-C	GSM8K	SQuAD	TriviaQA	Avg.	ARC-E	ARC-C	GSM8K	SQuAD	TriviaQA	Avg.
Zero-Shot		53.44	38.98	5.31	12.18	43.00	30.58	66.67	46.10	18.35	10.01	43.77	36.98
12.5% Samples													
Random	50%	59.25	41.02	8.11	12.75	48.75	33.98	70.55	48.14	22.74	19.57	52.63	42.73
Alpagasus	50%	58.91 ^{↓0.34}	41.69 ^{↑0.67}	6.67 ^{↓1.44}	12.81 ^{↑0.06}	46.28 ^{↓2.47}	33.27 ^{↓0.71}	72.84 ^{↑2.29}	48.47 ^{↑0.33}	20.55 ^{↓2.19}	20.41 ^{↑0.84}	52.21 ^{↓0.42}	42.90 ^{↑0.17}
	100%	58.91 ^{↓0.34}	41.36 ^{↑0.34}	7.58 ^{↓0.53}	13.45 ^{↑0.70}	46.96 ^{↓1.79}	33.65 ^{↓0.33}	70.72 ^{↑0.17}	48.81 ^{↑0.67}	23.35 ^{↑0.61}	24.30 ^{↑0.73}	52.71 ^{↑0.08}	43.98 ^{↑1.25}
Deita	50%	58.20 ^{↓1.05}	42.37 ^{↑1.35}	7.43 ^{↓0.68}	13.85 ^{↑1.10}	45.59 ^{↓3.16}	33.49 ^{↓0.49}	72.66 ^{↑2.11}	50.51^{↑2.37}	22.97 ^{↓0.23}	20.91 ^{↑1.34}	54.17 ^{↑1.54}	44.25 ^{↑1.52}
	100%	58.38 ^{↓0.87}	42.37 ^{↑1.35}	5.46 ^{↓2.65}	13.99 ^{↑1.24}	45.63 ^{↓3.12}	33.17 ^{↓0.81}	72.13 ^{↑1.58}	50.17 ^{↑2.03}	22.51 ^{↓0.23}	24.46 ^{↑0.89}	55.28 ^{↑2.65}	44.92 ^{↑2.19}
DS2	50%	58.73 ^{↓0.52}	42.37 ^{↑1.35}	5.99 ^{↓2.12}	15.05 ^{↑2.30}	45.51 ^{↓3.24}	33.53 ^{↓0.45}	71.25 ^{↑0.70}	47.46 ^{↓0.68}	21.23 ^{↓1.51}	21.57 ^{↑0.00}	51.23 ^{↓1.40}	42.55 ^{↓0.18}
	100%	59.08 ^{↓0.17}	41.36 ^{↑0.34}	6.67 ^{↓1.44}	15.14 ^{↑2.39}	45.69 ^{↓3.06}	33.59 ^{↓0.39}	72.13 ^{↑1.58}	49.83 ^{↑1.69}	23.35 ^{↑0.61}	23.99 ^{↑0.42}	52.85 ^{↑0.22}	44.43 ^{↑1.70}
LESS	50%	58.20 ^{↓1.05}	42.37 ^{↑1.35}	6.67 ^{↓1.44}	14.56 ^{↑1.81}	46.62 ^{↓2.13}	33.69 ^{↓0.29}	73.19 ^{↑2.64}	49.49 ^{↑1.35}	22.90 ^{↓0.16}	19.61 ^{↑0.04}	53.14 ^{↑0.51}	43.67 ^{↑0.94}
	100%	58.38 ^{↓0.87}	43.39^{↑2.37}	6.82 ^{↓1.29}	16.16 ^{↑3.41}	46.72 ^{↓2.03}	34.30 ^{↑0.32}	73.90^{↑3.38}	48.47 ^{↑0.33}	23.12 ^{↑0.38}	21.67 ^{↑2.10}	56.00 ^{↑3.37}	44.63 ^{↑1.90}
Q-Tuning (Ours)	50%	64.20^{↑4.95}	42.03 ^{↑1.01}	10.54^{↑2.43}	18.79^{↑6.84}	53.12^{↑3.37}	37.74^{↑3.76}	71.60 ^{↑1.05}	48.14 ^{↑0.00}	29.34^{↑6.60}	27.75^{↑8.18}	57.78^{↑5.18}	46.92^{↑4.19}
25% Samples													
Random	50%	60.32	41.69	5.76	13.43	48.41	33.92	70.19	46.10	20.62	24.07	53.74	42.95
Alpagasus	50%	59.61 ^{↓0.71}	42.71 ^{↑1.02}	5.91 ^{↑0.15}	14.12 ^{↑0.69}	47.30 ^{↓1.11}	33.93 ^{↑0.01}	70.37 ^{↑0.18}	47.80 ^{↑1.70}	20.85 ^{↑0.23}	22.38 ^{↓1.69}	52.87 ^{↓0.87}	42.78 ^{↓0.17}
	100%	60.49 ^{↑0.87}	43.39 ^{↑1.70}	7.51 ^{↑1.75}	13.16 ^{↓0.27}	48.41 ^{↑0.00}	34.59 ^{↑0.67}	71.60 ^{↑1.41}	50.17^{↑4.07}	23.35 ^{↑2.73}	22.38 ^{↓1.69}	55.73 ^{↑1.99}	44.65 ^{↑1.70}
Deita	50%	59.08 ^{↓1.24}	41.69 ^{↑0.00}	7.66 ^{↑1.90}	13.23 ^{↓0.20}	47.55 ^{↓0.86}	33.84 ^{↓0.08}	72.13 ^{↑1.94}	47.80 ^{↑1.70}	21.00 ^{↑0.38}	23.16 ^{↓0.91}	53.97 ^{↑0.23}	43.61 ^{↑0.66}
	100%	58.73 ^{↓1.59}	42.71 ^{↑1.02}	4.92 ^{↓0.84}	14.88 ^{↑1.45}	48.24 ^{↓0.17}	33.90 ^{↓0.02}	72.66^{↑2.47}	47.80 ^{↑1.70}	22.29 ^{↑1.67}	25.87 ^{↑1.80}	54.40 ^{↑0.66}	44.60 ^{↑1.65}
DS2	50%	59.44 ^{↓0.88}	42.03 ^{↑0.34}	6.44 ^{↑0.68}	15.59 ^{↑2.16}	47.34 ^{↓1.07}	34.17 ^{↑0.25}	70.72 ^{↑0.53}	48.47 ^{↑2.37}	22.74 ^{↑2.12}	21.15 ^{↓2.92}	53.49 ^{↓0.25}	43.32 ^{↑0.37}
	100%	59.26 ^{↓1.06}	41.36 ^{↓0.33}	7.05 ^{↑1.29}	15.20 ^{↑1.77}	48.76 ^{↑0.35}	34.33 ^{↑0.41}	71.43 ^{↑1.24}	48.47 ^{↑2.37}	22.44 ^{↑1.82}	21.72 ^{↓2.35}	53.76 ^{↑0.02}	43.57 ^{↑0.62}
LESS	50%	59.44 ^{↓0.88}	43.39^{↑1.70}	7.05 ^{↑1.29}	15.89 ^{↑2.46}	46.60 ^{↓1.81}	34.47 ^{↑0.55}	71.78 ^{↑1.59}	47.80 ^{↑1.70}	22.14 ^{↑1.52}	24.51 ^{↑0.44}	53.04 ^{↓0.70}	43.85 ^{↑0.90}
	100%	59.79 ^{↓0.53}	41.02 ^{↓0.67}	6.37 ^{↑0.61}	14.44 ^{↑1.01}	48.44 ^{↑0.03}	34.01 ^{↑0.69}	71.60 ^{↑1.41}	49.49 ^{↑3.39}	21.68 ^{↑1.06}	30.85 ^{↑6.78}	53.40 ^{↓0.34}	45.41 ^{↑2.46}
Q-Tuning (Ours)	50%	63.14^{↑2.82}	42.03 ^{↑0.34}	8.87^{↑3.11}	16.76^{↑3.33}	51.52^{↑3.11}	36.47^{↑2.55}	71.78 ^{↑1.59}	47.12 ^{↑1.02}	26.08^{↑5.46}	32.79^{↑8.72}	56.17^{↑2.43}	46.79^{↑3.84}
Full Dataset		61.55	42.37	8.64	13.80	50.45	35.36	71.25	45.76	26.68	31.81	53.67	45.84

Table 14: Evaluation of stronger task-relevant baselines under different sample ratios (12.5%, 25%, 100%) and fixed token ratio (50%) with the same sample pruner (*Random*), where \uparrow and \downarrow respectively denote improvements or degradations over the *Random* baseline.

Token Pruner	Sample Ratio	ARC-E	ARC-C	LLaMA2-7B				Mistral-7B					
		ARC-E	ARC-C	GSM8K	SQuAD	TriviaQA	Avg.	ARC-E	ARC-C	GSM8K	SQuAD	TriviaQA	Avg.
Zero-Shot		53.44	38.98	5.31	12.18	43.00	30.58	66.67	46.10	18.35	10.01	43.77	36.98
50% Tokens													
Random	12.5%	59.25	41.02	8.11	12.75	48.75	33.98	70.55	48.14	22.74	19.57	52.63	42.73
	25%	60.32 ^{↑1.07}	41.69 ^{↑0.67}	5.76 ^{↓2.35}	13.43 ^{↑0.68}	48.41 ^{↓0.34}	33.92 ^{↓0.06}	70.19 ^{↓0.36}	46.10 ^{↓2.04}	20.62 ^{↓2.12}	24.07 ^{↑0.50}	53.74 ^{↑1.11}	42.95 ^{↑0.22}
Rho-1	12.5%	53.97 ^{↓5.28}	40.34 ^{↓0.68}	6.67 ^{↓1.44}	13.27 ^{↑0.52}	46.62 ^{↓2.13}	32.17 ^{↓1.81}	71.60 ^{↑1.05}	47.46 ^{↓0.68}	20.62 ^{↓2.12}	23.74 ^{↑0.17}	53.71 ^{↑1.08}	43.43 ^{↑0.70}
	25%	59.08 ^{↓0.17}	41.02 ^{↑0.00}	6.67 ^{↓1.44}	13.99 ^{↑1.24}	47.19 ^{↓1.56}	33.59 ^{↓0.39}	70.37 ^{↑0.18}	46.78 ^{↓1.36}	21.76 ^{↓0.98}	20.85 ^{↑1.28}	52.22 ^{↓0.41}	42.40 ^{↓0.33}
	100%	60.67 ^{↑1.42}	41.36 ^{↑0.34}	7.28 ^{↓0.83}	14.97 ^{↑2.22}	49.82 ^{↑1.07}	34.82 ^{↑0.84}	63.67 ^{↓4.88}	43.39 ^{↓4.75}	11.68 ^{↓11.06}	15.46 ^{↓4.41}	31.84 ^{↓20.79}	33.21 ^{↓8.52}
TokenCleaning	12.5%	58.73 ^{↓0.52}	43.73^{↑2.71}	7.43 ^{↓0.68}	12.84 ^{↑0.09}	46.37 ^{↓2.38}	33.82 ^{↓0.16}	71.60 ^{↑1.05}	49.49^{↑3.88}	21.00 ^{↓1.74}	22.73 ^{↑3.16}	53.03 ^{↑0.40}	43.57 ^{↑0.84}
	25%	59.96 ^{↑0.71}	42.37 ^{↑1.35}	5.99 ^{↓2.12}	13.62 ^{↑0.87}	46.14 ^{↓2.61}	33.62 ^{↓0.36}	70.19 ^{↓0.36}	47.80 ^{↓0.34}	19.94 ^{↓2.80}	24.17 ^{↑4.60}	52.02 ^{↓0.61}	42.82 ^{↑0.09}
	100%	61.20 ^{↑1.95}	41.02 ^{↑0.00}	7.73 ^{↓0.38}	15.41 ^{↑2.66}	48.29 ^{↓0.46}	34.73 ^{↑0.75}	74.07^{↑3.52}	46.44 ^{↓1.70}	24.79 ^{↑2.05}	33.26 ^{↑13.69}	54.32 ^{↑1.69}	46.58 ^{↑3.85}
Q-Tuning (Ours)	12.5%	64.20^{↑4.95}	42.03 ^{↑1.01}	10.54^{↑2.43}	18.79^{↑6.84}	53.12^{↑3.37}	37.74^{↑3.76}	71.60 ^{↑1.05}	48.14 ^{↑0.00}	29.34^{↑6.60}	27.75^{↑8.18}	57.78^{↑5.18}	46.92^{↑4.19}
	25%	63.14 ^{↑3.89}	42.03 ^{↑1.01}	8.87 ^{↑3.11}	16.76 ^{↑3.33}	51.52 ^{↑3.11}	36.47 ^{↑2.49}	71.78 ^{↑1.23}	47.12 ^{↓1.02}	26.08 ^{↑3.34}	32.79^{↑13.22}	56.17 ^{↑3.54}	46.79 ^{↑4.06}
Full Dataset		61.55	42.37	8.64	13.80	50.45	35.36	71.25	45.76	26.68	31.81	53.67	45.84

C.5 FURTHER ABLATION STUDY

C.5.1 ABLATION WITH DIFFERENT HYPERPARAMETERS

Sensitivity of batch size in dynamic sample pruning.

Figure 9 presents whole ablation results on Mistral-7B, evaluating the effects of batch size (8/16/32) and neighbor awareness λ (0–1.0) under multiple keep-ratio settings across ARC-E, ARC-C, GSM8K, SQuAD, and TriviaQA, as well as the average over the five benchmarks.

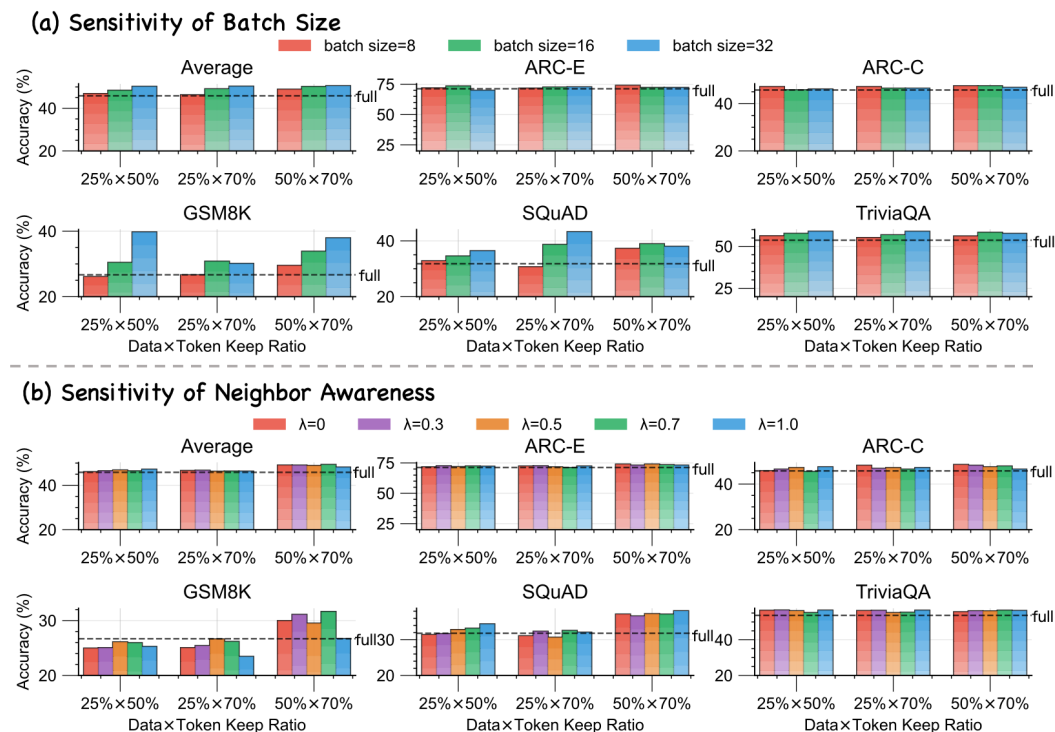


Figure 9: Effect of varying (a) batch size (8, 16, 32) and (b) neighbor awareness λ (0, 0.3, 0.5, 0.7, 1.0) for Mistral-7B under three data–token keep ratio configurations (25% \times 50%, 25% \times 70%, 50% \times 70%). Dashed lines marked “full” denote models trained on the full data without pruning.

C.5.2 ADDITIONAL ABLATION ACROSS DIFFERENT QUADRANTS

Table 15 presents additional quadrant-wise ablation results on the Wizard dataset, evaluated on LLaMA2-7B and Mistral-7B, examining how pruning Q1–Q4 at both the sample and token levels impacts downstream performance.

Table 15: Ablation on four quadrants under different sample and token ratios, where \checkmark indicates that the corresponding quadrant is pruned.

Sample pruning Q1	Token pruning Q3	Token pruning Q2	Token pruning Q4	LLaMA2-7B					Mistral-7B						
				ARC-E	ARC-C	GSM8K	SQuAD	TriviaQA	Avg.	ARC-E	ARC-C	GSM8K	SQuAD	TriviaQA	Avg.
12.5% Samples, 50% Tokens															
	\checkmark	\checkmark		60.14	41.02	8.04	13.11	50.23	34.51	71.43	47.80	28.13	31.09	54.27	46.54
\checkmark	\checkmark	\checkmark	\checkmark	64.02	41.02	9.86	19.26	52.00	37.23	71.60	48.47	28.89	26.46	57.88	46.66
\checkmark	\checkmark	\checkmark		64.20	42.03	10.54	18.79	53.12	37.74	71.60	48.14	29.34	27.75	57.78	46.92
25% Samples, 50% Tokens															
	\checkmark	\checkmark		59.61	41.02	7.51	13.74	50.25	34.43	71.25	46.10	29.04	31.55	54.84	46.56
\checkmark	\checkmark	\checkmark	\checkmark	61.02	40.68	9.10	18.66	51.68	36.23	70.72	47.12	28.58	31.85	56.33	46.92
\checkmark	\checkmark	\checkmark		63.14	42.03	8.87	16.76	51.52	36.47	71.78	47.12	26.08	32.79	56.17	46.79

D PSEUDOCODE OF THE PROPOSED METHOD

To facilitate clarity and reproducibility, we summarize the proposed Q-Tuning algorithm in Algorithm 1. The procedure unfolds within a single training iteration and consists of two tightly coupled stages: (i) *sample pruning*, where instances are dynamically selected based on their position in the error–uncertainty (EU) plane through an efficient bisection-based search of quantile thresholds, and (ii) *token pruning*, where retained samples undergo finer-grained filtering to preserve only the most informative subset of tokens.

Algorithm 1 Q-Tuning: dynamic data pruning in one iteration

```

1: Input: Mini-batch  $\mathcal{B}_t$ , model  $f_{\theta_t}$ , retention ratios  $r_{\text{sample}}, r_{\text{token}}$ , smoothing  $\lambda$ .
2: Output: Pruned mini-batch  $\tilde{\mathcal{B}}_t$ .
3: // Stage 1: Sample Pruning via EU Plane
4: Compute  $(\text{PPL}(x, y; f_{\theta}), \text{Ent}(x, y; f_{\theta}))$  for each  $x \in \mathcal{B}_t$ .
5: Initialize ranges  $\alpha_{\text{low}} = 0, \alpha_{\text{high}} = 0.49, \beta_{\text{low}} = 0, \beta_{\text{high}} = 0.49$ .
6: for  $k = 1$  to  $K_{\text{max}}$  do ▷ Bisection iterations on both axes
7:    $\alpha \leftarrow (\alpha_{\text{low}} + \alpha_{\text{high}})/2, \beta \leftarrow (\beta_{\text{low}} + \beta_{\text{high}})/2$ .
8:   Derive thresholds  $\text{ppl}_{\text{hi}} = Q_{1-\alpha}(\text{PPL}), \text{ppl}_{\text{lo}} = Q_{\alpha}(\text{PPL}), \text{ent}_{\text{lo}} = Q_{\beta}(\text{Ent}), \text{ent}_{\text{hi}} =$ 
9:    $Q_{1-\beta}(\text{Ent})$ .
10:  Partition samples into quadrants  $Q_1$ – $Q_4$ .
11:   $r \leftarrow \frac{|Q_2| + |Q_4|}{|\mathcal{B}_t|}$ .
12:  if  $r < r_{\text{sample}}$  then ▷ Too few kept, relax thresholds
13:     $\alpha_{\text{low}} \leftarrow \alpha, \beta_{\text{low}} \leftarrow \beta$ 
14:  else ▷ Too many kept, tighten thresholds
15:     $\alpha_{\text{high}} \leftarrow \alpha, \beta_{\text{high}} \leftarrow \beta$ 
16:  end if
17: end for
18:  $\mathcal{B}'_t \leftarrow Q_2 \cup Q_4$ .
19: if  $|\mathcal{B}'_t| < r_{\text{sample}} \cdot |\mathcal{B}_t|$  then ▷ Ensure target sample ratio
20:   Compute supp-score for each  $x \in Q_1 \cup Q_3$ :  $\text{supp}(x) = \max\{\hat{\text{PPL}}(x, y; f_{\theta}) -$ 
21:    $\hat{\text{Ent}}(x, y; f_{\theta}), \hat{\text{Ent}}(x, y; f_{\theta}) - \hat{\text{PPL}}(x, y; f_{\theta})\}$ , where  $\hat{\cdot}$  denotes min–max normalization.
22:   Select top-scoring samples to augment  $\mathcal{B}'_t$  until target size is met.
23: end if
24: // Stage 2: Token Pruning
25:  $\tilde{\mathcal{B}}_t \leftarrow \emptyset$ .
26: for each  $x \in \mathcal{B}'_t$  do
27:   if  $x \in Q_4$  then ▷ Calibration samples intact
28:     Keep all tokens of  $x$ .
29:   else if  $x \in Q_2$  then
30:     For each token  $i \in T(x)$ , compute
31:     
$$s_i(x, y; f_{\theta}) = (1 - \lambda) \text{PPL}_i(x, y; f_{\theta}) + \lambda [\text{PPL}_{i-1}(x, y; f_{\theta}) + \text{PPL}_{i+1}(x, y; f_{\theta})].$$

32:     Keep top- $r_{\text{token}}$  fraction of tokens ranked by  $s_i$ .
33:   end if
34:    $\tilde{\mathcal{B}}_t \leftarrow \tilde{\mathcal{B}}_t \cup \{m(x) \odot x\}$ .
35: end for
36: return  $\tilde{\mathcal{B}}_t$ .

```

1350 E FURTHER CASE STUDY

1351

1352 E.1 EXAMPLE SAMPLES FROM EACH QUADRANT

1353

1354 We report representative samples from each quadrant (Q1–Q4) at the beginning of training, including
 1355 their perplexity (PPL), entropy, and the corresponding instruction–output pairs.

1356

1357	PPL	Entropy	Example (Instruction / Output)
1358	Q1 Examples		
1359	8.33	2.12	Instruction: What are the three primary colors? Output: The three primary colors are red, blue, and yellow.
1360			
1361			
1362	8.04	1.73	Instruction: How can we reduce air pollution? Output: There are a number of ways to reduce air pollution, such as shifting to renewable energy sources...
1363			
1364			
1365			
1366	91.24	1.91	Instruction: Explain why the fraction 4/16 is equivalent to 1/4. Output: The fraction 4/16 is equivalent to 1/4 because both numerator and denominator are divisible by 4.
1367			
1368			
1369			
1370			
1371	Q2 Examples		
1372	8.65	1.54	Instruction: Write a short story in third person narration about a protagonist making an important career decision. Output: John was at a crossroads in his life. He had just graduate...
1373			
1374			
1375	11.02	1.48	Instruction: Write a short story in third person narration about a protagonist making an important career decision. Output: John was at a crossroads in his life. He had just graduate...
1376			
1377			
1378			
1379	15.74	1.50	Instruction: Identify the odd one out.\nTwitter, Instagram, Telegram Output: Telegram
1380			
1381			
1382	Q3 Examples		
1383	4.80	1.45	Instruction: Describe the structure of an atom. Output: An atom is made up of a nucleus containing protons and neutrons, surrounded by electrons...
1384			
1385			
1386	2.41	1.09	Instruction: How can we reduce air pollution? Output: There are a number of ways to reduce air pollution, such as shifting to renewable energy sources...
1387			
1388			
1389			
1390	4.73	1.49	Instruction: Identify the odd one out (Twitter, Instagram, Telegram). Output: Telegram
1391			
1392			
1393			
1394	Q4 Examples		
1395	6.03	1.83	Instruction: Give three tips for staying healthy. Output: 1. Eat a balanced diet... 2. Exercise regularly...
1396			
1397	8.04	1.71	Instruction: What are the three primary colors? Output: The three primary colors are red, blue, and yellow.
1398			
1399			
1400	5.37	1.79	Instruction: Write a short story in third person narration about a protagonist making an important career decision. Output: John was at a crossroads in his life. He had just graduate...
1401			
1402			
1403			

1404 E.2 PRUNED CASE
1405

1406 We present representative token-pruning cases from MATHINSTRUCT using the MISTRAL-7B model.
1407 For each example, we show the model’s original output and highlight in red the fine-grained
1408 segments that are *absent* in the corresponding pruned_output. This visualization makes the
1409 pruning behavior transparent at the character/substring level and helps diagnose which semantic
1410 components are preferentially removed (e.g., boilerplate, redundant explanations, or unstable numeric
1411 derivations).

1412 Unless otherwise stated, we keep the original formatting and line breaks of the model outputs to
1413 preserve exact alignment. The highlighting is produced by an LCS-based diff between output and
1414 pruned_output, ensuring that only truly deleted spans are marked while retained text remains
1415 uncolored. We do *not* introduce any additional edits beyond necessary \LaTeX escaping.

1416 These qualitative cases complement the aggregate statistics reported in the main text by revealing
1417 typical pruning patterns on mathematical reasoning tasks (definitions and restatements vs. compu-
1418 tational steps), and by surfacing edge cases where pruning removes numerically relevant content.
1419 Together, they provide an interpretable view of how token-level pruning affects the fidelity and brevity
1420 of generated solutions.

1421
1422 SAMPLE 0

1423 **Instruction** sandy is younger than molly by 20 years . if the ratio of their ages is 7 : 9 , how old
1424 is sandy ? Let’s write a program.

1425
1426 **Output**

```
1427 n0 = 20.0
1428 n1 = 7.0
1429 n2 = 9.0
1430 t0 = n2 - n1\n
1431 t1 = n1 / t0
1432 answer =n0 * t1
1433 print(answer)
```

1434
1435 SAMPLE 1

1436
1437 **Instruction** A manufacturing process produces 5000 units of a particular product. The probability
1438 of a certain type of defect occurring in a unit is 0.02. What is the probability that exactly 100
1439 units will have this defect? **Output (deletions highlighted)**

1440 To solve this problem, we can use the binomial probability
1441 formula, which is:

$$1442 P(X = k) = (nCk) * (p^k) * (1-p)^{(n-k)}$$

1443 where:

- 1444 - $P(X = k)$ is the probability of exactly k successes (defective
- 1445 units in this case) in n trials (total units produced)
- 1446 - nCk (read as " n choose k ") is the number of combinations of
- 1447 choosing k successes from n trials, calculated as $n! / (k! * (n-k)!)$
- 1448 - p is the probability of success (defect occurring) in a
- 1449 single trial
- 1450 - $1-p$ is the probability of failure (no defect occurring) in a
- 1451 single trial
- 1452 - n is the total number of trials (units produced)
- 1453 - k is the number of successes (defective units) we want to
- 1454 find the probability for
- 1455 In this problem, we have:
- 1456 - $n = 5000$ (total units produced)
- 1457 - $k = 100$ (defective units)

1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511

```
- p = 0.02 (probability of a defect occurring in a unit)
Now, we can plug these values into the binomial probability
formula:
P(X=100)=(5000C100)· (0.02)100· (0.98)4900)
First, let's calculate 5000C100:
5000C100 = 5000! / (100! * 4900!)
= 2.031759e+134 / (9.332622e+157 * 1.929713e+126)
= 2.031759e+134 / (1.801713e+284)
= 1.127271e-150
Now, let's calculate 0.02100 and 0.984900:
0.02100 = 1.267651e-70
0.984900 = 2.372472e-44
Finally, let's plug these values back into the formula:
P(X = 100) = 1.127271e-150 * 1.267651e-70 * 2.372472e-44
= 3.394e-264
The probability that exactly 100 units will have this defect is
approximately 3.394e-264, which is an extremely small
probability.
```

1512 F USE OF LLMs

1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565

This work did not involve the use of large language models (LLMs) at any stage. The design of experiments, data analysis, and manuscript preparation were conducted entirely by the authors through conventional computational methods and human expertise, without reliance on automated text generation or model-driven reasoning.