

TRAIN ON VALIDATION (ToV): FAST DATA SELECTION WITH APPLICATIONS TO FINE-TUNING

Ayush Jain*
Google Research
ayush.jain.2508@gmail.com

Eren Sasoglu*
Encord
eren.sasoglu@gmail.com

Andrea Montanari†
Granica Computing Inc.– (granica.ai)
andrea.montanari@granica.ai

ABSTRACT

State-of-the-art machine learning often follows a two-stage process: (i) pre-training on large, general-purpose datasets; (ii) fine-tuning on task-specific data. In fine-tuning, selecting training examples that closely reflect the target distribution is crucial. However, it is often the case that only a few samples are available from the target distribution. Existing data selection methods treat these target samples as a validation set and estimate the effect of adding or removing a single sample from the training pool by performing inference on the validation set.

We propose a simpler and faster alternative that inverts the usual role of train and validation: we perform inference on the training pool before and after fine-tuning on the validation set. We then select samples whose predictions change the most. Our key insight is that the training samples most affected by fine-tuning on a small validation set tend to be the most beneficial for reducing test loss on the target distribution. Experiments on instruction tuning and named entity recognition tasks show that, in most cases, our method achieves lower test log-loss than state-of-the-art approaches. We support our findings with theoretical analysis.

1 INTRODUCTION

While large language models (LLMs) are pretrained on internet-scale datasets, their downstream performance can be heavily dependent on the instruction-tuning stage in which they are fine-tuned on instruction/output pairs (Ouyang et al., 2022; Zhou et al., 2024; Longpre et al., 2023; Chung et al., 2024). These datasets are significantly smaller and are often gathered by using multiple heterogeneous sources. Instruction tuning becomes even more difficult when targeting a specialized use case (Wang et al., 2023). More generally, scarcity of domain-specific data is a ubiquitous challenge when fine-tuning foundation models.

This paper presents an easy-to-implement and low-complexity method for selecting a training dataset of prescribed size from heterogeneous sources to maximize the test time performance on the target distribution. Our method is motivated by the theory of influence functions (van der Vaart, 2000) yet avoids the computational burden of computing influence functions. Unlike influence-function–based data selection methods, ToV avoids per-example gradients, Hessian–vector products, and backpropagation through training dynamics. It instead estimates example utility using only forward loss evaluations before and after a small gradient step on the validation data. We validate this approach on two token-based learning tasks, instruction tuning and named entity recognition (NER), and show that in most cases it outperforms state-of-the-art data selection baselines. To illustrate its broad applicability, we show that it yields interesting results even for a simple logistic regression example (see Appendix D).

*Work done while at Granica Computing Inc.

†Department of Statistics and Department of Mathematics, Stanford University

To formalize the problem, assume access to two datasets: a small dataset from the target distribution \mathbb{P} on \mathcal{Z} and a larger one from possibly heterogeneous data sources. We refer to the dataset from the target as the ‘validation set’ $\mathbf{Z}^{\text{val}} := (z_1^{\text{val}}, \dots, z_{m_{\text{val}}}^{\text{val}})$ where z_i^{val} are i.i.d. samples from the target distribution \mathbb{P} , and to the larger heterogeneous dataset as ‘training pool’ $\mathbf{X} = (x_1, \dots, x_N)$, where $x_i \in \mathcal{Z}$. In general the distribution of the training pool differs from \mathbb{P} . Our goal is to minimize the test error on the target distribution with respect to the model parameters $\theta \in \mathbb{R}^p$:

$$R(\theta) := \mathbb{E}[\ell(\theta, z)], \quad (1)$$

where $\ell : \mathbb{R}^p \times \mathcal{Z} \rightarrow \mathbb{R}$ is a loss function. A separate target-distribution test set \mathbf{Z}^{st} (separate from \mathbf{Z}^{val}) is used to estimate $R(\theta)$ after fine-tuning.

We aim to achieve this by training (or fine-tuning) the model on a subset $S \subseteq [N]$ of the training pool, e.g. running stochastic gradient descent (SGD) with respect to the empirical risk:

$$\widehat{R}_S(\theta) := \frac{1}{|S|} \sum_{i \in S} \ell(\theta, x_i). \quad (2)$$

Let $\widehat{\theta}_S$ be the outcome of running SGD (or any specific training algorithm) on $\widehat{R}_S(\theta)$. We want to select the subset S (given a constraint on its size $|S|$) so that $\widehat{\theta}_S$ achieves a small test loss on the target distribution, i.e. as to minimize $R(\widehat{\theta}_S)$.

1.1 TRAIN ON VALIDATION: MOTIVATION AND ALGORITHM

To select the most helpful examples at model θ , we might score training examples by the decrease in validation loss induced by a single gradient step with respect to that example, then select those with the highest scores. Computing these scores directly requires $N+1$ full evaluations over the validation set. We derive an efficient approximation to these scores.

Consider a single gradient step with respect to a training example x :

$$\theta_x = \theta - \eta \nabla \ell(\theta, x). \quad (3)$$

The corresponding change in loss for a validation example z , $\ell(\theta, z) - \ell(\theta_x, z)$, can be approximated by a first-order Taylor expansion:

$$\ell(\theta, z) - \ell(\theta_x, z) \approx -\langle \nabla \ell(\theta, z), \theta_x - \theta \rangle = \eta \langle \nabla \ell(\theta, z), \nabla \ell(\theta, x) \rangle, \quad (4)$$

where the last step follows from Eq. (3). Pruthi et al. (2020) approximate the scores by computing gradients for each training and validation example and taking their dot products; Xia et al. (2024) extend this to token-based learning. Our method diverges from these approaches: it requires no per-example gradients.

Note the right-hand side is symmetric in x and z . In other words, the decrease in loss on z from a step on x is mirrored by the decrease in loss on x from a step on z . Our method exploits this train-validation symmetry. The change in overall validation loss for a single gradient step with respect to x is:

$$\frac{1}{m_{\text{val}}} \sum_{i=1}^{m_{\text{val}}} \left(\ell(\theta, z_i) - \ell(\theta_x, z_i) \right) \approx \frac{1}{m_{\text{val}}} \sum_{i=1}^{m_{\text{val}}} \eta \langle \nabla \ell(\theta, z_i), \nabla \ell(\theta, x) \rangle. \quad (5)$$

On the other hand, performing a batch gradient step at θ with respect to the validation set gives $\theta_{\mathbf{Z}^{\text{val}}} = \theta - \eta \frac{1}{m_{\text{val}}} \sum_{i=1}^{m_{\text{val}}} \nabla \ell(\theta, z_i)$. Combining this with Eq. (5), we get

$$\frac{1}{m_{\text{val}}} \sum_{i=1}^{m_{\text{val}}} \left(\ell(\theta, z_i) - \ell(\theta_x, z_i) \right) \approx \langle \theta - \theta_{\mathbf{Z}^{\text{val}}}, \nabla \ell(\theta, x) \rangle \approx \ell(\theta, x) - \ell(\theta_{\mathbf{Z}^{\text{val}}}, x). \quad (6)$$

In other words, the change in average validation loss from training on x can be approximated by the change in loss on x after training on the validation set \mathbf{Z}^{val} .

The approximation in Eq. (6) follows directly from a first-order Taylor expansion and therefore relies only on the update induced by the validation step being small and the loss landscape being locally

Algorithm 1 ToV Scoring Algorithm: *Interleaved ToV*.

-
- 1: **Input:** Pretrained model θ_0 , validation set Z^{val} , training pool $X = (\mathbf{x}_i : i \in [N])$, epochs L ,
 - 2: selected data count n , learning-rate schedule $\{\eta_k\}_{k=1}^L$, base model count m , $\varepsilon \in [0, 1]$
 - 3: **Output:** Set of examples $S \subset [N]$ of size n
 - 4: Sample base subset $U \subseteq [N]$ of size m randomly; define $X_U = (\mathbf{x}_i : i \in U)$
 - 5: Initialize model: $\hat{\theta}_0^{\text{bas}} \leftarrow \theta_0$; set scores $\phi_i \leftarrow 0$ for all $i \in [N] \setminus U$
 - 6: **for** $k = 1$ to L **do**
 - 7: Train $\hat{\theta}_{k-1}^{\text{bas}}$ on X_U for one epoch with learning rate η_k to obtain $\hat{\theta}_k^{\text{bas}}$
 - 8: Train $\hat{\theta}_k^{\text{bas}}$ for one epoch on Z^{val} with a learning rate $\varepsilon\eta_k$ to obtain $\hat{\theta}_k^{\text{val}}$
 - 9: **for each** $i \in [N] \setminus U$ **do**
 - 10: $\phi_i^{(k)} \leftarrow F(\ell(\hat{\theta}_k^{\text{val}}; \mathbf{x}_i) - \ell(\hat{\theta}_k^{\text{bas}}; \mathbf{x}_i))$ (see Section 2.1 for the definition of F)
 - 11: $\phi_i \leftarrow \phi_i + \phi_i^{(k)} / L$
 - 12: **end for**
 - 13: **end for**
 - 14: Return set $S \subseteq [N] \setminus U$ of size n on the basis of scores ϕ_i (see text)
-

smooth. If the update is too large or the landscape is highly non-linear, the approximation may degrade. In practice, we enforce this regime by using a validation learning rate that is smaller than the base learning rate. No other assumptions, such as independence between training and validation samples, are required.

Our main objective is to evaluate the left-hand side of Eq. (6) for all \mathbf{x} in the training set. The right-hand side provides a far more efficient route: (i) Compute the loss $\ell(\theta, \mathbf{x})$ for all training examples; (ii) fine-tune θ on the validation set to obtain $\theta_{Z^{\text{val}}}$; (iii) re-evaluate the new loss $\ell(\theta_{Z^{\text{val}}}, \mathbf{x})$ at each training sample \mathbf{x} , and approximate the effect of training on \mathbf{x} by computing the difference with the loss at point (i).

This requires one epoch of training on the validation set and two evaluations over the training pool, as opposed to N evaluations of the validation loss as suggested by a direct evaluation of the left-hand side of Eq. (6), and it does not require access to per-example gradients.

In the next sections, we use this idea to obtain a selection algorithm that alternates training on a subset of the training set and on the validation set. A specific implementation, which we refer to as *Interleaved ToV* (Method A), is given in Algorithm 1; a slightly different implementation, *Parallel ToV* (Method B), is given in Algorithm 2. In Method A, we start with a small random subset $U \subset [N]$ of the training pool. We train on U for L epochs, resulting in models $\hat{\theta}_1^{\text{bas}}, \dots, \hat{\theta}_L^{\text{bas}}$. For each epoch $k \in [L]$ we fine-tune $\hat{\theta}_k^{\text{bas}}$ for one epoch on the validation set, resulting in models $\hat{\theta}_k^{\text{val}}$. For each epoch, every remaining training example \mathbf{x}_i with $i \in [N] \setminus U$ is scored by the change in its loss between $\hat{\theta}_k^{\text{bas}}$ and $\hat{\theta}_k^{\text{val}}$, and scores are averaged across epochs.

The names reflect the mechanics: *Interleaved ToV* interleaves base-set training with a single validation update at each checkpoint, whereas *Parallel ToV* runs two training trajectories in parallel, only one of which is updated on the validation set (see Section 3).

After computing scores ϕ_i as in Algorithm 1, we select S using one of two strategies: (i) choose the n examples with the largest ϕ_i ; (ii) choose half from the highest-scoring examples and the other half uniformly at random from U to increase diversity.

Intuitively, large ϕ_i means that a small amount of training on the target distribution produces a large change in the model output at \mathbf{x}_i . Our working assumption, motivated by the heuristics above and formalized in Section 3, is that the converse also holds: training on \mathbf{x}_i will produce a large change in model output on the target distribution. Hence the scores ϕ_i can be used to select ‘important’ samples for the target.

An adaptation for token-based learning is described in Section 2, along with empirical results. Section 3 provides a mathematical justification that formalizes the argument above.

1.2 RELATED WORK

Our work relates to data selection and data attribution. The impact of a single example on the validation error can be approximated by a first-order Taylor expansion. This idea results in data selection methods based on influence functions (Wang et al., 2018; 2020; Ai et al., 2021; Kolossov et al., 2024). Classical influence functions estimate the effect of a single example on the empirical risk minimizer. Most closely related to our work are Pruthi et al. (2020); Bae et al. (2024); Xia et al. (2024), which instead estimate the influence of an example on the training dynamics. In particular, Bae et al. (2024) shows how to approximately propagate gradient changes at k -th epoch through all subsequent epochs. In contrast, Pruthi et al. (2020); Xia et al. (2024) make a crude approximation for this propagation. Limitations of influence-based methods are discussed in Schioppa et al. (2023).

Related ideas also were investigated in Liu et al. (2018), which however computes hypergradients for architecture search, rather than deriving per-example data-selection scores as we do. Concurrent work (Savani et al., 2025) also exploits a related directional-derivative symmetry, but in a different setting (token-level antidistillation) rather than per-example data selection.

The recent work of Xia et al. (2024) proposes LESS, a data selection method for instruction tuning that adapts influence ideas to Adam and long sequences. In particular, these authors emphasize the challenge of computing and storing gradients to compute influences. They address this problem via random projections and low-rank approximation. Engstrom et al. (2024) apply the datamodel framework (Ilyas et al., 2022; Park et al., 2023) to select pretraining data. Separately, a *replay* algorithm that stores only a logarithmic number of checkpoints is proposed in Engstrom et al. (2025). Methods that align training data distributions to a small target set include TSDS (Liu et al., 2024) and DSIR (Xie et al., 2023); domain/task-adaptive pretraining also improves transfer (Gururangan et al., 2020). Broader LLM data-efficiency work proposes LLM-guided quality scoring (Ask-LLM) and density sampling (Sachdeva et al., 2024), and clustering-based sensitivity sampling with provable guarantees (Axiotis et al., 2024). Finally, Data Filtering Networks (DFN) also leverage a held-out, high-quality set, but with a different goal and setup (Fang et al., 2023).

Our contribution differs by (i) inverting train/validation roles to approximate per-example influence using only forward losses and doesn't require per-example gradients, or Hessian-vector products, and (ii) showing that this simple, symmetry-based score is computationally inexpensive and outperforms recent data selection approaches for instruction tuning and NER.

2 DATA SELECTION FOR TOKEN-BASED LEARNING

In this section we describe our implementation of the general idea described in the introduction for token-based learning and present empirical results demonstrating its effectiveness. Since prediction takes place at the token level, while data selection takes place at the example level (e.g., instruction/output pair), we compute token scores and aggregate them as described in Section 2.1. Section 2.2 gives a brief overview of instruction-tuning and NER tasks. Experimental settings are introduced in Section 2.3. Empirical results are presented in Sections 2.4 and 2.5.

2.1 SCORE COMPUTATION FOR TOKEN-BASED LEARNING

Each example \mathbf{z} consists of an input \mathbf{z}^{in} and an output \mathbf{z}^{out} , both of which are strings and may differ in length. Let \mathcal{Z}^{out} denote the output vocabulary, and let $T(\mathbf{z})$ denote the length of the output string \mathbf{z}^{out} , which we write as $\mathbf{z}^{\text{out}} = (z^{\text{out}}(1), z^{\text{out}}(2), \dots, z^{\text{out}}(T(\mathbf{z})))$.

Given a model parameterized by θ , its prediction on example \mathbf{z} is a sequence of $T(\mathbf{z})$ conditional distributions, $\{p_t(\cdot | \mathbf{z}, \theta)\}_{t=1}^{T(\mathbf{z})}$, where each $p_t(\cdot | \mathbf{z}, \theta)$ denotes the model's predictive distribution over the output token at position t . Note that $p_t(\cdot | \mathbf{z}, \theta)$ depends on \mathbf{z} solely through \mathbf{z}^{in} and $z^{\text{out}}(1), \dots, z^{\text{out}}(t-1)$. We train models using the log-loss

$$\ell(\theta, \mathbf{z}) = -\frac{1}{T(\mathbf{z})} \sum_{t=1}^{T(\mathbf{z})} \log p_t(z^{\text{out}}(t) | \mathbf{z}; \theta). \quad (7)$$

To compare two models, θ and θ' , on example z , we define a per-token difference of log-loss

$$\Delta_t(z; \theta, \theta') = \log \frac{p_t(z^{\text{out}}(t) | z; \theta')}{p_t(z^{\text{out}}(t) | z; \theta)}. \quad (8)$$

Since our setting involves selecting entire examples rather than individual tokens, we aggregate the per-token differences into a single score per example. Specifically, we apply a transformation function $F: \mathbb{R} \rightarrow \mathbb{R}$ to each Δ_t before averaging across positions. The final score for example z is:

$$\phi(z; \theta, \theta') = \frac{1}{T(z)} \sum_{t=1}^{T(z)} F(\Delta_t(z; \theta, \theta')). \quad (9)$$

We consider three instantiations of the function F , leading to three different scoring methods:

MAXIMUM-IMPROVEMENT: $F(y) = y$ — emphasizes raw improvement.

MAXIMUM-ABSOLUTE CHANGE: $F(y) = |y|$ — captures the magnitude of change.

MAXIMUM-POSITIVE IMPROVEMENT: $F(y) = \max\{y, 0\}$ — ignores degradations.

The algorithm is therefore the same as in Algorithm 1, with the adaptation $\phi_i^{(k)} = \phi(x_i; \hat{\theta}_k^{\text{bas}}, \hat{\theta}_k^{\text{val}})$.

Given a budget of n examples, we choose $S \subseteq [N] \setminus U$, $|S| = n$ using one of these rules:

SCORE-ONLY: Choose the n examples $i \in [N] \setminus U$ that have the largest score ϕ_i .

SCORE+RANDOM: Choose the $n/2$ examples $i \in [N] \setminus U$ that have the largest score ϕ_i , and add $n/2$ more examples chosen uniformly at random (without replacement) from U .

Our scoring schemes tend to favor shorter examples due to their higher variance, which arises from having fewer tokens. To mitigate this bias, we partition the set $[N] \setminus U$ into 10 bins based on sequence length, ensuring each bin contains an equal number of examples. We then select an equal number of top-scoring examples from each bin.

After selecting S of size $|S| = n$, we train (or fine tune) a model on S to evaluate the selection scheme. We refer to this stage as *final training*.

We compare our schemes with three baselines:

RANDOM: The set S is selected uniformly at random subject to its size.

MAXIMUM UNCERTAINTY: Instead of the scores we defined, we use the following hardness score:

$$\psi_i := \frac{1}{T_i} \sum_{t=1}^{T_i} \log(p_t(z_i(t)|z_i; \hat{\theta}_L^{\text{bas}})(1 - p_t(z_i(t)|z_i; \hat{\theta}_L^{\text{bas}})), \quad (10)$$

This score extends the method of Ting & Brochu (2018); Wang et al. (2018); Ai et al. (2021); Kolossov et al. (2024) to token-based learning.

LESS: We used the publicly available implementation from Xia et al. (2024); see Appendix A.1.

2.2 PREDICTION TASKS

We evaluate our data selection framework in two distinct token-based tasks: instruction tuning (IT) and named entity recognition (NER). The framework we introduced above captures both tasks:

Instruction Tuning (IT) involves training a language model to follow natural language instructions. Each training example consists of:

Input z^{in} : a user instruction or prompt; **Output** z^{out} : the desired model response.

The output is typically multi-token and highly variable in content and length, depending on the instruction. The model learns to generate z^{out} conditioned on z^{in} . This naturally fits our framework, which models predictions as token-level distributions $p_t(\cdot | z, \theta)$.

Named Entity Recognition (NER) is a sequence labeling task where the model assigns a probability distribution over entity tags (e.g., PERSON, ORGANIZATION, ...) to each token. In this case: **Input** z^{in} : a tokenized sentence; **Output** z^{out} : a sequence of entity labels, aligned with the input.

In NER, predictions are computed as token-wise classification distributions and therefore output is of the same length as input sequence.¹ In this case, as a base model we take a pretrained language model and replace its prediction head with a classification head.

2.3 EXPERIMENTAL SETTING

In all of our experiments the training set consisted of $N = 36 \times 1024$ samples. For the base model training, we used $|U| = 4 \times 1024$ samples. The validation set size is $m_{\text{val}} = 1024$ and the test set size is $m_{\text{test}} = 10,000$. We vary the selected set size $n \in \{1, 2, 4, 8\} \times 1024$.

Number of epochs. Both for surrogate model training and final model training we determine the number of epochs by $L = (16 \times 1024)/n_{\text{tr}}$. We use a batch size of 16 whence the above ensures that the number batches used in training remains constant, and equal to 1024. In other words, all experiments in this section are at *constant compute*. Since base model training uses $|U| = m = 4 \times 1024$ samples, the number of epochs is $L = 4$.

Learning rate. The learning rate for both surrogate and final model training is selected using hyperparameter optimization for each selected set size n . The learning-rate optimization was carried out for random data selection hence placing our approach at a disadvantage.

We use linear learning rate scheduler and LoRA training [Hu et al. \(2022\)](#) with LoRA parameters $\alpha = 32$ and $\text{dropout} = 0.2$. For NER experiments, we used $\text{PEFTrank} = 1$ and for instruction tuning experiments, we used $\text{PEFTrank} = 256$. The learning rate for the validation examples is $\varepsilon = 1/10$ of the one for the base examples. We present here results with SCORE+RANDOM and refer to the appendix for SCORE-ONLY.

2.4 EXPERIMENTS FOR INSTRUCTION TUNING

For these experiments we used 3 different datasets, which we will refer to as $\mathcal{S} := \{\text{Slim Orca, Alpaca GPT-4, Alpaca GPT-3.5}\}$. As the foundation model, we use **Meta-Llama-3-8B**. Additional details of the model and datasets used are provided in the Appendix.

We designed five experimental setups. In each experiment, one dataset from \mathcal{S} is selected as the *target distribution*. We randomly sample validation and test sets, Z^{val} and Z^{test} , without replacement from the target dataset. These samples are excluded from further use. The *training pool* is then formed by randomly sampling an equal number of examples from one or more datasets in \mathcal{S} (excluding the validation and test samples), such that the total number of selected training samples is fixed at N . We denote by $\mathcal{S}_* \subseteq \mathcal{S}$ the datasets used to generate the training pool. The choices of the target dataset and of \mathcal{S}_* for each of the five experiments are summarized in Table 1. All reported results are averaged over 10 independent runs. In each run, we freshly sample the training, validation, and test sets. These experiments are designed to evaluate performance across a range of data configurations. In particular: in Experiments 1 and 4, the training set includes samples from both target distribution and other distributions; in Experiments 2 and 5, the training set includes samples only from non-target distributions; in Experiment 3, it includes only samples from the target distribution.

Figure 1 summarizes our results for instruction tuning for a fixed select size $n = 8 \times 1024$. We plot the improvement in test log-loss over random data selection for several data-selection strategies within the general framework described in Section 2.1, using Interleaved ToV scoring in algorithm 1 for scoring the examples and SCORE+RANDOM for selecting. We observe that the proposed strategies yield significantly better instruction tuning than random data selection or selecting by max-uncertainty. We observe an improvement (albeit a small one) even when both train and validation data are from Slim Orca (Exp 3), which is a case in which random selection should perform well. The proposed strategies also yield a significant improvement over LESS ([Xia et al., 2024](#)), with the exception of Experiment 2 in which LESS performs slightly better.

Figure 2 displays the evolution of test log loss with selected sample size n . We observe that a good choice of the data selection method results in model improvements that can be equivalent to or larger than doubling n . Plots show standard error (with scaling factor 1) for 10 runs.

¹In NER, typically token level probabilities are combined to assign labels to a whole word.

Table 1: Summary of instruction tuning experiments. Abbreviations: SO = Slim Orca, A4 = Alpaca GPT-4, A3.5 = Alpaca GPT-3.5.

Exp	Target	Training pool
1	SO	SO, A4, and A3.5
2	SO	A4 and A3.5
3	SO	SO
4	A4	SO, A4, and A3.5
5	A4	SO and A3.5

Table 2: Summary of named entity recognition experiments. Abbreviations: MN = Multinerd, A4p = Ai4p, C4 = C4, SB = Syn-big.

Exp	Target	Training pool
1	MN	MN, A4p, C4, and SB
2	MN	A4p, C4, and SB
3	MN	MN
4	A4p	MN, A4p, C4, and SB
5	A4p	MN, C4, and SB
6	A4p	A4p

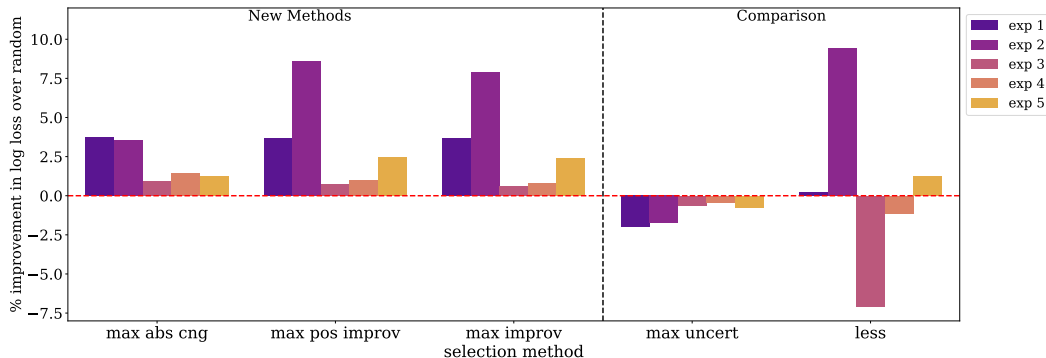


Figure 1: Test log-loss improvement (%) over random selection for instruction tuning with $n = 8 \times 1024$ samples. Each group of bars represents a data-selection strategy (maximum-uncertainty and LESS as baselines); colors show target/training pool configuration (Table 1). Results use Interleaved ToV scoring (Algorithm 1) with the SCORE+RANDOM selection strategy.

2.5 EXPERIMENTS FOR NAMED ENTITY RECOGNITION

The task is to classify whether a token is part of a person name or not. For these experiments we used 4 different labeled datasets, which we will refer to as $\mathcal{S} := \{\text{Multinerd, Ai4p, C4, Syn-big}\}$. We use **xlm-roberta-base** as the foundation model. Further details on the experiment, model and datasets used are presented in the Appendix.

We conducted six sets of experiments. As for the case of instruction tuning, for each set of experiments, we select one of the datasets \mathcal{S} as defining the target distribution, and one or more other datasets to define the training pool (denoted by \mathcal{S}_*). The choices of target datasets and \mathcal{S}_* are summarized in Table 2. The construction of train, test and validation sets is same as in instruction tuning.

Figure 3 summarizes our experiments with NER. We plot the improvement in test log-loss over random data selection for several score definitions. Throughout these experiments, we use SCORE+RANDOM. We observe that the strategies of Section 2.1 yield systematic improvements over random data-selection. Unlike in instruction tuning, maximum uncertainty also improves performance in most settings; however, ToV achieves the largest gains. In this case, LESS (Xia et al., 2024) does not improve over random selection.

We also evaluated token-level F1 on the NER task and observed trends consistent with the log-loss improvements across selection methods (App. A.4).

2.6 RUNTIME AND MEMORY.

We compared the computational cost of ToV and LESS under identical hardware and software settings (Appendix K.3) using the official LESS implementation. We report wall-clock runtime and

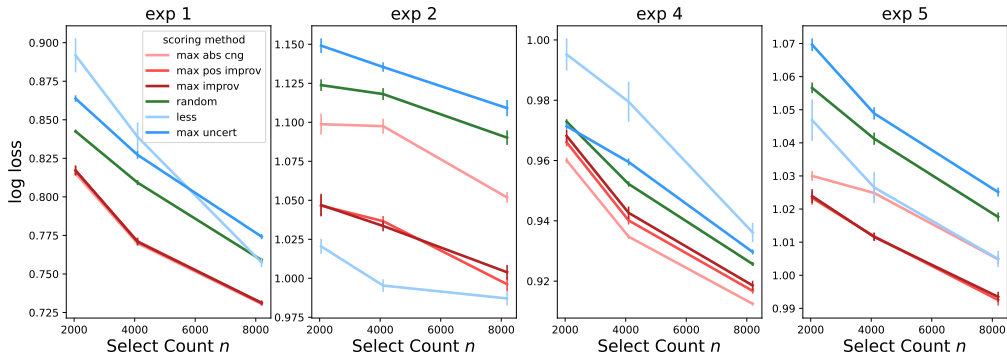


Figure 2: Test log-loss vs. number of selected samples n for instruction tuning. (Due to space limits, Exp. 3 plot is in the Appendix.) Lines show mean log-loss over 10 runs; error bars are ± 1 standard error. Results use Interleaved ToV scoring with the SCORE+RANDOM strategy.

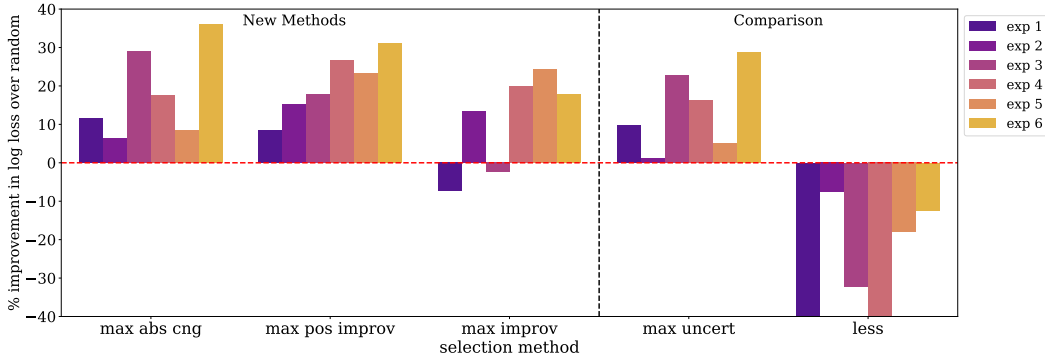


Figure 3: Test log-loss improvement (%) relative to random selection for NER at $n = 8 \times 1024$. Each group of bars represents a data-selection strategy; colors show target/training pool configuration (Table 2). Results use Interleaved ToV scoring (Algorithm 1) with the SCORE+RANDOM strategy

disk storage required by each method (e.g., checkpoints and auxiliary tensors). Reported values are averages over five runs per task. Across both instruction tuning and NER, ToV is substantially faster and more storage-efficient. Results are summarized in Table 3.

Table 3: Runtime and disk storage for ToV vs. LESS under identical hardware (mean of 5 runs).

Setting	Method	Runtime	Storage
Instruction tuning	LESS	4h 5m	4.9 GB
Instruction tuning	ToV	2h 3m	1.84 GB
NER	LESS	46m	4.1 GB
NER	ToV	8m	0.24 GB

ToV reduces runtime by 2–6 \times and disk storage by 2.5–16 \times relative to LESS across the evaluated tasks.

3 A FORMAL JUSTIFICATION

In this section, we present a mathematical analysis of our approach in the case of batch gradient descent (GD). We focus on the implementation *Parallel ToV* (Method B), described in Algorithm 2.

Parallel ToV differs from Interleaved ToV because at each training cycle k , training on the base set X_U is initialized with the output of the previous train-on-validation phase. Empirically Interleaved

Algorithm 2 ToV Scoring Algorithm: *Parallel ToV*

-
- 1: **Input:** Pretrained model θ_0 , validation set Z^{val} , training pool $X = (\mathbf{x}_i : i \in [N])$,
 - 2: selected data count $n \leq N$, base model count m
 - 3: **Output:** Set of examples $S \subset [N]$ of size n
 - 4: Sample base subset $U \subseteq [N]$ of size m randomly; define $X_U = (\mathbf{x}_i : i \in U)$
 - 5: Initialize models: $\hat{\theta}_0^{\text{bas},+} \leftarrow \theta_0$, $\hat{\theta}_0^{\text{bas}} \leftarrow \theta_0$; set scores $\Upsilon_i \leftarrow 0$ for all $i \in [N] \setminus U$
 - 6: **for** $k = 1$ to L **do**
 - 7: Train for one epoch on X_U with learn. rate η_k and init. $\hat{\theta}_{k-1}^{\text{bas},+}$. Denote the output by $\hat{\theta}_{0,k}^{\text{bas},+}$
 - 8: Train for one epoch on Z^{val} with learn. rate $\varepsilon \cdot \eta_k$ and init. $\hat{\theta}_{0,k}^{\text{bas},+}$. Denote the output by $\hat{\theta}_k^{\text{bas},+}$
 - 9: Train for one epoch on X_U with learn. rate η_k and init. $\hat{\theta}_{k-1}^{\text{bas}}$. Denote the output by $\hat{\theta}_k^{\text{bas}}$
 - 10: **for** each $i \in [N] \setminus U$ **do**
 - 11: $\Upsilon_i^{(k)} \leftarrow \ell(\hat{\theta}_k^{\text{bas}}, \mathbf{x}_i) - \ell(\hat{\theta}_k^{\text{bas},+}, \mathbf{x}_i)$
 - 12: $\Upsilon_i \leftarrow \Upsilon_i + \Upsilon_i^{(k)} / L$
 - 13: **end for**
 - 14: **end for**
 - 15: Select $S \subseteq [N] \setminus U$ with size $|S| = n$ using scores Υ_i
-

ToV performs somewhat better than Parallel ToV, see Appendix C. We use Parallel ToV for analysis just because the resulting mathematical expressions are simpler.

We find empirically that the ToV works well beyond token-based learning, and hence our focus will be to understand it in a generic learning problem. Appendix D demonstrates this point by considering a simple logistic regression problem.

3.1 IDEAL SCORES, LINEARIZATION, INFLUENCE FUNCTIONS

In order to estimate the model improvement produced by sample $i \in [N] \setminus U$ we could train a model on two training sets X_U and $X_{U \cup i}$, using empirical risk functions $\hat{R}_U(\theta)$, $\hat{R}_{U \cup i}(\theta)$. We thus would run GD for L steps, with initialization $\hat{\theta}_0^{\text{bas}} = \hat{\theta}_0^{\text{bas}+i} = \theta_0$:

$$\hat{\theta}_{k+1}^{\text{bas}} = \hat{\theta}_k^{\text{bas}} - \eta m \nabla \hat{R}_U(\hat{\theta}_k^{\text{bas}}), \quad \hat{\theta}_{k+1}^{\text{bas}+i} = \hat{\theta}_k^{\text{bas}+i} - \eta(m+1) \nabla \hat{R}_{U \cup i}(\hat{\theta}_k^{\text{bas}+i}). \quad (11)$$

At iteration k , we have thus two models $\hat{\theta}_k^{\text{bas}}$ and $\hat{\theta}_k^{\text{bas}+i}$ that differ uniquely in whether sample i is used or not. We define the *ideal score* to be the difference in validation error between these two models, averaged over epochs

$$S_i := \frac{1}{L} \sum_{s=1}^L [\hat{R}_{\text{val}}(\hat{\theta}_s^{\text{bas}}) - \hat{R}_{\text{val}}(\hat{\theta}_s^{\text{bas}+i})] = \frac{1}{m_{\text{val}} L} \sum_{s=1}^L \sum_{j=1}^{m_{\text{val}}} \{ \ell(\hat{\theta}_s^{\text{bas}}, \mathbf{z}_j^{\text{val}}) - \ell(\hat{\theta}_s^{\text{bas}+i}, \mathbf{z}_j^{\text{val}}) \}. \quad (12)$$

Evaluating this score is computationally expensive, hence several groups (Pruthi et al., 2020; Bae et al., 2024; Xia et al., 2024) proposed to use a first order Taylor expansion to approximate the difference between the two models. Expanding S_i with respect to the contribution of $\ell(\cdot; \mathbf{x}_i)$ yields

$$S_i^{\text{lin}} = \frac{\eta}{L} \sum_{0 \leq s < t \leq L} \langle \nabla \hat{R}_{\text{val}}(\hat{\theta}_t^{\text{bas}}), \mathbf{M}_{t,s+1} \nabla \ell(\hat{\theta}_s^{\text{bas}}, \mathbf{x}_i) \rangle. \quad (13)$$

where $\mathbf{M}_{t,t} = \mathbf{I}_d$ and $\mathbf{M}_{t,r}$ captures the propagation of perturbations along the GD trajectory:

$$\mathbf{M}_{t,r} := \mathbf{H}_{t-1} \cdot \mathbf{H}_{t-2} \cdots \mathbf{H}_r, \quad \mathbf{H}_k := \mathbf{I} - \eta m \nabla^2 \hat{R}_U(\hat{\theta}_k^{\text{bas}}). \quad (14)$$

The next result shows that S_i^{lin} approximates well S_i in a quantitative way, under local convexity.

Proposition 1. *Assume there exist $c_0, C_1, M > 0$ such that $\nabla^2 \hat{R}_U(\hat{\theta}_k^{\text{bas}}) \succeq c_0 \mathbf{I}_d$, $\|\nabla \ell(\hat{\theta}_k^{\text{bas}}, \mathbf{x}_i)\| \leq C_1$ for all k and, for all θ_1, θ_2 , $\|\nabla^2 \hat{R}_U(\theta_1) - \nabla^2 \hat{R}_U(\theta_2)\|_{\text{op}} \leq M \|\theta_1 - \theta_2\|_2$, $\|\nabla \ell(\theta_1; \mathbf{x}_i) - \nabla \ell(\theta_2; \mathbf{x}_i)\|_{\text{op}} \leq M \|\theta_1 - \theta_2\|_2$. Further assume that $\|\nabla^2 \hat{R}_{\text{val}}(\hat{\theta}_k^{\text{bas}})\|_{\text{op}} \leq C_1$ and $\|\nabla^2 \hat{R}_{\text{val}}(\theta_1) - \nabla^2 \hat{R}_{\text{val}}(\theta_2)\|_{\text{op}} \leq M \|\theta_1 - \theta_2\|_2$ for all θ_1, θ_2 as well. Finally, assume there exists a constant C_η such that $\eta_k = \eta \leq C_\eta / m \forall k$. Then there exists $C = C(c_0, C_1, C_\eta, M)$ such that*

$$|S_i - S_i^{\text{lin}}| \leq C / m^2. \quad (15)$$

The assumption $\eta \leq C_\eta/m$ is justified by the fact that we expect the Hessian of $\widehat{R}_U(\cdot)$ to be of order one, and hence the stepsize for this objective (which is given by ηm see Eq. (11)) should be of order one. As shown in the proof, the typical size of S_i^{lin} is of order $1/m$, and hence Eq. (15) establishes that the difference $|S_i - S_i^{\text{lin}}|$ is negligible.

We emphasize that Proposition 1 does not make any assumption about the data distribution and instead entirely relies on regularity properties of the risk function.

3.2 TRAIN-VALIDATION DUALITY

We consider Interleaved ToV and Parallel ToV defined in Algorithms 1, 2. We emphasize the dependence on ε by writing $\phi_i = \phi_i(\varepsilon)$ and $\Upsilon_i = \Upsilon_i(\varepsilon)$. It is easy to derive the small ε asymptotics $\phi_i = \phi_i^{\text{lin}}\varepsilon + o(\varepsilon)$, $\Upsilon_i = \Upsilon_i^{\text{lin}}\varepsilon + o(\varepsilon)$, where, for $\mathbf{g}_{s,i} := \nabla \ell(\widehat{\boldsymbol{\theta}}_s^{\text{bas}}; \mathbf{x}_i)$,

$$\phi_i^{\text{lin}} := \frac{\eta m_{\text{val}}}{L} \sum_{s=1}^L \langle \nabla \widehat{R}_{\text{val}}(\widehat{\boldsymbol{\theta}}_s^{\text{bas}}), \mathbf{g}_{s,i} \rangle, \quad \Upsilon_i^{\text{lin}} := \frac{\eta m_{\text{val}}}{L} \sum_{0 \leq t < s \leq L} \langle \nabla \widehat{R}_{\text{val}}(\widehat{\boldsymbol{\theta}}_{t+1}^{\text{bas}}), \mathbf{M}_{s,t+1}^\top \mathbf{g}_{s,i} \rangle. \quad (16)$$

We show that these are good approximations of $\Upsilon_i(\varepsilon)$, $\phi_i(\varepsilon)$ uniformly in dimension, sample size.

Theorem 1. Consider Algorithms 1, 2 with fixed stepsize $\eta_k = \eta$ (and $F(x) = -x$ in Algorithm 1). Under the assumptions of Proposition 1, further assume $\|\nabla \widehat{R}_{\text{val}}(\widehat{\boldsymbol{\theta}}_k^{\text{bas}})\| \leq C_1$ for all k , and $\|\nabla_{\boldsymbol{\theta}}^2 \ell(\boldsymbol{\theta}; \mathbf{x})\|_{\text{op}} \leq C_1$. Then there exist $c_* = c_*(c_0, M, C_1)$, $C = C(c_0, M, C_1)$ such that, for $\varepsilon m_{\text{val}}/m \leq c_*$,

$$|\Upsilon_i(\varepsilon) - \Upsilon_i^{\text{lin}}\varepsilon| \leq C(\varepsilon m_{\text{val}}/m)^2, \quad |\phi_i(\varepsilon) - \phi_i^{\text{lin}}\varepsilon| \leq C(\varepsilon m_{\text{val}}/m)^2. \quad (17)$$

Note that Υ_i^{lin} differ from S_i^{lin} because of: (i) The different order of s and t ; (ii) The fact that $\mathbf{M}_{t,s+1}$ is replaced by its transpose in Eq. (16). Υ_i^{lin} measures the influence of training on validation data when making inference at \mathbf{x}_i , while S_i^{lin} measures the influence of training on \mathbf{x}_i data when making inference on validation. These two measures of ‘influence’ differ by the replacement of $\mathbf{M}_{t,s+1}$ by $\mathbf{M}_{s,t}^\top$. However, in a number of cases we expect these two matrices to be not too different, and hence the two scores to yield similar results. We can prove that Υ_i^{lin} and S_i^{lin} coincide (for large L) under local convexity conditions.

Theorem 2. Assume $\boldsymbol{\theta} \mapsto \ell(\boldsymbol{\theta}; \mathbf{x})$ to be twice continuously differentiable and that $\|\nabla \widehat{R}_{\text{val}}(\widehat{\boldsymbol{\theta}}_k^{\text{bas}})\| \leq C_1$, $\|\nabla \ell(\widehat{\boldsymbol{\theta}}_k^{\text{bas}}; \mathbf{x}_i)\| \leq C_1$ for all k . Further assume that gradient descent iterates $(\widehat{\boldsymbol{\theta}}_k^{\text{bas}} : k \geq 0)$ converge to $\widehat{\boldsymbol{\theta}}_\infty^{\text{bas}} = \lim_{k \rightarrow \infty} \widehat{\boldsymbol{\theta}}_k^{\text{bas}}$ which is a local minimum of $\widehat{R}_U(\boldsymbol{\theta})$ with $\mathbf{Q}_\infty := \nabla^2 \widehat{R}_U(\widehat{\boldsymbol{\theta}}_\infty^{\text{bas}}) \succ \mathbf{0}$ (strictly). Then

$$\lim_{L \rightarrow \infty} \frac{1}{m_{\text{val}}} \Upsilon_i^{\text{lin}}(L) = \lim_{L \rightarrow \infty} S_i^{\text{lin}}(L) = \frac{1}{m} \langle \nabla \widehat{R}_{\text{val}}(\widehat{\boldsymbol{\theta}}_\infty^{\text{bas}}), \mathbf{Q}_\infty^{-1} \nabla \ell(\widehat{\boldsymbol{\theta}}_\infty^{\text{bas}}; \mathbf{x}_i) \rangle := S_{i,\infty}^{\text{lin}}. \quad (18)$$

The last expression in Eq. (18) (denoted by $S_{i,\infty}^{\text{lin}}$) is the classical formula for influence functions of M-estimators (van der Vaart, 2000). Both our approach and the dynamical influence function $S_i^{\text{lin}}(L)$ can be regarded as approximations of $S_{i,\infty}^{\text{lin}}$ in this case.

In fine tuning, the model is likely to be overparametrized, and it is unrealistic to assume convergence to a strict minimum (with $\nabla^2 \widehat{R}_U(\widehat{\boldsymbol{\theta}}_\infty^{\text{bas}}) \succ \mathbf{0}$). On the other hand, the weights will not change significantly during this phase and it is reasonable to approximate fine-tuning as fitting an overparametrized linear model with respect to the empirical neural tangent features learnt in the pre-training phase.

Theorem 3. Consider the loss function $\ell(\boldsymbol{\theta}; \mathbf{x}) = (y(\mathbf{x}) - \langle \boldsymbol{\psi}(\mathbf{x}), \boldsymbol{\theta} \rangle)^2/2$ for some response variables $y(\mathbf{x})$, and featurization map $\boldsymbol{\psi} : \mathbb{R}^d \rightarrow \mathbb{R}^p$, $p > m$. Let $\boldsymbol{\Psi} \in \mathbb{R}^{|U| \times p}$ be the matrix with rows $(\boldsymbol{\psi}(\mathbf{x}_j) : j \in U)$, $\boldsymbol{\Psi}_{\text{val}} \in \mathbb{R}^{m_{\text{val}} \times p}$ be the matrix with rows $(\boldsymbol{\psi}(\mathbf{z}_j^{\text{val}}) : j \leq m_{\text{val}})$, $\mathbf{P}_{\boldsymbol{\Psi}}$ the projector to the kernel of $\boldsymbol{\Psi}$, $\mathbf{y} = (y(\mathbf{x}_j) : j \in U)$, $\widehat{\boldsymbol{\theta}} := \boldsymbol{\Psi}^\dagger \mathbf{y}$, $\mathbf{r}^{\text{val}} := (y(\mathbf{z}_j^{\text{val}}) - \langle \widehat{\boldsymbol{\theta}}, \boldsymbol{\psi}(\mathbf{z}_j^{\text{val}}) \rangle : j \leq m_{\text{val}})$, $r(i) := y(\mathbf{x}_i) - \langle \widehat{\boldsymbol{\theta}}, \boldsymbol{\psi}(\mathbf{x}_i) \rangle$. If GD is initialized with $\boldsymbol{\theta}_0 = \mathbf{0}$, and we use constant stepsize $\eta < \|\boldsymbol{\Psi}\|_{\text{op}}^2/2$, then

$$\lim_{L \rightarrow \infty} \frac{1}{L m_{\text{val}}} \Upsilon_i^{\text{lin}}(L) = \lim_{L \rightarrow \infty} \frac{1}{L} S_i^{\text{lin}}(L) = \frac{\eta}{2} r(i) \langle \mathbf{r}^{\text{val}}, \boldsymbol{\Psi}_{\text{val}}^\top \mathbf{P}_{\boldsymbol{\Psi}} \boldsymbol{\psi}(\mathbf{x}_i) \rangle. \quad (19)$$

ACKNOWLEDGEMENTS

We are grateful to Neeraja Abhyankar, Alankrita Bhatt, Joseph Gardi, Mukur Gupta, Germain Kolossov, Marc Laugharn, Rahul Ponnala, Sahasrajit Sarmasarkar, Andreas Santucci and Pulkit Tandon, for several conversations about this work.

REFERENCES

- Mingyao Ai, Jun Yu, Huiming Zhang, and HaiYing Wang. Optimal subsampling algorithms for big data regressions. *Statistica Sinica*, 31(2):749–772, 2021.
- AI@Meta. Llama 3 model card. 2024. URL https://github.com/meta-llama/llama3/blob/main/MODEL_CARD.md.
- Kyriakos Axiotis, Vincent Cohen-Addad, Monika Henzinger, Sammy Jerome, Vahab Mirrokni, David Saulpic, David Woodruff, and Michael Wunder. Data-efficient learning via clustering-based sensitivity sampling: Foundation models and beyond. *arXiv preprint arXiv:2402.17327*, 2024.
- Juhan Bae, Wu Lin, Jonathan Lorraine, and Roger Grosse. Training data attribution via approximate unrolled differentiation. *arXiv:2405.12186*, 2024.
- Peter L Bartlett, Andrea Montanari, and Alexander Rakhlin. Deep learning: a statistical viewpoint. *Acta numerica*, 30:87–201, 2021.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned language models. *Journal of Machine Learning Research*, 25(70):1–53, 2024.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. Un-supervised cross-lingual representation learning at scale. *CoRR*, abs/1911.02116, 2019. URL <http://arxiv.org/abs/1911.02116>.
- Logan Engstrom, Axel Feldmann, and Aleksander Madry. Dsdm: Model-aware dataset selection with datamodels. *arXiv preprint arXiv:2401.12926*, 2024.
- Logan Engstrom, Andrew Ilyas, Benjamin Chen, Axel Feldmann, William Moses, and Aleksander Madry. Optimizing ml training with metagradient descent. *arXiv preprint arXiv:2503.13751*, 2025.
- Alex Fang, Albin Madappally Jose, Amit Jain, Ludwig Schmidt, Alexander Toshev, and Vaishaal Shankar. Data filtering networks. *arXiv preprint arXiv:2309.17425*, 2023.
- Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A Smith. Don’t stop pretraining: Adapt language models to domains and tasks. *arXiv preprint arXiv:2004.10964*, 2020.
- Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022.
- Andrew Ilyas, Sung Min Park, Logan Engstrom, Guillaume Leclerc, and Aleksander Madry. Data-models: Predicting predictions from training data. In *Proceedings of the 39th International Conference on Machine Learning*, 2022.
- Germain Kolossov, Andrea Montanari, and Pulkit Tandon. Towards a statistical theory of data selection under weak supervision. In *The Twelfth International Conference on Learning Representations*, 2024.
- Wing Lian, Guan Wang, Bleys Goodson, Eugene Pentland, Austin Cook, Chanvichet Vong, and "Teknium". Slimorca: An open dataset of gpt-4 augmented flan reasoning traces, with verification, 2023. URL <https://huggingface.co/Open-Orca/SlimOrca>.

- Hanxiao Liu, Karen Simonyan, and Yiming Yang. Darts: Differentiable architecture search. *arXiv preprint arXiv:1806.09055*, 2018.
- Zifan Liu, Amin Karbasi, and Theodoros Rekatsinas. Tsds: Data selection for task-specific model finetuning. *Advances in Neural Information Processing Systems*, 37:10117–10147, 2024.
- Shayne Longpre, Le Hou, Tu Vu, Albert Webson, Hyung Won Chung, Yi Tay, Denny Zhou, Quoc V Le, Barret Zoph, Jason Wei, et al. The flan collection: Designing data and methods for effective instruction tuning. In *International Conference on Machine Learning*, pp. 22631–22648. PMLR, 2023.
- Subhabrata Mukherjee, Arindam Mitra, Ganesh Jawahar, Sahaj Agarwal, Hamid Palangi, and Ahmed Awadallah. Orca: Progressive learning from complex explanation traces of gpt-4, 2023.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35: 27730–27744, 2022.
- Sung Min Park, Kristian Georgiev, Andrew Ilyas, Guillaume Leclerc, and Aleksander Madry. Trak: Attributing model behavior at scale. *arXiv preprint arXiv:2303.14186*, 2023.
- Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. Instruction tuning with gpt-4. *arXiv preprint arXiv:2304.03277*, 2023.
- Garima Pruthi, Frederick Liu, Satyen Kale, and Mukund Sundararajan. Estimating training data influence by tracing gradient descent. *Advances in Neural Information Processing Systems*, 33: 19920–19930, 2020.
- Noveen Sachdeva, Benjamin Coleman, Wang-Cheng Kang, Jianmo Ni, Lichan Hong, Ed H Chi, James Caverlee, Julian McAuley, and Derek Zhiyuan Cheng. How to train data-efficient llms. *arXiv preprint arXiv:2402.09668*, 2024.
- Yash Savani, Asher Trockman, Zhili Feng, Yixuan Even Xu, Avi Schwarzschild, Alexander Robey, Marc Finzi, and J Zico Kolter. Antidistillation sampling. *arXiv preprint arXiv:2504.13146*, 2025.
- Andrea Schioppa, Katja Filippova, Ivan Titov, and Polina Zablotskaia. Theoretical and practical perspectives on what influence functions do. *Advances in Neural Information Processing Systems*, 36:27560–27581, 2023.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca, 2023.
- Simone Tedeschi and Roberto Navigli. MultiNERD: A multilingual, multi-genre and fine-grained dataset for named entity recognition (and disambiguation). In *Findings of the Association for Computational Linguistics: NAACL 2022*, pp. 801–812, Seattle, United States, July 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.findings-naacl.60. URL <https://aclanthology.org/2022.findings-naacl.60>.
- Daniel Ting and Eric Brochu. Optimal subsampling with influence functions. *Advances in neural information processing systems*, 31, 2018.
- Aaad W van der Vaart. *Asymptotic Statistics*. Cambridge University Press, 2000.
- HaiYing Wang, Rong Zhu, and Ping Ma. Optimal subsampling for large sample logistic regression. *Journal of the American Statistical Association*, 113(522):829–844, 2018.
- Yizhong Wang, Hamish Ivison, Pradeep Dasigi, Jack Hessel, Tushar Khot, Khyathi Chandu, David Wadden, Kelsey MacMillan, Noah A Smith, Iz Beltagy, et al. How far can camels go? exploring the state of instruction tuning on open resources. *Advances in Neural Information Processing Systems*, 36:74764–74786, 2023.

Zifeng Wang, Hong Zhu, Zhenhua Dong, Xiuqiang He, and Shao-Lun Huang. Less is better: Unweighted data subsampling via influence function. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pp. 6340–6347, 2020.

Mengzhou Xia, Sadhika Malladi, Suchin Gururangan, Sanjeev Arora, and Danqi Chen. LESS: selecting influential data for targeted instruction tuning. In *Proceedings of the 41st International Conference on Machine Learning*, pp. 54104–54132, 2024.

Sang Michael Xie, Shibani Santurkar, Tengyu Ma, and Percy S Liang. Data selection for language models via importance resampling. *Advances in Neural Information Processing Systems*, 36: 34201–34227, 2023.

Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, et al. Lima: Less is more for alignment. *Advances in Neural Information Processing Systems*, 36, 2024.

A ADDITIONAL EXPERIMENTAL DETAILS AND RESULTS

A.1 EXPERIMENTS FOR TOKEN-BASED LEARNING

In these experiments, we used pretrained models as base models and constructed training, validation, and test sets from real-world datasets. Details of the datasets and models are provided in Section K.

For each training example count—both for surrogate model training (used for scoring) and final model training—we selected the learning rate from the following grid:

$$[3e-6, 1e-5, 3e-5, 1e-4, 3e-4, 1e-3, 3e-3, 1e-2].$$

The optimal learning rate was determined by training models on randomly sampled subsets from the training pool and evaluating their test log-loss. For each learning rate, the loss was averaged over 10 runs, with a new random subset used in each run. The best-performing learning rate was selected separately for each experimental configuration listed in Table 1 and Table 2.

Implementation details for Less (Xia et al., 2024) We used the public implementation from the authors’ GitHub repository. The projection dimension was set to 8192. Learning rate and other hyperparameters were tuned identically for all approaches. For both our method and LESS, the surrogate model used the same number of samples and was trained for four epochs, matching the settings in the LESS paper. Following the original LESS procedure, we selected the top-scoring examples.

A.2 EXPANDED RESULTS FOR INSTRUCTION TUNING

In the main paper, we compared our scoring methods for the SCORE+RANDOM strategy. Due to space constraints, Figure 2 omitted results for Experiment 3. In Figure 4, we provide an expanded version that includes results for Experiment 3 as well.

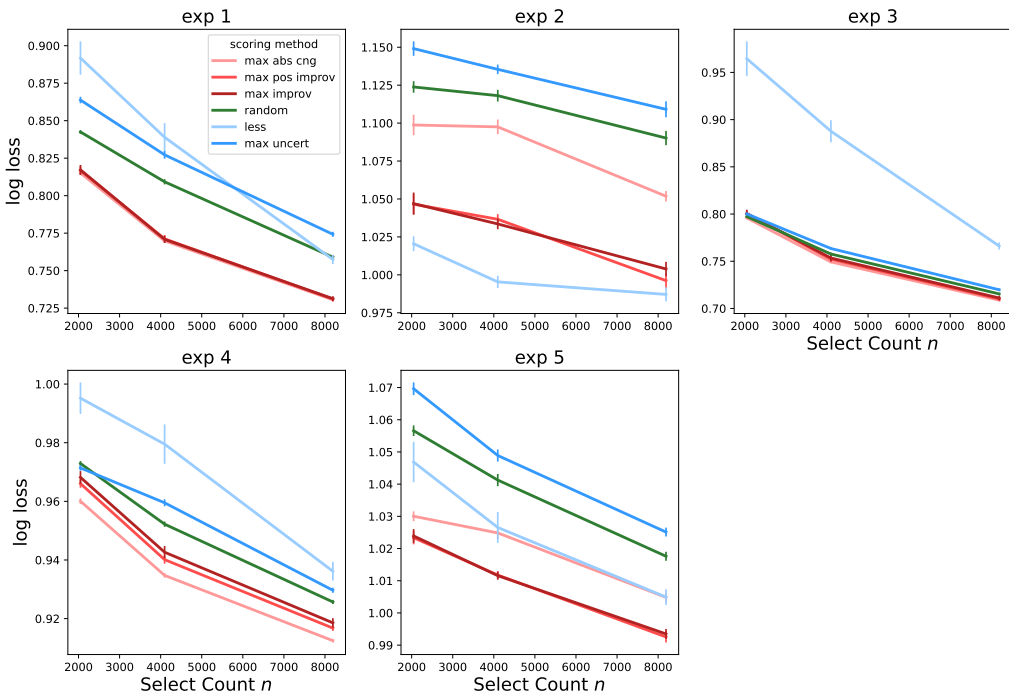


Figure 4: Expanded version of Figure 2 including the Experiment 3 plot.

A.3 EXPANDED RESULTS FOR NAMED ENTITY RECOGNITION

In Figure 3 of the main paper, we reported results for SCORE+RANDOM using a fixed selected sample size of $n = 8 \times 1024$, across all experiment configurations in Table 2.

In Figure 5, we show how the test log-loss varies with the selected sample size n for different scoring methods under the SCORE+RANDOM strategy, and how these compare to random selection.

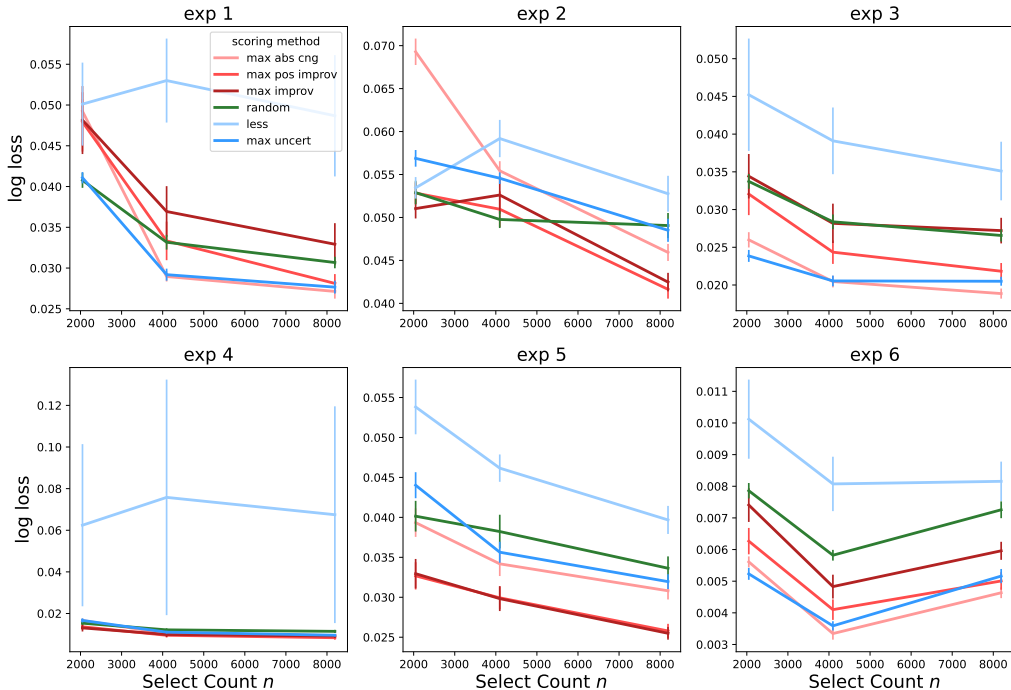


Figure 5: Test log-loss vs. number of selected samples n for NER. Lines show mean log-loss over 10 runs; error bars are ± 1 standard error. Results use Interleaved ToV with the SCORE+RANDOM strategy.

A.4 ADDITIONAL NER RESULTS: TOKEN-LEVEL F1

To complement the log-loss evaluation reported above and in the main paper, we also measure token-level F1 on the NER task. Figure 6 shows test 1 – F1 as a function of the number of selected samples n , using the same experimental setup as Fig. 5. Lower values therefore, indicate better sequence labeling performance.

We observe trends consistent with the log-loss results: our selection methods outperform random selection and competing baselines in several settings, confirming that cross-entropy gains translate into improved token-level F1.

B COMPARISON OF SCORE+RANDOM AND SCORE-ONLY SELECTION

In this section, we examine how the performance of our strategies changes when all training examples are selected from the top-scoring set (SCORE-ONLY) instead of selecting only half of them from the top and the other half at random (SCORE+RANDOM).

Recall that our scores approximate how much benefit each example provides when added to a randomly chosen pool of training data. A higher score therefore indicates an example expected to be

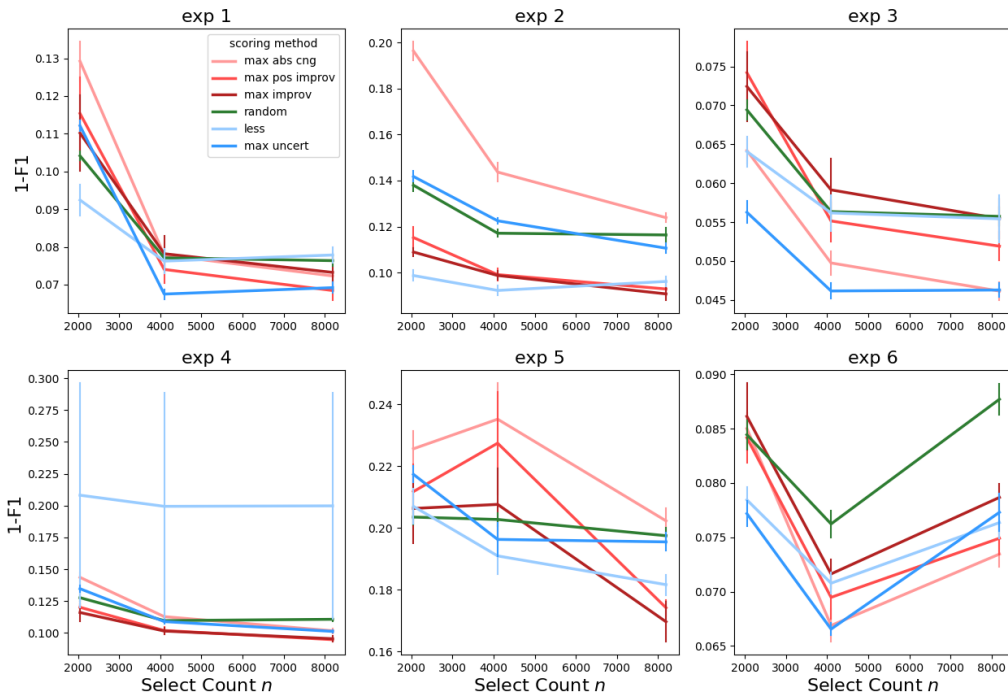


Figure 6: Test 1 – F1 vs. number of selected samples n for NER. Lines show mean 1 – F1 over 10 runs; error bars denote ± 1 standard error. Results use Interleaved ToV with the SCORE+RANDOM strategy.

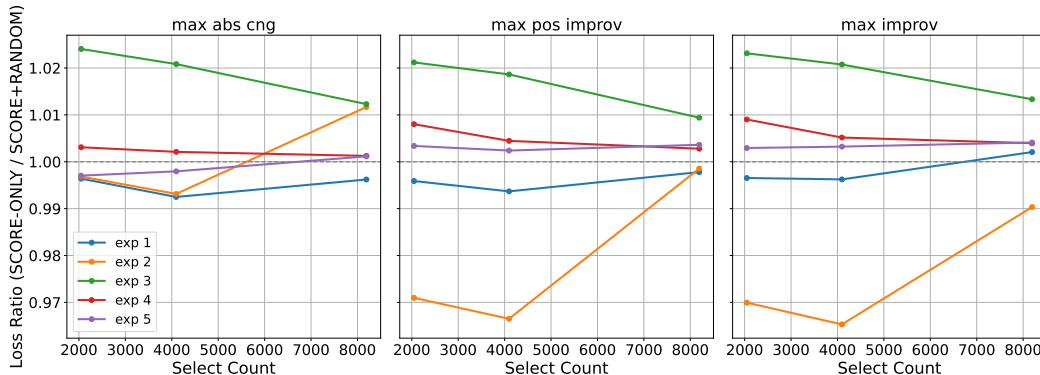


Figure 7: Ratio of log loss for SCORE-ONLY versus SCORE+RANDOM across our three scoring strategies and all instruction-tuning setups in Table 1. Scores are computed using Algorithm 1.

more helpful in that setting. SCORE+RANDOM selects half of the final training set from the highest-scoring examples and fills the rest with random examples, whereas SCORE-ONLY takes only the top-scoring examples. This design creates a trade-off:

- **Pure exploitation:** Selecting only top-scoring examples can maximize immediate gain because every chosen example has a high estimated contribution.
- **Score validity and diversity:** The scores are defined relative to adding examples to a random pool. If we select only top examples, the resulting set may differ substantially from the

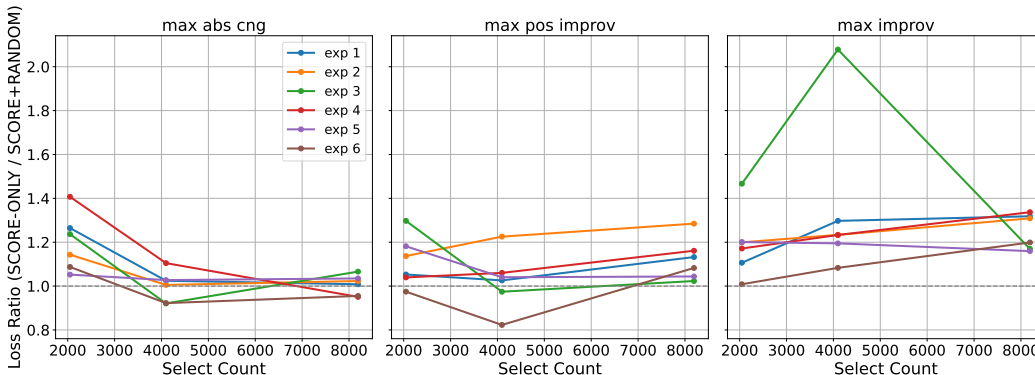


Figure 8: Ratio of log-loss for SCORE-ONLY versus SCORE+RANDOM across our three scoring strategies and all NER setups in Table 2. Scores are computed using Algorithm 1.

random reference, making the scores a less accurate guide. Randomly adding half the examples keeps the final set closer to the conditions under which the scores were computed and also protects against loss of diversity.

Which effect dominates varies by task.

Figures 7 and 8 show the ratio of log-loss for the two selection strategies in instruction tuning and NER respectively. In each figure the three subplots correspond to our three scoring strategies; different lines indicate the various experimental setups. Algorithm 1 is used to obtain the scores.

For instruction tuning, SCORE+RANDOM performs better in three setups (3, 4, 5), while SCORE-ONLY is better in the remaining two (1, 2) across most selection sizes and scoring methods. For NER, SCORE+RANDOM tends to outperform SCORE-ONLY more often, particularly for the Max-Improvement scores.

C PARALLEL TOV VS INTERLEAVED TOV

All previous plots used Interleaved ToV (Algorithm 1) for scoring. Here we compare the performance of the two scoring methods: Interleaved ToV (Algorithm 1) and Parallel ToV (Algorithm 2)—across our experiments, using the SCORE+RANDOM selection strategy for both.

Figures 9 and 10 show the ratio of test log-loss obtained with Parallel ToV relative to Interleaved ToV for instruction tuning and NER, respectively. Each figure contains three subplots corresponding to our three scoring strategies, and different lines represent the various experimental setups.

The results indicate that for instruction tuning, Interleaved ToV is most often superior, while for NER there is no consistent winner. A possible explanation is that Parallel ToV uses two distinct training trajectories. Our analysis assumes that the resulting models differ only slightly, but in practice, the two training trajectories can diverge substantially. This effect is likely to be stronger with large and highly overparameterized models such as Meta-Llama-3-8B, which we used for instruction tuning. We expect the larger distance between the two models to result in less accurate score estimation in Parallel ToV, as compared to Interleaved ToV.

D LOGISTIC REGRESSION EXPERIMENTS

In these experiments, we synthetically generated the training pool, validation set, and test set. We begin by defining a parametric family of distributions used to construct the data.

For a given $p > 0$ and parameter vector $\theta \in \mathbb{R}^p$, we define a distribution \mathcal{P}_θ over pairs (x, y) , where $x \in \mathbb{R}^p$ and $y \in \{0, 1\}$. The features are sampled as $x \sim \mathcal{N}(0, I)$, and the label y is drawn according

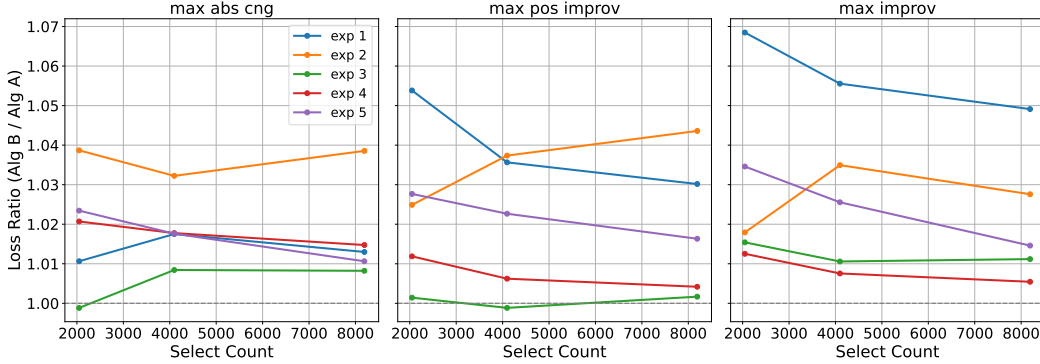


Figure 9: Ratio of test log-loss for Parallel ToV (Method B; Algorithm 2) relative to Interleaved ToV (Method A; Algorithm 1) for instruction tuning. Results use the SCORE+RANDOM selection strategy. Each subplot corresponds to one scoring strategy; lines denote different experimental setups in Table 1.

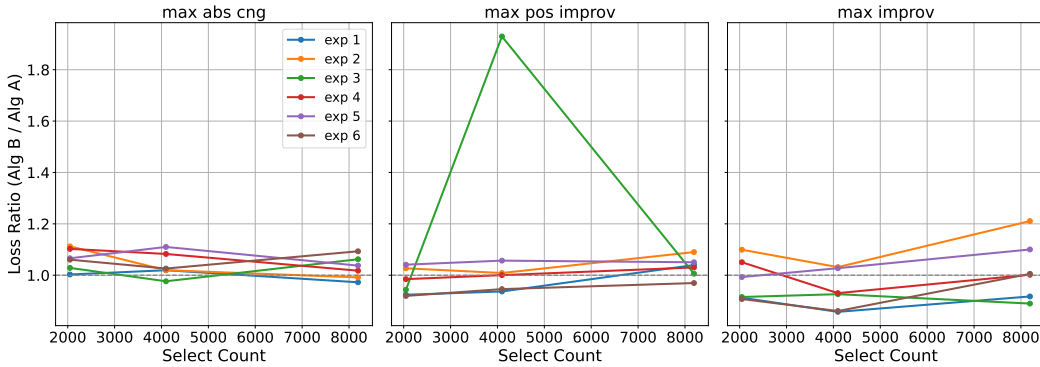


Figure 10: Ratio of test log-loss for Parallel ToV (Method B; Algorithm 2) relative to Interleaved ToV (Method A; Algorithm 1) for NER. Results use the SCORE+RANDOM selection strategy. Each subplot corresponds to one scoring strategy; lines denote different experimental setups in Table 2.

to a logistic model:

$$\Pr(y = 1 | x) = \frac{1}{1 + \exp(-x \cdot \theta)}, \quad \Pr(y = 0 | x) = 1 - \Pr(y = 1 | x).$$

We randomly sample a unit vector θ^* from the unit sphere to serve as the target direction. A second unit vector θ' is then drawn such that it lies at an angle γ from θ^* . In our experiments, we set $p = 10$ and $\gamma = \pi/2$.

The training pool consists of $N = 128 \times 1024$ samples, drawn independently from the mixture distribution:

$$\mathcal{D}_{\text{train}} = \frac{1}{2} \mathcal{P}_{\theta^*} + \frac{1}{2} \mathcal{P}_{\theta'}.$$

The validation and test sets contain $m_{\text{val}} = 1024$ and $m_{\text{tst}} = 10,000$ samples respectively, both drawn i.i.d. from the target distribution \mathcal{P}_{θ^*} .

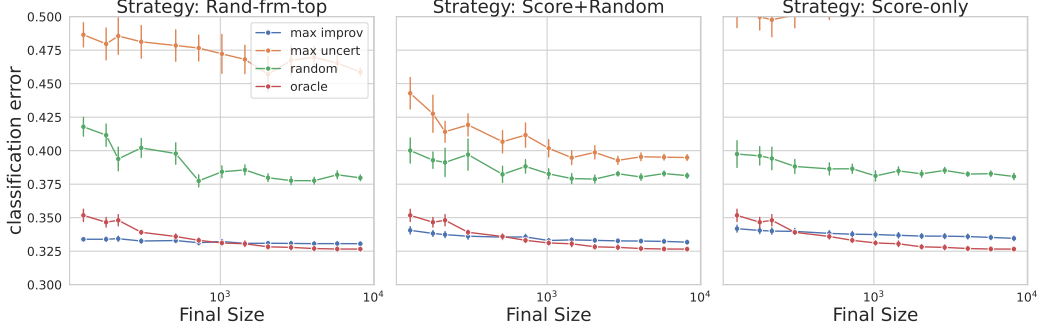


Figure 11: Data selection experiments with Parallel ToV for logistic regression on synthetic data in $d = 10$ dimensions (4 epochs of training). Each color corresponds to a distinct method to score data in the training pool, and each frame to a distinct method to use the score to form the selected set. Each symbol corresponds to the average of 10 experiments.

For scoring, we used Parallel ToV as described in Algorithm 2; similar experiments with Interleaved ToV produced comparable results, so we report only Parallel ToV here. Algorithm 2 does not specify the method to select data on the basis of scores. In Figure 11 we compare SCORE-ONLY and SCORE+RANDOM (already introduced above) with a third one RANDOM-FROM-TOP that selects at random from the top 50% subset of data with highest scores.

The RANDOM-FROM-TOP method is included only for these synthetic logistic-regression experiments, for theoretical interest as by construction, the training pool contains half of its examples from the target distribution.

The base set used for initial training contains $|U| = 4 \times 1024$ examples.

For both the scoring model and the final model training, we used 4 epochs of batch gradient descent with a linear decay learning rate scheduler. The initial learning rate was set to 0.5. We used $\epsilon = \frac{1}{10}$ for adjusting the learning rate on validation examples.

The selected subset size n was varied from 128 to 8192 in multiplicative steps of $\sqrt{2}$. All results are averaged over 10 independent runs. The final performance curves are presented in Figure 11.

E PROOF OF PROPOSITION 1

Throughout this proof, we denote by C a generic constant that can depend on c_0, C_1, M, C_η and whose value is allowed to change from line to line.

Letting $\Delta_s(i) = \hat{\theta}_s^{\text{bas}+i} - \hat{\theta}_s^{\text{bas}}$, Eq. (11) yields

$$\begin{aligned} \Delta_{k+1}(i) &= \Delta_k(i) - \eta m \nabla^2 \widehat{R}_U(\hat{\theta}_k^{\text{bas}}) \Delta_k(i) - \eta \nabla \ell(\hat{\theta}_k^{\text{bas}}; \mathbf{x}_i) + \text{err}_k(i) \\ &= \mathbf{H}_k \Delta_k(i) - \eta \nabla \ell(\hat{\theta}_k^{\text{bas}}; \mathbf{x}_i) + \text{err}_k(i), \end{aligned} \quad (20)$$

where \mathbf{H}_k is defined as in Eq. (14) and

$$\text{err}_k(i) := -\eta [\nabla \ell(\hat{\theta}_k^{\text{bas}+i}; \mathbf{x}_i) - \nabla \ell(\hat{\theta}_k^{\text{bas}}; \mathbf{x}_i)] - \eta m \int_0^1 [\nabla^2 \widehat{R}_U(\bar{\theta}_k(z)) - \nabla^2 \widehat{R}_U(\hat{\theta}_k^{\text{bas}})] \Delta_k(i) dz,$$

where $\bar{\theta}_k(z) = (1-z)\hat{\theta}_k^{\text{bas}} + z\hat{\theta}_k^{\text{bas}+i}$. By assumption $\theta \mapsto \nabla \ell(\theta; \mathbf{x}_i)$ and $\theta \mapsto \nabla^2 \widehat{R}_U(\theta)$ are M -Lipschitz, whence

$$\begin{aligned} \|\text{err}_k(i)\| &\leq \eta M \|\hat{\theta}_k^{\text{bas}+i} - \hat{\theta}_k^{\text{bas}}\| + \eta m M \|\hat{\theta}_k^{\text{bas}+i} - \hat{\theta}_k^{\text{bas}}\| \|\Delta_k(i)\| \\ &= \eta M \|\Delta_k(i)\| + \eta m M \|\Delta_k(i)\|^2. \end{aligned} \quad (21)$$

Define $\Delta_k^{\text{lin}}(i)$ by letting $\Delta_k^{\text{lin}}(i) = \mathbf{0}$ and, for $k \geq 0$,

$$\Delta_{k+1}^{\text{lin}}(i) = \mathbf{H}_k \Delta_k^{\text{lin}}(i) - \eta \nabla \ell(\hat{\theta}_k^{\text{bas}}; \mathbf{x}_i). \quad (22)$$

Comparing with Eq. (20), we obtain

$$\begin{aligned} (\Delta_{k+1}(i) - \Delta_{k+1}^{\text{lin}}(i)) &= \mathbf{H}_k(\Delta_{k+1}(i) - \Delta_{k+1}^{\text{lin}}(i)) + \text{err}_k(\eta, m), \\ \Rightarrow \Delta_t(i) - \Delta_t^{\text{lin}}(i) &= \sum_{s=0}^{t-1} \mathbf{M}_{t,s+1} \text{err}_s(\eta, m). \end{aligned} \quad (23)$$

Since $\nabla^2 \widehat{R}_U(\hat{\theta}_k^{\text{bas}}) \succeq c_0 \mathbf{I}_d$, we have $\|\mathbf{H}_k\|_{\text{op}} \leq (1 - c_0 m \eta)$, and therefore

$$\begin{aligned} \|\Delta_t(i) - \Delta_t^{\text{lin}}(i)\| &\leq \sum_{s=0}^{t-1} \|\mathbf{M}_{t,s+1}\|_{\text{op}} \|\text{err}_s(\eta, m)\| \\ &\leq \sum_{s=0}^{t-1} (1 - c_0 m \eta)^{t-s-1} \|\text{err}_s(\eta, m)\|. \end{aligned} \quad (24)$$

Further, from Eq. (22), and using $\|\nabla \ell(\hat{\theta}_k^{\text{bas}}; \mathbf{x}_i)\| \leq C_1$, we get

$$\begin{aligned} \Delta_t^{\text{lin}}(i) &= -\eta \sum_{s=0}^{t-1} \mathbf{M}_{t,s+1} \nabla \ell(\hat{\theta}_s^{\text{bas}}; \mathbf{x}_i), \\ \Rightarrow \|\Delta_t^{\text{lin}}(i)\| &\leq C_1 \eta \sum_{s=0}^{t-1} (1 - c_0 m \eta)^{t-s-1} \leq \frac{C}{m}. \end{aligned} \quad (25)$$

Let $D_t(i) := \max_{s \leq t} \|\Delta_s(i)\|$, $E_t(i) := \max_{s \leq t} \|\text{err}_s(i)\|$. Using Eqs. (21), (24) and (25), we get

$$\begin{aligned} D_t(i) &\leq \frac{C}{m} + \frac{1}{c_0 m \eta} E_{t-1}(i), \\ E_t(i) &\leq \eta M D_t(i) + \eta m M D_t(i)^2. \end{aligned}$$

Using these inequalities together, we obtain, for all $m \geq m_0$ (and eventually adjusting the constant C)

$$D_t(i) \leq \frac{C}{m}, \quad E_t(i) \leq \frac{C \eta}{m}, \quad (26)$$

whence, using again Eq. (24), we get

$$\|\Delta_t(i) - \Delta_t^{\text{lin}}(i)\| \leq \frac{C}{m^2}. \quad (27)$$

Notice that we can rewrite

$$S_i^{\text{lin}} = -\frac{1}{L} \sum_{s=1}^L \langle \nabla \widehat{R}_{\text{val}}(\hat{\theta}_s^{\text{bas}}), \Delta_s^{\text{lin}}(i) \rangle, \quad (28)$$

whence, using the fact that $\|\nabla^2 \widehat{R}_{\text{val}}(\bar{\theta})\|_{\text{op}} \leq C$ for all $\bar{\theta} \in [\hat{\theta}_k^{\text{bas}}, \hat{\theta}_k^{\text{bas}+i}]$ (this follows from the assumed bound $\|\nabla^2 \widehat{R}_{\text{val}}(\hat{\theta}_k^{\text{bas}})\|_{\text{op}} \leq C_1$ and the Lipschitz property of $\theta \mapsto \nabla^2 \widehat{R}_{\text{val}}(\hat{\theta})$), we get

$$\begin{aligned} |S_i - S_i^{\text{lin}}| &\leq C \max_{s \leq L} \|\Delta_s(i)\|^2 + \frac{1}{L} \sum_{s=1}^L |\langle \nabla \widehat{R}_{\text{val}}(\hat{\theta}_s^{\text{bas}}), \Delta_s(i) - \Delta_s^{\text{lin}}(i) \rangle| \\ &\leq C \max_{s \leq L} \|\Delta_s(i)\|^2 + C \max_{s \leq L} \|\Delta_s(i) - \Delta_s^{\text{lin}}(i)\| \\ &\leq \frac{C}{m^2}, \end{aligned}$$

and this completes the proof.

F PROOF OF THEOREM 1

Throughout this proof, we denote by C a generic constant that can depend on c_0, C_1, M, C_η and whose value is allowed to change from line to line.

F.1 BOUND ON Υ_i

The iteration for $\hat{\theta}_k^{\text{bas}}$ and $\hat{\theta}_k^{\text{bas},+}$, as specified by Algorithm 2, reads

$$\hat{\theta}_{k+1}^{\text{bas}} = \hat{\theta}_k^{\text{bas}} - \eta m \nabla \widehat{R}_U(\hat{\theta}_k^{\text{bas}}), \quad (29)$$

$$\hat{\theta}_{k+1}^{\text{bas},+} = \hat{\theta}_{0,k+1}^{\text{bas},+} - \varepsilon \eta m_{\text{val}} \nabla \widehat{R}_{\text{val}}(\hat{\theta}_{0,k+1}^{\text{bas},+}), \quad \hat{\theta}_{0,k+1}^{\text{bas},+} = \hat{\theta}_k^{\text{bas},+} - \eta m \nabla \widehat{R}_U(\hat{\theta}_k^{\text{bas},+}). \quad (30)$$

Letting $\Delta_k := \hat{\theta}_k^{\text{bas},+} - \hat{\theta}_k^{\text{bas}}$, and $\Delta_{0,k} := \hat{\theta}_{0,k}^{\text{bas},+} - \hat{\theta}_k^{\text{bas}}$, we obtain

$$\Delta_{k+1} = \Delta_{0,k+1} - \varepsilon \eta m_{\text{val}} \nabla \widehat{R}_{\text{val}}(\hat{\theta}_{k+1}^{\text{bas}}) + \text{err}_{k+1}^{(1)}, \quad (31)$$

$$\Delta_{0,k+1} = \mathbf{H}_k \Delta_k + \text{err}_k^{(2)}. \quad (32)$$

where, letting $\bar{\theta}_{0,k+1}(z) = (1-z)\hat{\theta}_{k+1}^{\text{bas}} + z\hat{\theta}_{0,k+1}^{\text{bas},+}$ and $\bar{\theta}_k(z) = (1-z)\hat{\theta}_k^{\text{bas}} + z\hat{\theta}_k^{\text{bas},+}$, we have

$$\begin{aligned} \text{err}_{k+1}^{(1)} &:= -\varepsilon \eta m_{\text{val}} \int_0^1 \nabla^2 \widehat{R}_{\text{val}}(\bar{\theta}_{0,k+1}(z)) \Delta_{0,k+1} dz, \\ \text{err}_k^{(2)} &:= -\eta m \int_0^1 [\nabla^2 \widehat{R}_U(\bar{\theta}_k(z)) - \nabla^2 \widehat{R}_U(\hat{\theta}_k^{\text{bas}})] \Delta_k dz. \end{aligned}$$

Using the assumption that $\|\nabla^2 \widehat{R}_{\text{val}}(\hat{\theta}_{k+1}^{\text{bas}}(z))\|_{\text{op}} \leq C$ and $\theta \mapsto \nabla^2 \widehat{R}_{\text{val}}(\theta)$ is M -Lipschitz, we get:

$$\|\text{err}_{k+1}^{(1)}\| \leq C \varepsilon \eta m_{\text{val}} \{ \|\Delta_{0,k+1}\| + \|\Delta_{0,k+1}\|^2 \}. \quad (33)$$

On the other hand, since $\theta \mapsto \nabla^2 \widehat{R}_U(\theta)$ is also M -Lipschitz, we have

$$\|\text{err}_k^{(2)}\| \leq \eta m M \|\Delta_k\|^2, \quad (34)$$

whence, using Eq. (32) and $\|\mathbf{H}_k\|_{\text{op}} \leq 1$

$$\begin{aligned} \|\Delta_{0,k+1}\| &\leq \|\Delta_k\| + \eta m M \|\Delta_k\|^2 \\ \Rightarrow \|\text{err}_{k+1}^{(1)}\| &\leq C \varepsilon \eta m_{\text{val}} \{ \|\Delta_k\| + \|\Delta_k\|^2 + \eta^2 m^2 \|\Delta_k\|^4 \}, \end{aligned} \quad (35)$$

where in the last line we used the assumption that $\eta m \leq C_\eta$.

Substituting Eqs. (34) and (35) in Eq. (31), (32), we obtain (using again $\eta m \leq C_\eta$)

$$\Delta_{k+1} = \mathbf{H}_k \Delta_k - \varepsilon \eta m_{\text{val}} \nabla \widehat{R}_{\text{val}}(\hat{\theta}_{k+1}^{\text{bas}}) + \text{err}_k, \quad (36)$$

$$\|\text{err}_k\| \leq C \eta \varepsilon m_{\text{val}} (\|\Delta_k\| + \|\Delta_k\|^4) + C \eta m \|\Delta_k\|^2. \quad (37)$$

We define $\Delta_k^{\text{lin}} = \mathbf{0}$ and, for $k \geq 0$,

$$\Delta_{k+1}^{\text{lin}} = \mathbf{H}_k \Delta_k^{\text{lin}} - \varepsilon \eta m_{\text{val}} \nabla \widehat{R}_{\text{val}}(\hat{\theta}_{k+1}^{\text{bas}}), \quad (38)$$

whence

$$\Delta_t^{\text{lin}} = -\varepsilon \eta m_{\text{val}} \sum_{s=0}^{t-1} \mathbf{M}_{t,s+1} \nabla \widehat{R}_{\text{val}}(\hat{\theta}_{s+1}^{\text{bas}}), \quad \Delta_t - \Delta_t^{\text{lin}} = \sum_{s=0}^{t-1} \mathbf{M}_{t,s+1} \text{err}_s. \quad (39)$$

Define $D_t := \max_{s \leq t} \|\Delta_s\|$, $E_t := \max_{s \leq t} \|\text{err}_s\|$. Using the fact that $\|\mathbf{M}_{t,s+1}\|_{\text{op}} \leq (1 - c_0 m \eta)^{t-s-1}$ and the assumption $\|\nabla \widehat{R}_{\text{val}}(\hat{\theta}_k^{\text{bas}})\| \leq C_1$, we get, from Eqs. (37), (39),

$$D_{t+1} \leq \frac{C \varepsilon m_{\text{val}}}{m} + \frac{C}{m \eta} E_t, \quad (40)$$

$$E_t \leq C \eta \varepsilon m_{\text{val}} (D_t + D_t^4) + C \eta m D_t^2, \quad (41)$$

Using the assumption $\varepsilon m_{\text{val}}/m \leq c_*$, this is easily seen to imply

$$D_t \leq C \frac{\varepsilon m_{\text{val}}}{m}, \quad E_t \leq C \frac{(\varepsilon m_{\text{val}})^2}{m} \eta. \quad (42)$$

Substituting in Eq. (39), we get

$$\|\Delta_t - \Delta_t^{\text{lin}}\| \leq \sum_{s=0}^{t-1} (1 - c_0 m \eta)^{t-s-1} \|\text{err}_s\| \leq C \left(\frac{\varepsilon m_{\text{val}}}{m} \right)^2. \quad (43)$$

The linearized score of Eq. (16) can be rewritten as

$$\Upsilon_i^{\text{lin}} \varepsilon = -\frac{1}{L} \sum_{s=1}^L \langle \nabla \ell(\hat{\theta}_s^{\text{bas}}; \mathbf{x}_i), \Delta_s^{\text{lin}} \rangle. \quad (44)$$

Using the fact that $\|\nabla \ell(\hat{\theta}_k^{\text{bas}}; \mathbf{x}_i)\|, \|\nabla \ell(\hat{\theta}_k^{\text{bas}}; \mathbf{x}_i)\|_{\text{op}} \leq C_1$, we get

$$|\Upsilon_i(\varepsilon) - \Upsilon_i^{\text{lin}} \varepsilon| \leq \frac{1}{L} \sum_{s=1}^L \left| \ell(\hat{\theta}_k^{\text{bas}}; \mathbf{x}_i) - \ell(\hat{\theta}_k^{\text{bas},+}; \mathbf{x}_i) + \langle \nabla \ell(\hat{\theta}_s^{\text{bas}}; \mathbf{x}_i), \Delta_s^{\text{lin}} \rangle \right| \quad (45)$$

$$\leq \frac{C}{L} \sum_{s=1}^L \|\Delta_s\|^2 + \frac{C}{L} \sum_{s=1}^L \|\Delta_s - \Delta_s^{\text{lin}}\| \quad (46)$$

$$\leq C \left(\frac{\varepsilon m_{\text{val}}}{m} \right)^2, \quad (47)$$

F.2 BOUND ON ϕ_i

The iteration for $\hat{\theta}_k^{\text{bas}}$ and $\hat{\theta}_k^{\text{bas},+}$, as specified by Algorithm 2, reads

$$\hat{\theta}_{k+1}^{\text{bas}} = \hat{\theta}_k^{\text{bas}} - \eta m \nabla \widehat{R}_U(\hat{\theta}_k^{\text{bas}}), \quad (48)$$

$$\hat{\theta}_{k+1}^{\text{val}} = \hat{\theta}_{k+1}^{\text{bas}} - \varepsilon \eta m_{\text{val}} \nabla \widehat{R}_{\text{val}}(\hat{\theta}_{k+1}^{\text{bas}}). \quad (49)$$

Hence, we can rewrite

$$\phi_i^{\text{lin}} \varepsilon = -\frac{1}{L} \sum_{s=1}^L \langle \nabla \ell(\hat{\theta}_s^{\text{bas}}; \mathbf{x}_i), \hat{\theta}_s^{\text{val}} - \hat{\theta}_s^{\text{bas}} \rangle.$$

Using the assumptions $\|\nabla^2 \ell(\hat{\theta}_s^{\text{bas}}; \mathbf{x}_i)\|_{\text{op}} \leq C_1, \|\widehat{R}_{\text{val}}(\hat{\theta}_k^{\text{bas}})\| \leq C_1$, we obtain

$$|\phi_i(\varepsilon) - \phi_i^{\text{lin}} \varepsilon| \leq \frac{1}{L} \sum_{s=1}^L \left| \ell(\hat{\theta}_s^{\text{val}}; \mathbf{x}_i) - \ell(\hat{\theta}_s^{\text{bas}}; \mathbf{x}_i) - \langle \nabla \ell(\hat{\theta}_s^{\text{bas}}; \mathbf{x}_i), \hat{\theta}_s^{\text{val}} - \hat{\theta}_s^{\text{bas}} \rangle \right| \quad (50)$$

$$\leq \frac{C}{L} \sum_{s=1}^L \|\hat{\theta}_s^{\text{val}} - \hat{\theta}_s^{\text{bas}}\|^2 \leq C(\varepsilon \eta m_{\text{val}})^2. \quad (51)$$

The claim thus follows by recalling that $\eta \leq C_\eta/m$.

G PROOF OF THEOREM 2

To lighten notation, we define $\mathbf{r}_k := \nabla \widehat{R}_{\text{val}}(\hat{\theta}_k^{\text{bas}})$ and $\mathbf{v}_k(i) := \nabla \ell(\hat{\theta}_k^{\text{bas}}; \mathbf{x}_i)$.

For any $L, L_1 \in \mathbb{Z}$, we have

$$\begin{aligned}
\Upsilon_i^{\text{lin}}(L) &= \Upsilon_i^{\text{lin},0}(L) + \Upsilon_i^{\text{lin},1}(L) + \Upsilon_i^{\text{lin},2}(L) + \Upsilon_i^{\text{lin},3}(L), \\
\Upsilon_i^{\text{lin},0}(L) &:= \frac{\eta m_{\text{val}}}{L} \sum_{0 \leq t < s \leq L_0} \langle \mathbf{r}_{t+1}, \mathbf{M}_{s,t+1}^\top \mathbf{v}_s(i) \rangle, \\
\Upsilon_i^{\text{lin},1}(L) &:= \frac{\eta m_{\text{val}}}{L} \sum_{0 \leq t \leq L_0, L_0 < s \leq L} \langle \mathbf{r}_{t+1}, \mathbf{M}_{s,t+1}^\top \mathbf{v}_s(i) \rangle, \\
\Upsilon_i^{\text{lin},2}(L) &:= \frac{\eta m_{\text{val}}}{L} \sum_{L_0 < t < s \leq L_0: |s-t| \geq L_1} \langle \mathbf{r}_{t+1}, \mathbf{M}_{s,t+1}^\top \mathbf{v}_s(i) \rangle, \\
\Upsilon_i^{\text{lin},2}(L) &:= \frac{\eta m_{\text{val}}}{L} \sum_{L_0 < t < s \leq L_0: |s-t| \geq L_1} \langle \mathbf{r}_{t+1}, \mathbf{M}_{s,t+1}^\top \mathbf{v}_s(i) \rangle, \\
\Upsilon_i^{\text{lin},3}(L) &:= \frac{\eta m_{\text{val}}}{L} \sum_{L_0 < t < s \leq L_0: |s-t| < L_1} \langle \mathbf{r}_{t+1}, \mathbf{M}_{s,t+1}^\top \mathbf{v}_s(i) \rangle.
\end{aligned}$$

Since by continuity we have $\lim_{k \rightarrow \infty} \nabla^2 \widehat{R}_U(\hat{\theta}_k^{\text{bas}}) = \mathbf{Q}_\infty$, for any $\delta \in (0, 1/2)$, we can choose L_0 large enough so that $(1 - \delta)\mathbf{Q}_\infty \preceq \nabla^2 \widehat{R}_U(\hat{\theta}_k^{\text{bas}}) \preceq (1 + \delta)\mathbf{Q}_\infty$ for all $k > L_0$. In particular there exists $c_0 > 0$ (independent of ε) such that $\|\mathbf{H}_k\|_{\text{op}} \leq (1 - c_0 m \eta)$ for all $k > L_0$.

Clearly $|\Upsilon_i^{\text{lin},0}(L)| \leq C(L_0)/L \rightarrow 0$ as $L \rightarrow \infty$. Further

$$\begin{aligned}
|\Upsilon_i^{\text{lin},1}(L)| &\leq \frac{C\eta m_{\text{val}}}{L} \sum_{0 \leq t \leq L_0, L_0 < s \leq L} \|\mathbf{M}_{s,t+1}\|_{\text{op}} \\
&\leq \frac{C\eta m_{\text{val}}}{L} \sum_{0 \leq t \leq L_0, L_0 < s \leq L} (1 - c_0 m \eta)^{s-t-1} \\
&\leq \frac{C\eta m_{\text{val}}}{L} \frac{L_0}{c_0 m \eta} \rightarrow 0.
\end{aligned}$$

Finally, by increasing L_0 , we can ensure that, for $k > L_0$, $\|\mathbf{H}_k - \mathbf{H}_\infty\|_{\text{op}} \leq \delta$, $\|\mathbf{r}_k - \mathbf{r}_\infty\| \leq \delta$, $\|\mathbf{v}_k(i) - \mathbf{v}_\infty(i)\| \leq \delta$ (where $\mathbf{H}_\infty = \mathbf{I} - \eta m \mathbf{Q}_\infty$ and $\mathbf{r}_\infty, \mathbf{v}_\infty(i)$). Hence

$$|\langle \mathbf{r}_{t+1}, \mathbf{M}_{s,t+1}^\top \mathbf{v}_s(i) \rangle - \langle \mathbf{r}_\infty, \mathbf{H}_\infty^{s-t-1} \mathbf{v}_\infty(i) \rangle| \leq C|t-s+1|(1 - c_0 m \eta)^{s-t-1} \delta.$$

Therefore, letting

$$\tilde{\Upsilon}_i^{\text{lin},2}(L) := \frac{\eta m_{\text{val}}}{L} \sum_{L_0 < t < s \leq L} \langle \mathbf{r}_\infty, \mathbf{H}_\infty^{s-t-1} \mathbf{v}_\infty(i) \rangle, \quad (52)$$

we have

$$\begin{aligned}
|\Upsilon_i^{\text{lin},2}(L) - \tilde{\Upsilon}_i^{\text{lin},2}(L)| &\leq \frac{\eta m_{\text{val}}}{L} \sum_{L_0 < t < s \leq L} |\langle \mathbf{r}_{t+1}, \mathbf{M}_{s,t+1}^\top \mathbf{v}_s(i) \rangle - \langle \mathbf{r}_\infty, \mathbf{H}_\infty^{s-t-1} \mathbf{v}_\infty(i) \rangle| \\
&\leq \frac{\eta m_{\text{val}}}{L} \sum_{L_0 < t < s \leq L} C|t-s+1|(1 - c_0 m \eta)^{s-t-1} \delta \\
&\leq \eta m_{\text{val}} \cdot \frac{1}{(c_0 m \eta)^2} \delta.
\end{aligned}$$

Finally, using again $|\langle \mathbf{r}_\infty, \mathbf{H}_\infty^{s-t-1} \mathbf{v}_\infty(i) \rangle| \leq (1 - c_0 m \eta)^{s-t-1}$, we have

$$\begin{aligned}
\lim_{L \rightarrow \infty} \tilde{\Upsilon}_i^{\text{lin},2}(L) &= \lim_{L \rightarrow \infty} \frac{\eta m_{\text{val}}}{L} \sum_{L_0 < t \leq L} \sum_{s=t+1}^{\infty} \langle \mathbf{r}_\infty, \mathbf{H}_\infty^{s-t-1} \mathbf{v}_\infty(i) \rangle \\
&= \eta m_{\text{val}} \sum_{k=0}^{\infty} \langle \mathbf{r}_\infty, \mathbf{H}_\infty^k \mathbf{v}_\infty(i) \rangle \\
&= \eta m_{\text{val}} \langle \mathbf{r}_\infty, (\mathbf{I} - \mathbf{H}_\infty)^{-1} \mathbf{v}_\infty(i) \rangle \\
&= \frac{m_{\text{val}}}{m} \langle \mathbf{r}_\infty, \mathbf{Q}_\infty^{-1} \mathbf{v}_\infty(i) \rangle.
\end{aligned}$$

This finishes the proof of the part of Eq. (18) which concerns the limit of Υ_i^{lin} . The calculation of $\lim_{L \rightarrow \infty} S_i^{\text{lin}}(L)$ is completely analogous and we omit it.

H PROOF OF THEOREM 3

To simplify notations, we write $y_j = y(\mathbf{x}_j)$ for the response variables and $\boldsymbol{\psi}_j = \boldsymbol{\psi}(\mathbf{x}_j)$ for the feature vectors. Similarly, for the $y_j^{\text{val}} = y(\mathbf{z}_j^{\text{val}})$, $\boldsymbol{\psi}(\mathbf{z}_j^{\text{val}}) = \boldsymbol{\psi}_j^{\text{val}}$.

With these notations, we have $\mathbf{H}_k = \mathbf{H}$ independent of k and

$$\nabla \ell(\hat{\boldsymbol{\theta}}; \mathbf{x}_i) = -(y_i - \langle \boldsymbol{\psi}_i, \boldsymbol{\theta} \rangle) \boldsymbol{\psi}_i, \quad (53)$$

$$\nabla \hat{R}_{\text{val}}(\boldsymbol{\theta}) = -\frac{1}{m} \boldsymbol{\Psi}^\top (\mathbf{y} - \boldsymbol{\Psi} \boldsymbol{\theta}), \quad (54)$$

$$\mathbf{H} = \mathbf{I} - \eta \boldsymbol{\Psi}^\top \boldsymbol{\Psi}. \quad (55)$$

Hence,

$$\Upsilon_i^{\text{lin}} = \frac{\eta m_{\text{val}}}{L} \sum_{0 \leq t < s \leq L} r_s(i) \langle \boldsymbol{\Psi} \mathbf{r}_{t+1}^{\text{val}}, \mathbf{H}^{s-t-1} \boldsymbol{\psi}_i \rangle, \quad (56)$$

$$\mathbf{r}_t^{\text{val}} := \mathbf{y}^{\text{val}} - \boldsymbol{\psi}^{\text{val}} \hat{\boldsymbol{\theta}}_{t+1}^{\text{bas}}, \quad r_s(i) := y_i - \langle \boldsymbol{\psi}_i, \hat{\boldsymbol{\theta}}_s^{\text{bas}} \rangle \quad (57)$$

Since \mathbf{P}_Ψ is the projector onto the null-space of \mathbf{H} , and by our choice of η , we have $\mathbf{H} = \mathbf{P}_\Psi + \mathbf{H}_\perp$, where the row/column space of \mathbf{H}_\perp is orthogonal to the one of \mathbf{P}_Ψ and $\|\mathbf{H}_\perp\|_{\text{op}} = (1 - c_\psi \eta) \in [0, 1)$. As a consequence $\|\mathbf{H}^{s-t-1} - \mathbf{P}_\Psi\|_{\text{op}} \leq (1 - c_\psi \eta)^{s-t-1}$.

Define

$$\tilde{\Upsilon}_i^{\text{lin}} := \frac{\eta m_{\text{val}}}{L} \sum_{0 \leq t < s \leq L} r_s(i) \langle \boldsymbol{\Psi} \mathbf{r}_{t+1}^{\text{val}}, \mathbf{P}_\Psi \boldsymbol{\psi}_i \rangle. \quad (58)$$

Then we have

$$\begin{aligned} \left| \frac{1}{L} \Upsilon_i^{\text{lin}} - \frac{1}{L} \tilde{\Upsilon}_i^{\text{lin}} \right| &\leq \frac{\eta m_{\text{val}}}{L^2} \sum_{0 \leq t < s \leq L} \left| r_s(i) \langle \boldsymbol{\Psi} \mathbf{r}_{t+1}^{\text{val}}, \mathbf{P}_\Psi \boldsymbol{\psi}_i \rangle \right| \\ &\leq \frac{\eta m_{\text{val}}}{L^2} \sum_{0 \leq t < s \leq L} |r_s(i)| \|\boldsymbol{\Psi} \mathbf{r}_{t+1}^{\text{val}}\| \|\mathbf{H}^{s-t-1} - \mathbf{P}_\Psi\|_{\text{op}} \|\boldsymbol{\psi}_i\| \\ &\stackrel{(a)}{\leq} C \frac{\eta m_{\text{val}}}{L^2} \sum_{0 \leq t < s \leq L} (1 - c_\psi \eta)^{s-t-1} \\ &\leq C \frac{\eta m_{\text{val}}}{L} \frac{1}{c_\psi \eta} \xrightarrow{L \rightarrow \infty} 0. \end{aligned}$$

In (a), we used the fact that $\lim_{t \rightarrow \infty} \hat{\boldsymbol{\theta}}_t^{\text{bas}} = \hat{\boldsymbol{\theta}}$ (Bartlett et al., 2021), and therefore $|r_s(i)|$, $\|\boldsymbol{\Psi} \mathbf{r}_{t+1}^{\text{val}}\|$ remain bounded as $s, t \rightarrow \infty$.

In view of the above, $\lim_{L \rightarrow \infty} \Upsilon_i^{\text{lin}}/L = \lim_{L \rightarrow \infty} \tilde{\Upsilon}_i^{\text{lin}}/L$. For the latter, we have

$$\begin{aligned} \lim_{L \rightarrow \infty} \frac{1}{L} \tilde{\Upsilon}_i^{\text{lin}} &= \lim_{L \rightarrow \infty} \frac{\eta m_{\text{val}}}{L^2} \sum_{0 \leq t < s \leq L} r_s(i) \langle \boldsymbol{\Psi} \mathbf{r}_{t+1}^{\text{val}}, \mathbf{P}_\Psi \boldsymbol{\psi}_i \rangle \\ &= \lim_{L \rightarrow \infty} \frac{\eta m_{\text{val}}}{L^2} \sum_{L_0 \leq t < s \leq L} r_s(i) \langle \boldsymbol{\Psi} \mathbf{r}_{t+1}^{\text{val}}, \mathbf{P}_\Psi \boldsymbol{\psi}_i \rangle \\ &= \lim_{L \rightarrow \infty} \frac{\eta m_{\text{val}}}{L^2} \sum_{L_0 \leq t < s \leq L} r(i) \langle \boldsymbol{\Psi} \mathbf{r}^{\text{val}}, \mathbf{P}_\Psi \boldsymbol{\psi}_i \rangle + \text{err}(L_0, L), \end{aligned}$$

where

$$|\text{err}(L_0, L)| \leq C \sup_{s \geq L_0} |r_s(i) - r(i)| + C \sup_{t \geq L_0} \|\mathbf{r}_t^{\text{val}} - \mathbf{r}_{t+1}^{\text{val}}\|. \quad (59)$$

Since $\lim_{t \rightarrow \infty} \hat{\theta}_t^{\text{bas}} = \hat{\theta}$, we have $\lim_{L_0 \rightarrow \infty} \limsup_{L \rightarrow \infty} \text{err}(L_0, L) = 0$. Therefore,

$$\begin{aligned} \lim_{L \rightarrow \infty} \frac{1}{L} \tilde{\Upsilon}_i^{\text{lin}} &= \lim_{L_0 \rightarrow \infty} \lim_{L \rightarrow \infty} \frac{\eta m_{\text{val}}}{L^2} \sum_{L_0 \leq t < s \leq L} r(i) \langle \Psi \mathbf{r}^{\text{val}}, \mathbf{P}_{\Psi} \psi_i \rangle \\ &= \frac{1}{2} \eta m_{\text{val}} r(i) \langle \Psi \mathbf{r}^{\text{val}}, \mathbf{P}_{\Psi} \psi_i \rangle. \end{aligned}$$

This proves the limit for $\Upsilon_i^{\text{lin}}(L)$ in Eq. (19).

The limit of $S_i^{\text{lin}}(L)$ is computed essentially by the same argument and we omit the derivation.

I GRADIENT ACCUMULATION IN TOV

A potential concern in Train-on-Validation (ToV) scoring is that repeated updates on the validation set across scoring checkpoints could accumulate on the same model parameters and lead to overfitting of the selection policy to \mathbf{Z}^{val} . In Interleaved ToV, however, each scoring model is initialized from a base-model checkpoint that has not previously been updated on the validation set and is then updated on \mathbf{Z}^{val} only once. Thus, validation gradients do not accumulate on the same parameters across epochs. In Parallel ToV, validation-updated models continue training and may, in principle, accumulate such effects; in practice, this is mitigated by using a much smaller learning rate for validation training. Nevertheless, Parallel ToV may remain more susceptible to validation overfitting, which could contribute to its slightly inferior performance relative to Interleaved ToV.

J LIMITATION

The core idea of “train-on-validation” impacting training examples is general, but the specific scoring function $F(\cdot)$ and aggregation strategy might need adaptation for different problem settings.

The SCORE+RANDOM selection strategy often outperformed SCORE-ONLY in our experiments, suggesting that diversity plays an important role beyond simply selecting the “most affected” examples. While this is a practical improvement, it also indicates that our current scoring mechanism might not fully capture the optimal diversity or coverage needed for effective generalization. It will be interesting to explore more sophisticated diversity-aware scoring or selection mechanisms that explicitly balance our scoring methods with representation across the data space.

Although we mitigated bias toward shorter examples through length-based binning, a more refined length-normalization or task-specific weighting might further enhance the selection process. Furthermore, it will be interesting to see if the performance of our strategies improves further compared to random selection when the learning rate is also tuned for these strategies, rather than just for random selection.

Finally, our theoretical analysis relies on stylized settings that are plausible for simple models but may not hold in many large-scale applications.

K MODELS AND DATASETS INFORMATION

K.1 DATASET INFORMATION

- Slim Orca:
 - [Link](#)
 - Citations: [Longpre et al. \(2023\)](#); [Mukherjee et al. \(2023\)](#); [Lian et al. \(2023\)](#)
 - Licence: mit
- Alpaca GPT-4:
 - Paper: [Peng et al. \(2023\)](#)
 - [Repository](#)
 - [Link](#)

- Licence: cc-by-nc-4.0
- Alpaca GPT-3.5:
 - Paper: [Taori et al. \(2023\)](#)
 - [Link](#)
 - Licence: cc-by-nc-4.0
- Multinerd:
 - Paper: [Tedeschi & Navigli \(2022\)](#)
 - [Link](#)
 - Licence: cc-by-nc-sa-4.0
- Ai4p:
 - [Link](#)
 - Licence: [link](#)
- C4 dataset:
 - [Link](#)
 - Labeled for NER task using llms.
 - Licence: [terms of use](#)
- Syn-Big:
 - Synthetically generated by us using llms.
 - Proprietary dataset

K.2 PRETRAINED MODEL INFORMATION

- Meta-Llama-3-8B [AI@Meta \(2024\)](#)
 - [Link](#)
 - License: llama3
- xlm-roberta-base [Conneau et al. \(2019\)](#)
 - [Link](#)
 - License: mit

K.3 COMPUTE RESOURCES INFORMATION

All experiments were conducted on a single machine equipped with 8 NVIDIA A100 (80 GB) GPUs and 128 AMD EPYC 9354 (32-core) CPUs.

Instruction-tuning experiments across all five setups, with 10 runs per setup and including hyperparameter search, required approximately two weeks of wall-clock time using parallel execution across GPUs. NER experiments completed within 4–5 days under the same setup. Logistic regression experiments were lightweight and required less than 30 minutes in total.