DISENTANGLING THE ROLES OF REPRESENTATION AND SELECTION IN DATA PRUNING (FOR FINE-TUNING)

Anonymous authors

Paper under double-blind review

ABSTRACT

Data pruning, the process of carefully selecting a small subset of training data, has been shown to improve both training efficiency and performance. It typically involves two steps: (1) obtaining a representation for each instance, and (2) applying a selection algorithm using these representations. However, the distinct roles of these two steps, as well as their interactions, remain unclear. To address this, we conduct a systematic study of data pruning, focusing on NLP fine-tuning. Our theoretical and empirical findings reveal that data representation often plays a more fundamental role than the selection algorithm: gradients, despite being computationally expensive, provide stronger pruning signals than other representations, making gradient-based methods consistently outperform cheaper alternatives. We also demonstrate that different selection algorithms excel in specific scenarios but are heavily influenced by the chosen representation. These insights provide clear guidelines for future research and practical applications.

023 024 025

026

044

045

046

004

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

The remarkable success of modern deep learning is largely driven by vast training datasets (Kaplan et al., 2020; Hoffmann et al., 2022; Sardana et al., 2024). However, scaling the size of datasets comes with great computational costs. Fortunately, recent studies have shown that by carefully selecting a small subset of the original large dataset, a process known as *data pruning*, it is possible to improve training efficiency while also improving generalization performance (Paul et al., 2021; Sorscher et al., 2022; Du et al., 2023; Xia et al., 2024). Moreover, effective pruning methods can offer insights into learning algorithms and the roles of training data (Koh & Liang, 2017; Ilyas et al., 2022).

Most data pruning methods involve two steps: First, obtaining a *representation* for each instance in
the original training set (e.g., hidden states from a pretrained language model); second, selecting
instances based on these representations given a *data budget* (e.g., 30% of the training set), according
to a *selection algorithm*. Despite the success of data pruning, *existing studies treat these two steps as a single process*, leaving fundamental questions about the roles of representations and selection
algorithms unanswered: *Which data representations and selection algorithms are most effective, and what is their suitability for different tasks? How do different representations impact the data points selected by each selection algorithm?* To address these questions, we conduct a thorough study of
data pruning through both theoretical and empirical lenses, focusing on NLP fine-tuning tasks. Our
contributions are as follows:

- 1. We offer a systematic overview of existing data pruning methods (§2). From this overview, we identify three common ways to create representations: based on *training dynamics, hid-den states*, and *gradients*. We also identify three common objectives of selection algorithms: maximizing *difficulty*, maximizing *diversity*, and maximizing *validation performance*.
- We both theoretically analyze the pruning signals contained in different representations (§3), and empirically validate that the representations of stronger signals, despite being computationally more expensive, are more effective (§4). Specifically, we find that *data pruning based on gradients often performs the best*. In contrast, data pruning based on hidden states, which is computationally cheaper, often does not perform better than random selection, especially with low data budgets. Our experiments evaluate six representative methods across common representations and selection algorithms, covering a wide range of

and summarization (text generation).

data budget to perform well, maximizing validation performance excels when train-test distributions differ, and that maximizing diversity in general does not work very well.
4. Through both interpretable experiments with synthetic data (§3) and common NLP tasks (§4), we show that with the same selection objective, different data representations can lead to drastically different selections of instances. For instance, when the aim is to select the most difficult instances, using hidden states results in selecting instances far from the decision boundary, while using gradients select instances close to it. Moreover, we show that the sensitivity of selection algorithms towards the change of representations vary, and

NLP tasks: hate speech detection (classification), commonsense reasoning (multiple choice),

3. Our empirical experiments (§4) further show that maximizing difficulty requires a higher

065 066

054

056

Our findings highlight the strengths and limitations of current data pruning methods, offering clear
guidelines for selecting appropriate representations and algorithms given specific tasks and constraints.
We recommend prioritizing the development of efficient gradient-based methods, due to their superior
performance and better interpretability (Pruthi et al., 2020; Park et al., 2023).

that data representations often matter more than selection algorithms.

071 072

073

2 OVERVIEW OF DATA PRUNING METHODS

Various data pruning methods have been proposed. They typically follow two steps for data selection:
(1) obtaining a representation for each instance, and (2) using an algorithm to select instances based
on their representations. Regarding representations, most methods rely on one of the three types:
training dynamics, hidden states, and gradients. Our overview is structured accordingly. While recent
studies have explored external large language models, by prompting and distillation, particularly in
instruction tuning (Sachdeva et al., 2024; Chen et al., 2024; Lu et al., 2024; Liu et al., 2024), we focus
on representations from the model that we are training. This allows us to analyze signals that directly
reflect its learning behavior.

Regarding selection algorithms, most methods are built on one or more of the following objectives:
 retaining the most difficult data instances, maximizing the diversity of selected data, and retaining the
 training instances that improve a model's performance on held-out data the most.

085

087

2.1 TRAINING-DYNAMIC-BASED METHODS

Most data pruning methods based on training dynamics aim to maximize the ratio of difficult instances being selected, where the difficulty is defined by heuristics like low prediction confidence and high 089 training loss. For example, Toneva et al. (2019) determine the difficulty of examples by the earliest 090 epoch after which that example is always correctly classified; Jiang et al. (2019) keep the examples 091 with the highest training loss; Swayamdipta et al. (2020) and Du et al. (2023) use the standard 092 deviation and mean of prediction probabilities of the correct class across different epochs; Paul et al. (2021) keep the examples with the largest error norm; Baldock et al. (2021) select examples that need 094 to pass through more layers before being correctly classified; and Marion et al. (2023) and Kwok 095 et al. (2024) quantify difficulty by perplexity. Moreover, Maini et al. (2022) show that examples that 096 are forgotten slower when training on a held-out subset are more rare and thus worth keeping. In the context of instruction tuning, Li et al. (2024) quantify the difficulty of a response by the ratio of its 097 losses between generations with and without its instruction. 098

Other selection objectives have also been explored. For instance, Mindermann et al. (2022) prioritize training on learnable, worth-learning, and not-yet-learned examples, quantified by the loss difference between the model itself, and a small reference model trained on a held-out dataset. Furthermore,
Yang et al. (2024) aim to diversify the training data by clustering instances based on their training losses across different epochs and sampling evenly from each cluster.

104 105

- 2.2 HIDDEN-STATE-BASED METHODS
- 107 Hidden-state-based methods often exploit similarities between data samples to realize various selection objectives. For example, Sorscher et al. (2022) first perform a *k*-means clustering of hidden

states, then select the ones that lie farther from their cluster centroids, because these instances are
less prototypical and therefore more likely to be difficult. Abbas et al. (2023) and Tirumala et al.
(2023) extend this approach by identifying and removing semantic duplication after clustering. This
de-duplication step helps balance data diversity and difficulty by discarding redundant instances.

112

113 2.3 GRADIENT-BASED METHODS

Another line of methods uses gradient information of each training instance to estimate its importance.
Specifically, most gradient-based methods estimate the impact of removing instances by simulating
re-training scenarios. These approaches use gradient-based tools like influence functions (Koh &
Liang, 2017; Pruthi et al., 2020) and datamodels (Ilyas et al., 2022; Park et al., 2023).

The most common objective is to maximize or maintain validation performance. For example, Xia et al. (2024) and Engstrom et al. (2024) discard training instances with low contribution towards the validation performance of the target task. They approximate this influence by using the TracIn influence function (Pruthi et al., 2020) and datamodels (Ilyas et al., 2022). Moreover, Killamsetty et al. (2021) and Yang et al. (2023) keep instances that likely result in similar models as when training on the full dataset.

Gradients are also often used to select difficult instances. For example, Feldman & Zhang (2020) demonstrate that data instances of high self-influence (i.e., training on itself contributes more to its prediction) are more difficult and help generalization. Using the TracIn influence function, this self-influence score can be estimated as the gradient norm. Similarly, Thakkar et al. (2023) prioritizes data instances of different self-influences at different stages of pretraining, to limit the influence of noisy data while focusing on higher-quality ones; and Bejan et al. (2023) uses automatic curriculum learning to filter noisy data.

- 132
- 133 134

3 A TALE OF DATA REPRESENTATIONS AND SELECTION ALGORITHMS

This section studies the distinct roles of data representations, selection algorithms, and their interactions. Specifically, we focus on six representative methods that span all three major types of data representations and selection objectives. First, we provide in-depth explanations of each method (§3.1). Next, we present a theoretical analysis of the signals each representation offers (§3.2). Finally, we conduct interpretable synthetic experiments to study how the data instances favored by each selection algorithm are shaped by the chosen data representations (§3.3).

141 142

Notation We now summarize the notation used in this paper. We denote the original training set with N instances as $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$. The selected subset of data is represented by $\mathcal{S} \subset \mathcal{D}$. We use B to denote the data budget (e.g., $B = |\mathcal{S}| = 0.2N$). For a data point (x_i, y_i) and a model with parameters θ , we use $f_{\theta}(x_i)$ to denote the model's logits, $\ell(f_{\theta}(x_i), y_i)$ to denote the loss, and $p_{\theta}(y_i|x_i)$ to denote the prediction probability of the *correct class* or token. Moreover, we use $h_{\theta}(x_i)$ to denote the last layer hidden state of input x_i , and $g_{\theta}(x_i, y_i) = \nabla_{\theta} \ell(f_{\theta}(x_i), y_i)$ to represent the gradient. Moreover, we train each model for T epochs, use $p_t(y_i|x_i)$ to denote the prediction probability of the correct class at epoch $t \leq T$, and use η_t to denote the learning rate at epoch t.

150 151

152

3.1 PRELIMINARIES: DETAILS OF DATA PRUNING METHODS

Data pruning typically involves two steps under a given data budget (e.g., 20% of the original dataset).
First, a *reference model* is used to obtain representations for each instance, such as hidden states from a pretrained model. Second, a selection algorithm chooses a subset based on these representations, e.g., the ones that are farthest from the clustering centroids, following a selection objective, e.g., selecting the most difficult instances. The selected instances are then used to a train new *main model* which is the final model of interest. We summarize the representations and selection objectives of these methods in Appendix B.

We first describe two training-dynamic-based methods, Hard-to-Learn and SmallToLarge (S2L).
 These training dynamics (e.g., training losses) are collected from a reference training run, during which a reference model is trained on the original dataset.

162 Hard-to-Learn The Hard-to-Learn method is based on a simple intuition: training instances that 163 are difficult for models to fit often contain fewer regularities, and these instances can improve model 164 generalization (Swayamdipta et al., 2020; Jiang et al., 2021). Concretely, it represents each training 165 instance by a score computed from the **training dynamics** of a reference model. In classification tasks, 166 for a given instance (x_i, y_i) , its score is defined as the average prediction probability of the correct label across different epochs, i.e., $s_{\text{hard}}(x_i, y_i) = \frac{1}{T} \sum_{t=1}^{T} p_t(y_i|x_i)$. Instances with the lowest scores are selected for training the main model, i.e., $S = \arg \min_{|S|=B} \sum_{(x_i, y_i) \in S} s_{\text{hard}}(x_i, y_i)$. To 167 168 169 extend this concept to generation tasks, Bhatnagar et al. (2022) and Ince et al. (2023) replace s_{hard} 170 with the inverse perplexity. Additionally, Jiang et al. (2021) show that s_{hard} correlates well with the 171 expected accuracy of an instance when it is excluded from the training data. This correlation also suggests that Hard-to-Learn instances contain fewer regularities, supporting the intuition. 172

173

SmallToLarge Also based on **training dynamics**, SmallToLarge (S2L) aims to select **diverse** instances. Specifically, Yang et al. (2024) observe that training instances with different loss trajectories across epochs likely require different knowledge. To make sure different skills are evenly represented in the subset S, S2L performs three steps. Initially, it represents each training instance by its *cross entropy loss trajectory*. Then, it performs a k-means clustering on these trajectories, and sorts these clusters by size in an ascending order. Finally, it iteratively samples (B - |S|)/(K - k + 1) instances from each cluster¹, where K is the number of clusters, and k is the current cluster index. This sample size choice helps S2L prioritize smaller clusters to increase the diversity of the selected data.

Next, we discuss two hidden-state-based methods, **Prototypicality** and **SemDeDup**. Hidden states are usually computed using a reference pre-trained model, and thus require no further training, making them computationally more efficient than other methods.

Prototypicality The Prototypicality method (Sorscher et al., 2022) aims to select difficult instances by exploiting their *similarities*: it measures difficulty based on how *prototypical* an instance is in the dataset. Specifically, after representing instances by their hidden states, prototypicality applies *k*-means clustering and ranks instances based on their distances to their cluster centroids. Instances with larger distances are considered less prototypical, more difficult, and selected for further training.

SemDeDup Building on Prototypicality, SemDeDup adds one more step to account for data
 diversity besides difficulty (Abbas et al., 2023). Concretely, after the clustering step of Prototypicality,
 SemDeDup identifies semantically duplicate pairs of instances within each cluster using cosine
 similarities of their hidden states. For each identified duplicate pair, it retains the instance that lies
 farther from the cluster centroid, thereby prioritizing diversity while maintaining difficulty.

Finally, we discuss two gradient-based methods, LESS and Memorization. They require training a reference model on the original dataset, and computing gradients using different checkpoints. This makes them very costly, because they require performing back-propagation with batch sizes of 1 to obtain the gradients, which have a high dimensionality, equal to the number of model parameters.

201 **LESS** Unlike the above methods that only rely on representations of the training data, LESS assumes the availability of a validation set. It aims to select instances that maximize the validation 202 performance, measured by validation loss reduction (Xia et al., 2024). Naively, this estimation 203 requires retraining models on a large number of random training subsets, which is prohibitively 204 expensive for modern neural models. Therefore, LESS employs influence functions for approxima-205 tion (Pruthi et al., 2020). Specifically, given the t-th reference model checkpoint θ_t , LESS estimates 206 the loss reduction of a training instance (x_i, y_i) w.r.t. a validation instance (x_{val}, y_{val}) by the dot 207 product of their (normalized) gradients. For multiple checkpoints, LESS performs weighted averag-ing by learning rates. Formally, $s_{\text{LESS}}(x_i, y_i) = \sum_{t=1}^{T} \eta_t g_{\theta_t}(x_i, y_i) \cdot g_{\theta_t}(x_{\text{val}}, y_{\text{val}})$. The instances with the highest scores are kept. Moreover, for efficiency, LESS uses LoRA (Hu et al., 2022) and 208 209 210 random projection (Park et al., 2023; Johnson & Lindenstrauss, 1984) for dimensionality reduction. 211 Additionally, LESS considers optimizer states when computing training gradients. 212

213 Memorization Feldman & Zhang (2020) define memorization in training as the loss increase of a
 214 training instance before and after it is removed from the training set, i.e., self-influence. Following a
 215

¹For simplicity, here we reload |S| as the number of already sampled instances.

similar intuition as Hard-to-Learn, Feldman (2020) argues that instances with high memorization scores are usually more **difficult** and thus contribute more to generalization. Similar to LESS, memorization is also usually approximated by influence functions using **gradients** (Thakkar et al., 2023; Bejan et al., 2023; Pruthi et al., 2020). In this work, following Pruthi et al. (2020), we compute the memorization score of a training instance (x_i, y_i) by $\sum_{t=1}^{T} \eta_t g_{\theta_t}(x_i, y_i) \cdot g_{\theta_t}(x_i, y_i)$. Inspired by Xia et al. (2024), we also adopt LoRA and random projection for efficiency.

3.2 WHAT DO DIFFERENT REPRESENTATIONS REVEAL?

Data representation is central to data pruning, yet it remains unclear *how different representations vary in the signals they retain.* This section offers theoretical insights into this question. Specifically, as described in §3.1, most data pruning methods rely on the similarity between instances based on their representations (e.g., S2L, Prototypicality, and LESS). We therefore ask: *How does the notion of similarity between two instances change across different representation spaces?* Our analyses build insights into the pruning signals from different representations, which we will further study in the following synthetic (§3.3) and NLP task (§4) experiments.

For clarity of analysis, we consider a binary classification task with labels $y \in \{-1, +1\}$, optimized 232 with binary cross-entropy loss. For the representations, we focus on the prediction probability of the 233 correct class, the hidden states before the classification layer (parameterized by w, i.e., the classifier 234 or language modeling head), and the gradients of this classification layer.² Given a training instance 235 (x_i, y_i) and a model θ , let $h(x_i; \theta)$ be the hidden states. We can derive the prediction probability of the 236 correct class $p_{\theta}(y_i|x_i) = \sigma(y_i w^T h(x_i; \theta))$, where σ is the sigmoid function. We can also compute 237 the gradients as $g_w(x_i, y_i) = \nabla_w \ell(f_\theta(x_i), y_i) = (1 - \sigma(y_i w^T h(x_i; \theta))) y_i h(x_i; \theta)$. Specifically, we 238 use Euclidean distance as the similarity metric, and consider hidden states as the fundamental units, 239 given their involvement in both prediction probability and gradient expressions. 240

We first consider the correct class prediction probability difference. Formally, given two training instances (x_i, y_i) and (x_j, y_j) , this is $|p_\theta(y_i|x_i) - p_\theta(y_j|x_j)| = |\sigma(y_i w^T h(x_i)) - \sigma(y_j w^T h(x_j))|$. Because the Sigmoid function is smooth and monotonically increasing, when this difference is small, $w^T |y_i h(x_i) - y_j h(x_j)|$ should also be small. We can identify two cases: when the labels agree, i.e., $y_i = y_j$, $h(x_i)$ and $h(x_j)$ should be close in both direction and length (i.e., L2 norm) *after being projected onto* w; when labels do not agree, i.e., $y_i \neq y_j$, $h(x_i)$ and $h(x_j)$ should be opposite but of similar length. In other words, in contrast to distances between hidden states, i.e., $||h(x_i) - h(x_j)||$, the difference in prediction probabilities also considers label agreement.

Let us now analyze the Euclidean distance between instances represented as gradients³,

$$\begin{aligned} \|g_w(x_i, y_i) - g_w(x_j, y_j)\|_2 &= \sqrt{\|g_w(x_i, y_i)\|_2^2 + \|g_w(x_j, y_j)\|_2^2 - 2g_w(x_i, y_i)^T g_w(x_j, y_j)} \\ &= \sqrt{\operatorname{err}_i^2 \|h_i\|_2^2 + \operatorname{err}_j^2 \|h_j\|_2^2 - 2\operatorname{err}_i \operatorname{err}_j h(x_i)^T h(x_j) y_i y_j}, \quad \text{where } \operatorname{err}_i = 1 - \sigma(y_i w^T h(x_i)). \end{aligned}$$

Here err_i is the prediction error of instance (x_i, y_i) . We discuss two scenarios. First, if the gradients are normalized (as in LESS), the first two terms above are both 1. The third term depends on the cosine similarity (equivalent to the dot product) of hidden states and label agreement. To minimize gradient distance, when $y_i = y_j$, $h(x_i)$ and $h(x_j)$ should have high cosine similarities; when $y_i \neq y_j$, $h(x_i)$ and $h(x_j)$ should be of low cosine similarities. This is similar to the prediction probability case, except that the similarity between hidden states is independent of both their lengths and w.

260 Second, when gradients are not normalized, the gradient distance is $\|\operatorname{err}_i h(x_i) - \operatorname{err}_j h(x_j)\|_2$ and 261 $\|\operatorname{err}_i h(x_i) + \operatorname{err}_j h(x_j)\|_2$ respectively when $y_i = y_j$ and $y_i \neq y_j$. Similar to the case for prediction 262 probability, when the gradient distance is small, hidden states should be similar when their labels 263 agree, and dissimilar when their labels disagree. However, here the hidden states are scaled by their 264 prediction errors: if we assume that the length of hidden states are similar (which holds in practice), 265 we require their prediction errors to be similar as well.

In summary, similarities based on training dynamics and gradients are closely related to hidden states.
 Besides, they both implicitly encode *label agreement*, and gradients also encode *prediction errors*.

268 269

248

249 250 251

253

222

²Classification layer gradient is a popular choice for large-scale models, e.g., Pruthi et al. (2020).

³Our conclusions apply to cosine similarity, as used in LESS (equates to using normalized gradients here).



Figure 1: Interactions between data representations and selection algorithms. We generate 600 data points from a 2D Gaussian mixture model and compare different methods to select 30% (180) of the data points. The color represents the ground-truth label, and the red Xs are the selected data points. (1) Using different representations with the same selection objective or algorithm leads to drastically different subset selections (a–c, g–h); (2) The sensitivity of selection algorithms towards the used representations varies (e.g., compare d–f with b–c).

296

288

289

290

291

3.3 How do representations and selection algorithms interact?

297 This section analyzes how different representations affect the outcomes of selection algorithms. 298 Specifically, we aim to answer two questions: *does a given selection algorithm choose different subset* 299 *selections when using different representations? Which selection algorithms are more sensitive to* 300 *representation changes?* For this purpose, we first conduct an interpretable analysis using a synthetic 301 dataset. Building on these insights, we will empirically study more complex scenarios in §4.

To compare selection objectives across different representations, we focus on three selection algorithms that are representation-agnostic: prototypicality (prioritizing difficulty), S2L (prioritizing diversity), and LESS (prioritizing influence on validation set performance). We also include Hardto-Learn for comparison. We generate 600 data points from a 2D Gaussian mixture model and use each method to select 30% (180) of the data points (Figures 1 and 6). We train a logistic regression classifier to serve as the reference model to obtain training dynamics and gradients, and the original data points are used as hidden states. Our observations are as follows.

First, even when data pruning methods have the same objective, the representations and selection algorithms used can result in drastically different subset selections. For example, both Hard-to-Learn and Prototypicality aim to select difficult instances. To achieve this, Hard-to-Learn selects instances with low correct class probabilities, which favors the ones that are located near the decision boundary (Figure 1a). In contrast, Prototypicality clusters data instances based on hidden states, and selects those that are far from the cluster centroids. As a result, it selects instances from the sparsely populated regions (Figure 1b).

Second, even the same selection algorithm can prefer different data points, when using different representations. As discussed in §3.2, instances with similar hidden states but different labels are far apart in gradient space (§3.2). For example, when using gradients, Prototypicality favors instances near the decision boundary as they are far from centroids (Figure 1c). In contrast, with hidden states, it selects instances far from the decision boundary, as mentioned above (Figure 1b).

Third, *selection algorithms vary in sensitivity towards representations*. For example, S2L selects similar instances with different representations (Figures 1d–f). Indeed, diversity-preserving algorithms are less affected by the choice of representation because they sample evenly from different regions. In contrast, both Prototypicality (Figure 1b–c) and LESS (Figure 1g–h) select remarkably different



346 Figure 2: Spearman correlation of scores calculated by different methods (2a-2b, all methods select 347 instances with the highest scores to retrain) and their model performance under different data budgets 348 (2c-2f). Here Hard, Proto, and Mem refer to Hard-to-Learn, Prototypicality, and Memorization. 349 LESS-OOD refers to LESS using DynaHate as the validation set. We also experimented with forced label matching. Regarding correlations, we identify two pairs of methods with moderate to high 350 correlations: Hard-to-Learn and Memorization, and Prototypicality and SemDedup. Regarding model 351 performance, we show that 1) data pruning is not always effective: hidden-state-based methods do 352 not outperform random selection, while LESS (with label matching) is competitive on all tasks; 353 2) different objectives suit different scenarios: maximizing difficulty performs well with high data 354 budgets, and prioritizing validation performance works better when there is a mismatch between the 355 train and test distributions; preserving diversity works less well. 356

358 359

360 361 362

363

364 365

366

367

368

4 EMPIRICAL EVALUATIONS

4.1 EXPERIMENTAL SETUP

Building on our theoretical and synthetic analyses (§3), we now study the impact of various data representations and selection algorithms from §3.1 on NLP tasks. We evaluate their consistency (§4.2) and performance across different data budgets (§4.3). We also perform ablation studies to examine the role of each component, as well as their effectiveness in handling label noise (§4.4).

data points using hidden states and gradients, because instance similarities computed by gradients

additionally encode their label agreement and the magnitudes of their prediction errors (§3.2).

We evaluate on three diverse NLP tasks: CAD (hate speech detection, binary classification (Vidgen et al., 2021)), WinoGrande (common sense reasoning, multiple choice (Sakaguchi et al., 2021)), and DialogSum (summarization, generation (Chen et al., 2021)). For CAD, we include Dyna-Hate (Kiela et al., 2021) as an OOD test set. For CAD and WinoGrande we use DeBERTaV3_{Large} and DeBERTaV3_{Base} (He et al., 2023), and for DialogSum we use OPT-125M and OPT-350M (Zhang et al., 2022).⁴ We experiment with six data budgets: 5%, 15%, 30%, 50%, 70%, and 100% of the original dataset, and fine-tune all models for 15 epochs. For fair comparisons, we use the same

⁴We use relatively small models to avoid huge computation during both training (trained 1200+ models for controlled comparisons) and gradient projection (can take > 10 times because of the high dimensionality).

reference models as the main models (more details in Appendix A), although recent works (Du et al., 2023; Yang et al., 2024) have shown the promising performance of using efficient reference models.
We also experimented with forced label matching to match the original dataset (LESS-balanced in Figure 2d, and Figure 2f), to address highly skewed pruned datasets (details in §4.3).

382

384

4.2 CONSISTENCY BETWEEN DATA PRUNING METHODS

This section analyzes whether different pruning methods rank data points similarly. We use Spearman's r to compare the scores assigned to each instance by various methods⁵ because (1) rank-based correlation are more suitable for comparing scores of different ranges, and (2) it allows for evaluation without setting a fixed data budget. Results for DeBERTaV3_{Large} on WinoGrande and OPT-350M on DialogSum are shown in Figures 2a and 2b, with additional results in Appendix C.2.

390 Consistent with §3.3, we observe that the same objective can lead to different data selections: For 391 instance, although both Prototypicality and Hard-to-Learn aim to select difficult instances, they show almost no correlation. However, the pairs Prototypicality - SemDedup and Hard-to-Learn - Memo-392 *rization* show moderate consistency across datasets and models. This aligns with our expectations: 393 Both Prototypicality and SemDedup are based on clustering hidden states, with SemDedup adding a 394 step for semantic deduplication. Meanwhile, both Hard-to-Learn and Memorization select difficult 395 instances whose predictions are barely improved by training on other instances (Jiang et al., 2021). 396 Since Memorization requires gradients, it is far more costly, making Hard-to-Learn the more scalable 397 option, although the moderate correlation suggests there are still differences in their data selections.

398 399 400

431

4.3 PERFORMANCE UNDER DIFFERENT DATA BUDGETS

401 Our previous analyses focus on how different components in-402 fluence data selection. This section further analyzes their impli-403 cations on model performance under different data budgets, to 404 answer questions regarding which data pruning method to use in specific scenarios. Three baselines are considered: random se-405 lection (Rand), the full original dataset (100% data budget), and 406 a dummy predictor (Dummy, the better one between a random-407 ized predictor and a majority class predictor), as performance 408 below it is considered as failed (Mosbach et al., 2021). We 409 make three observations (Figures 2c to 2f, and Appendix C.3). 410

First, *data pruning is not always effective*. For example, Hardto-Learn and S2L perform underperform random selection on
DialogSum (Figure 2c) and CAD (Figure 2d). Surprisingly,
hidden-state-based methods (Prototypicality and SemDedup),
perform worse or similar to random selection on all tasks, sug-



Figure 3: DeBERTaV3_{Large} on CAD: the ratio of hateful instances in selected instances, under different data budgets.

gesting clustering pretrained hidden states is not suitable for our setting.⁶

Second, *higher data budgets are needed for methods that select difficult instances*, i.e., Hard-to-Learn and Memorization. They achieve good performance with > 30% data budgets on CAD (Figure 2d) and WinoGrande (Figure 8e), but they both struggle with lower data budgets (< 30%) (e.g., worse or comparable to Dummy). This is consistent with Swayamdipta et al. (2020): including only the most difficult instances will make models fail to converge. Moreover, due to their consistency in both scores (§4.2) and classification performance, future studies may prioritize Hard-to-Learn over Memorization for selecting difficult instances in classification tasks, aiming for better efficiency.

Third, gradient-based methods (LESS and Memorization) are competitive (Figure 2c-2e), indicating
that gradients are effective data representations. Between them, LESS performs better but requires *label matching for highly skewed datasets*. For example, on CAD, where 90% of instances are
non-hateful, directly applying LESS yields poor results (Figure 2d).

 ⁵All methods, except S2L, assign scores to instances and retain the highest-scoring ones. Since S2L only assigns cluster labels, we exclude it from this analysis. We also experimented with overlaps under different data budgets and obtain consistent observations.

⁶Note that our experiments differ from previous studies that use hidden states, as they focused on high data budget settings (Sorscher et al., 2022) and noisy pretraining datasets (Abbas et al., 2023).



(a) Pretrained vs. Fine-Tuned (FT) Hidden States: OPT-350M on DialogSum (Rouge-L)



Hate (F1)

Figure 4: Ablation studies on pretrained vs. fine-tuned hidden states (4a), and using different data representations with the same selection algorithm (4b-4c).

We hypothesize the reason to be as follows. LESS ranks training instances based on their gradient 448 similarity to validation data. However, gradient similarities depend on labels: since hidden states 449 usually have > 0 cosine similarities (known as anisotropic space), instances with the same label have 450 higher gradient similarities (§3.2). This makes LESS over-select instances with the majority label 451 when using unbalanced validation sets, e.g., non-hateful instances in CAD, leading to pruned datasets 452 with even more skewed label distributions. 453

We validate this by plotting the ratio of hateful labels in selected instances across different methods 454 and data budgets in Figure 3. We confirm that non-hateful instances have much higher LESS scores 455 than hateful ones, with very few hateful instances being selected until the data budget reaches 90%, 456 unlike other methods.¹ To address this, we enforced the selected label ratio to match the original 457 dataset, which substantially improves LESS performance (Figure 2f). However, this constraint can 458 harm methods like Hard-to-Learn and memorization, which favor hateful instances. 459

460 4.4 **ABLATION STUDIES** 461

462 This section presents ablation studies to explore the impact of key components, including fine-tuned 463 versus pretrained hidden states, the interaction between representations and selection algorithms (as 464 in §3.3 but focusing on task performance), and the ability of these methods to handle label noises. We include our results in Figure 4 and Appendix C.4. 465

467 **Fine-Tuned Hidden States** Our previous results show that hidden-state-based methods perform comparably to random selection. However, fine-tuning could potentially encode task-specific and 468 label information into hidden states, helping identify prototypical and diverse training instances. 469

470 We therefore study whether using fine-tuned hidden states can improve model performance. Specif-471 ically, we experiment with two types of fine-tuned hidden states: early (fine-tuned for one epoch, 472 retaining more pretrained knowledge) and late (fine-tuned for 15 epochs, encoding more task-specific and label information). The results of both prototypicality and SemDedup are shown in Figure 4a. 473 *Fine-tuned hidden states can barely improve model performance:* there is little difference between 474 using different hidden states, and they all perform comparably to random selection. 475

476

466

441

442

443 444

445 446 447

Representation-Selection Relation It remains unclear whether the effectiveness of different data 477 pruning methods stems primarily from the representations they use, the selection algorithms they use, 478 or a specific combination of both. Similar to §3.3, to disentangle the contributions from these different 479 components, we conduct a systematic study that combines all three different data representations with 480 the selection algorithms of S2L, Prototypicality, and LESS, which respectively prioritize difficulty, 481 diversity, and validation performance. We show the results for DeBERTaV3_{Large} on WinoGrande and 482 DynaHate in Figure 4b and Figure 4c, and make two observations.

⁴⁸⁴ ⁷Sorscher et al. (2022) reported a similar imbalance for Prototypicality. We also observe marginally more 485 skewed pruned datasets, and forced matching slightly improve their performance (especially with high data budgets, which is their original setting). However, their performance is overall comparable to random selection.

First, *similar representations often lead to similar performance*. For example, Figure 4b shows that
 the performance of using the same representation with different selection algorithms are more similar
 than using different representations with the same selection algorithm. This indicates that, *compared to selection algorithm, representation plays a more fundamental role.*

Second, *the optimal selection algorithm is task-dependent* and still plays a key role. For example,
 Figure 4c shows that methods using the LESS selection algorithm outperform others on DynaHate.
 One possible explanation is that DynaHate is an OOD test set that presents a substantial distribution
 shift from the training data (CAD), while the LESS selection algorithm compensates for this shift by
 selecting instances matching the validation set. This observation suggests that future practitioners
 should choose algorithms based on goals: prioritizing difficulty for better robustness, and prioritizing
 contribution to validation performance when adapting models to a different domain.

497

498 **Noise Detection** An important applications of data pruning is 499 removing noise from training data (Swayamdipta et al., 2020; 500 Paul et al., 2021). To evaluate the noise detection capabilities of different methods, we create two noisy versions of each dataset 501 by changing 2% and 5% of the training labels. Concretely, we 502 flip the labels for CAD and WinoGrande (both datasets have 503 binary labels), and substitute the summaries of DialogSum with 504 a randomly selected one. We then train reference models on 505 these noisy datasets, and use different methods to select data. 506 We use fine-tuned hidden states (both early and late) instead 507 of pretrained hidden states, as pretrained hidden states do not 508 encode label-relevant information. 509



We analyze the percentage of selected noisy instances out of all noisy instances across data budgets. Figure 5 shows results for DeBERTaV3_{Large} on WinoGrande (2% noise rate), with

Figure 5: The percentage of selected noisy instances in all noisy instances of DeBERTaV3_{Large} on WinoGrande, with 2% noise rate.

additional results in Appendix C.4.⁸ Except for Hard-to-Learn, which is proposed for noise detection (Swayamdipta et al., 2020), none of the methods effectively filter out noisy instances. However, Hard-to-Learn tends to select more noisy instances, which conflicts with our goal of identifying clean data. Consequently, *none of the data pruning methods are suitable for training with noisy data*. Notably, although Memorization shows moderate consistency with Hard-to-Learn (§4.2), it only ranks noisy instances marginally higher than clean instances.

518 519 520

521

522

523

524

525

526

527

528

5 CONCLUSION

Summary Despite the success of data pruning, the roles of its two key components – data representation and selection algorithm – and their interactions, are not well understood. In this work, we have conducted both theoretical and empirical analyses on these two choices in fine-tuning tasks. Our results highlight the importance of using rich data representations, showing that gradient-based methods consistently outperform computationally cheaper alternatives. Additionally, the optimal selection algorithm varies depending on specific use cases, although the outcomes are heavily influenced by the chosen representations. Moving forward, our findings stress the need for the development of **scalable supervised representations**, i.e., representations that encode label-relevant information, as more effective alternatives to the current unsupervised ones, e.g., pretrained hidden states.

529 530 531

Limitations The most notable limitation of our work is its focus on task-specific fine-tuning that leaves multi-task instruction tuning unexplored. This is largely due to (1) the huge amount of computation required to conduct rigorous controlled studies as ours, and (2) the challenges in scalable and low-cost evaluation (Zheng et al., 2023). Future studies could explore instruction-tuning with synthetic data, which has been recently shown to be effective for proof-of-concept studies (Allen-Zhu & Li, 2024). Moreover, we focus on methods that do not require external models (e.g., prompting language models to evaluate example quality). Future work can extend our analysis to them.

⁸S2L samples instances from different clusters in a balanced way (§3.1). Therefore, even with a data budget of 100%, it does not select all data points, making the percentage does not reach 100%.

540 REFERENCES

- Amro Kamal Mohamed Abbas, Kushal Tirumala, Daniel Simig, Surya Ganguli, and Ari S. Morcos.
 Semdedup: Data-efficient learning at web-scale through semantic deduplication. In *ICLR 2023 Workshop on Mathematical and Empirical Understanding of Foundation Models*, 2023. URL https://openreview.net/forum?id=4vlGm9gv6c.
- Zeyuan Allen-Zhu and Yuanzhi Li. Physics of language models: Part 3.1, knowledge storage and
 extraction. In *Forty-first International Conference on Machine Learning*, 2024. URL https:
 //openreview.net/forum?id=5x788rqbcj.
- Robert John Nicholas Baldock, Hartmut Maennel, and Behnam Neyshabur. Deep learning through the lens of example difficulty. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*, 2021. URL https://openreview. net/forum?id=WWRBHhH158K.
- Irina Bejan, Artem Sokolov, and Katja Filippova. Make every example count: On the stability and utility of self-influence for learning from noisy NLP datasets. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 10107–10121, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.625.
 //aclanthology.org/2023.emnlp-main.625.
- Rajat Bhatnagar, Ananya Ganesh, and Katharina Kann. CHIA: CHoosing instances to annotate for machine translation. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2022*, pp. 7299–7315, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.findings-emnlp.540. URL https://aclanthology.org/2022.findings-emnlp.540.
- Lichang Chen, Shiyang Li, Jun Yan, Hai Wang, Kalpa Gunaratna, Vikas Yadav, Zheng Tang, Vijay Srinivasan, Tianyi Zhou, Heng Huang, and Hongxia Jin. Alpagasus: Training a better alpaca with fewer data. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=FdVXgSJhvz.
- Yulong Chen, Yang Liu, Liang Chen, and Yue Zhang. DialogSum: A real-life scenario dialogue summarization dataset. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pp. 5062–5074, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.findings-acl.449. URL https://aclanthology.org/2021.findings-acl.449.
- Yupei Du, Albert Gatt, and Dong Nguyen. Ftft: efficient and robust fine-tuning by transferring training dynamics. *arXiv preprint arXiv:2310.06588*, 2023. URL https://arxiv.org/abs/2310.06588.
- Logan Engstrom, Axel Feldmann, and Aleksander Madry. Dsdm: Model-aware dataset selection
 with datamodels. In *Forty-first International Conference on Machine Learning*, 2024. URL
 https://openreview.net/forum?id=GC8HkKeH8s.
- Vitaly Feldman. Does learning require memorization? a short tale about a long tail. In *Proceedings* of the 52nd Annual ACM SIGACT Symposium on Theory of Computing, STOC 2020, pp. 954–959, New York, NY, USA, 2020. Association for Computing Machinery. ISBN 9781450369794. doi: 10.1145/3357713.3384290. URL https://doi.org/10.1145/3357713.3384290.
- Vitaly Feldman and Chiyuan Zhang. What neural networks memorize and why: Discovering the long tail via influence estimation. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), Advances in Neural Information Processing Systems, volume 33, pp. 2881–2891. Curran Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper_files/paper/2020/file/1e14bfe2714193e7af5abc64ecbd6b46-Paper.pdf.
- Pengcheng He, Jianfeng Gao, and Weizhu Chen. DeBERTav3: Improving deBERTa using ELECTRA style pre-training with gradient-disentangled embedding sharing. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?
 id=sE7-XhLxHA.

616

- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre. Training compute-optimal large language models, 2022.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022. URL https://openreview.net/forum? id=nZeVKeeFYf9.
- Andrew Ilyas, Sung Min Park, Logan Engstrom, Guillaume Leclerc, and Aleksander Madry.
 Datamodels: Understanding predictions with data and data with predictions. In Kamalika
 Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato (eds.),
 Proceedings of the 39th International Conference on Machine Learning, volume 162 of *Proceedings of Machine Learning Research*, pp. 9525–9587. PMLR, 17–23 Jul 2022. URL
 https://proceedings.mlr.press/v162/ilyas22a.html.
- Osman Batur İnce, Tanin Zeraati, Semih Yagcioglu, Yadollah Yaghoobzadeh, Erkut Erdem, and Aykut Erdem. Harnessing dataset cartography for improved compositional generalization in transformers. In *The 2023 Conference on Empirical Methods in Natural Language Processing*, 2023. URL https://openreview.net/forum?id=BMIjPXooNq.
 - Hamish Ivison, Yizhong Wang, Valentina Pyatkin, Nathan Lambert, Matthew Peters, Pradeep Dasigi, Joel Jang, David Wadden, Noah A. Smith, Iz Beltagy, and Hannaneh Hajishirzi. Camels in a changing climate: Enhancing Im adaptation with tulu 2, 2023.
- Angela H. Jiang, Daniel L.-K. Wong, Giulio Zhou, David G. Andersen, Jeff Dean, Gregory R. Ganger,
 Gauri Joshi, Michael Kaminsky, Michael A. Kozuch, Zachary Chase Lipton, and Padmanabhan
 Pillai. Accelerating deep learning by focusing on the biggest losers. *ArXiv*, abs/1910.00762, 2019.
 URL https://api.semanticscholar.org/CorpusID:203626838.
- Ziheng Jiang, Chiyuan Zhang, Kunal Talwar, and Michael C Mozer. Characterizing structural regularities of labeled data in overparameterized models. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 5034–5044. PMLR, 18–24 Jul 2021. URL https://proceedings.mlr.press/v139/jiang21k.html.
- William B. Johnson and Joram Lindenstrauss. Extensions of Lipschitz mappings into a Hilbert space.
 pp. 189–206. 1984. doi: 10.1090/conm/026/737400. URL http://www.ams.org/conm/026/.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- 635 Douwe Kiela, Max Bartolo, Yixin Nie, Divyansh Kaushik, Atticus Geiger, Zhengxuan Wu, Bertie Vidgen, Grusha Prasad, Amanpreet Singh, Pratik Ringshia, Zhiyi Ma, Tristan Thrush, Sebastian 636 Riedel, Zeerak Waseem, Pontus Stenetorp, Robin Jia, Mohit Bansal, Christopher Potts, and Adina 637 Williams. Dynabench: Rethinking benchmarking in NLP. In Kristina Toutanova, Anna Rumshisky, 638 Luke Zettlemoyer, Dilek Hakkani-Tur, Iz Beltagy, Steven Bethard, Ryan Cotterell, Tanmoy 639 Chakraborty, and Yichao Zhou (eds.), Proceedings of the 2021 Conference of the North American 640 Chapter of the Association for Computational Linguistics: Human Language Technologies, pp. 641 4110–4124, Online, June 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021. 642 naacl-main.324. URL https://aclanthology.org/2021.naacl-main.324. 643
- Krishnateja Killamsetty, Durga S, Ganesh Ramakrishnan, Abir De, and Rishabh Iyer. Grad-match:
 Gradient matching based data subset selection for efficient deep model training. In Marina Meila
 and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning*,
 volume 139 of *Proceedings of Machine Learning Research*, pp. 5464–5474. PMLR, 18–24 Jul
 2021. URL https://proceedings.mlr.press/v139/killamsetty21a.html.

- 648 Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. 649 In Doina Precup and Yee Whye Teh (eds.), Proceedings of the 34th International Conference 650 on Machine Learning, volume 70 of Proceedings of Machine Learning Research, pp. 1885-651 1894. PMLR, 06-11 Aug 2017. URL https://proceedings.mlr.press/v70/koh17a. 652 html.
- Devin Kwok, Nikhil Anand, Jonathan Frankle, Gintare Karolina Dziugaite, and David Rolnick. 654 Dataset difficulty and the role of inductive bias, 2024. 655
- 656 Ming Li, Yong Zhang, Zhitao Li, Jiuhai Chen, Lichang Chen, Ning Cheng, Jianzong Wang, 657 Tianyi Zhou, and Jing Xiao. From quantity to quality: Boosting LLM performance with self-658 guided data selection for instruction tuning. In Kevin Duh, Helena Gomez, and Steven Bethard 659 (eds.), Proceedings of the 2024 Conference of the North American Chapter of the Association 660 for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pp. 7595–7628, Mexico City, Mexico, June 2024. Association for Computational Linguistics. URL 661 https://aclanthology.org/2024.naacl-long.421. 662
- 663 Wei Liu, Weihao Zeng, Keqing He, Yong Jiang, and Junxian He. What makes good data for 664 alignment? a comprehensive study of automatic data selection in instruction tuning. In The Twelfth 665 International Conference on Learning Representations, 2024. URL https://openreview. 666 net/forum?id=BTKAeLqLMw. 667
- Keming Lu, Hongyi Yuan, Zheng Yuan, Runji Lin, Junyang Lin, Chuanqi Tan, Chang Zhou, and 668 Jingren Zhou. #instag: Instruction tagging for analyzing supervised fine-tuning of large language 669 models. In The Twelfth International Conference on Learning Representations, 2024. URL 670 https://openreview.net/forum?id=pszewhybU9. 671
- 672 Pratyush Maini, Saurabh Garg, Zachary Chase Lipton, and J Zico Kolter. Characterizing datapoints 673 via second-split forgetting. In ICML 2022: Workshop on Spurious Correlations, Invariance and 674 Stability, 2022. URL https://openreview.net/forum?id=bCdztvpaEUG. 675
- Max Marion, Ahmet Üstün, Luiza Pozzobon, Alex Wang, Marzieh Fadaee, and Sara Hooker. When 676 less is more: Investigating data pruning for pretraining llms at scale, 2023. 677
- 678 Sören Mindermann, Jan M Brauner, Muhammed T Razzak, Mrinank Sharma, Andreas Kirsch, 679 Winnie Xu, Benedikt Höltgen, Aidan N Gomez, Adrien Morisot, Sebastian Farquhar, and Yarin 680 Gal. Prioritized training on points that are learnable, worth learning, and not yet learnt. In 681 Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato 682 (eds.), Proceedings of the 39th International Conference on Machine Learning, volume 162 of 683 Proceedings of Machine Learning Research, pp. 15630–15649. PMLR, 17–23 Jul 2022. URL https://proceedings.mlr.press/v162/mindermann22a.html. 684
- Marius Mosbach, Maksym Andriushchenko, and Dietrich Klakow. On the stability of fine-686 tuning {bert}: Misconceptions, explanations, and strong baselines. In International Conference on Learning Representations, 2021. URL https://openreview.net/forum?id= 688 nzpLWnVAyah. 689
- 690 Sung Min Park, Kristian Georgiev, Andrew Ilyas, Guillaume Leclerc, and Aleksander Madry. TRAK: 691 Attributing model behavior at scale. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), Proceedings of the 40th International 692 Conference on Machine Learning, volume 202 of Proceedings of Machine Learning Research, pp. 693 27074-27113. PMLR, 23-29 Jul 2023. URL https://proceedings.mlr.press/v202/ 694 park23c.html.
- 696 Mansheej Paul, Surya Ganguli, and Gintare Karolina Dziugaite. Deep learning on a data diet: 697 Finding important examples early in training. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan (eds.), Advances in Neural Information Processing Systems, 2021. URL 699 https://openreview.net/forum?id=Uj7pF-D-YvT.
- 700

687

653

Garima Pruthi, Frederick Liu, Satyen Kale, and Mukund Sundararajan. Estimating training data influ-701 ence by tracing gradient descent. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin

702 (eds.), Advances in Neural Information Processing Systems, volume 33, pp. 19920–19930. Cur-703 ran Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper_files/ 704 paper/2020/file/e6385d39ec9394f2f3a354d9d2b88eec-Paper.pdf. 705 Noveen Sachdeva, Benjamin Coleman, Wang-Cheng Kang, Jianmo Ni, Lichan Hong, Ed H. Chi, 706 James Caverlee, Julian McAuley, and Derek Zhiyuan Cheng. How to train data-efficient llms, 707 2024. 708 709 Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: an ad-710 versarial winograd schema challenge at scale. Commun. ACM, 64(9):99-106, aug 2021. ISSN 711 0001-0782. doi: 10.1145/3474381. URL https://doi.org/10.1145/3474381. 712 713 Nikhil Sardana, Jacob Portes, Sasha Doubov, and Jonathan Frankle. Beyond chinchilla-optimal: 714 Accounting for inference in language model scaling laws. In Forty-first International Conference on Machine Learning, 2024. URL https://openreview.net/forum?id=0bmXrtTDUu. 715 716 Ben Sorscher, Robert Geirhos, Shashank Shekhar, Surya Ganguli, and Ari S. Morcos. Beyond neural 717 scaling laws: beating power law scaling via data pruning. In Alice H. Oh, Alekh Agarwal, Danielle 718 Belgrave, and Kyunghyun Cho (eds.), Advances in Neural Information Processing Systems, 2022. 719 URL https://openreview.net/forum?id=UmvSlP-PyV. 720 721 Swabha Swayamdipta, Roy Schwartz, Nicholas Lourie, Yizhong Wang, Hannaneh Hajishirzi, Noah A. 722 Smith, and Yejin Choi. Dataset cartography: Mapping and diagnosing datasets with training 723 dynamics. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu (eds.), Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pp. 9275-724 9293, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020. 725 emnlp-main.746. URL https://aclanthology.org/2020.emnlp-main.746. 726 727 Megh Thakkar, Tolga Bolukbasi, Sriram Ganapathy, Shikhar Vashishth, Sarath Chandar, and Partha 728 Talukdar. Self-influence guided data reweighting for language model pre-training. In Houda 729 Bouamor, Juan Pino, and Kalika Bali (eds.), Proceedings of the 2023 Conference on Empir-730 ical Methods in Natural Language Processing, pp. 2033–2045, Singapore, December 2023. 731 Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.125. URL 732 https://aclanthology.org/2023.emnlp-main.125. 733 Kushal Tirumala, Daniel Simig, Armen Aghajanyan, and Ari Morcos. D4: Improving 734 llm pretraining via document de-duplication and diversification. In A. Oh, T. Neumann, 735 A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), Advances in Neural Informa-736 tion Processing Systems, volume 36, pp. 53983-53995. Curran Associates, Inc., 2023. 737 https://proceedings.neurips.cc/paper_files/paper/2023/file/ URL 738 a8f8cbd7f7a5fb2c837e578c75e5b615-Paper-Datasets_and_Benchmarks. 739 pdf. 740 741 Mariya Toneva, Alessandro Sordoni, Remi Tachet des Combes, Adam Trischler, Yoshua Bengio, and 742 Geoffrey J. Gordon. An empirical study of example forgetting during deep neural network learning. 743 In International Conference on Learning Representations, 2019. URL https://openreview. net/forum?id=BJlxm30cKm. 744 745 Bertie Vidgen, Dong Nguyen, Helen Margetts, Patricia Rossini, and Rebekah Tromble. Introducing 746 CAD: the contextual abuse dataset. In Kristina Toutanova, Anna Rumshisky, Luke Zettlemoyer, 747 Dilek Hakkani-Tur, Iz Beltagy, Steven Bethard, Ryan Cotterell, Tanmoy Chakraborty, and Yichao 748 Zhou (eds.), Proceedings of the 2021 Conference of the North American Chapter of the Association 749 for Computational Linguistics: Human Language Technologies, pp. 2289–2303, Online, June 750 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.naacl-main.182. URL 751 https://aclanthology.org/2021.naacl-main.182. 752 753 Mengzhou Xia, Sadhika Malladi, Suchin Gururangan, Sanjeev Arora, and Danqi Chen. LESS: Selecting influential data for targeted instruction tuning. In ICLR 2024 Workshop on Navigating 754

 and Addressing Data Problems for Foundation Models, 2024. URL https://openreview. net/forum?id=Kw3ckB2Kfc.

756	Shuo Yang, Zeke Xie, Hanyu Peng, Min Xu, Mingming Sun, and Ping Li. Dataset pruning: Reduc-
757	ing training data by examining generalization influence. In The Eleventh International Confer-
758	ence on Learning Representations, 2023. URL https://openreview.net/forum?id=
759	4wZiAXD29TQ.
760	$\mathbf{X} = \mathbf{X} = \mathbf{X} + \mathbf{M} + $
761	Yu Yang, Siddhartha Mishra, Jeffrey N Chiang, and Banaran Mirzasoleiman. Smalltolarge (s21):
762	of small models. ArYiv abs/2403 07384, 2024 LIPL https://api.somanti.co.cholar
763	org/CorpusID:268363364
764	019/001pu31D.200505504.
765	Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher
766	Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer language
767	models. arXiv preprint arXiv:2205.01068, 2022.
768	Lianmin Zheng Wei-Lin Chiang Ying Sheng Siyuan Zhuang Zhanghao Wu Yonghao Zhuang
769	Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang, Joseph E, Gonzalez, and Jon Stoica.
770	Judging LLM-as-a-judge with MT-bench and chatbot arena. In <i>Thirty-seventh Conference on</i>
771	Neural Information Processing Systems Datasets and Benchmarks Track, 2023. URL https:
772	//openreview.net/forum?id=uccHPGDlao.
773	
774	
775	
776	
777	
778	
779	
780	
781	
782	
783	
784	
785	
786	
787	
788	
789	
790	
791	
792	
793	
794	
795	
796	
797	
798	
799	
800	
801	
802	
803	
804	
805	
806	
807	
808	
809	

810 A EXPERIMENTAL DETAILS

812 **Implementation Details** All experiments were conducted using the AdamW optimizer. For most 813 models, we set the learning rate to 2e-5, except for DeBERTaV3_{Large}, where we followed He et al. 814 (2023) and used 1e-5. Additionally, we warmup the learning rate for the first 10% of training steps. 815 For gradient-based pruning, the reference models were trained with LoRA, using a higher learning 816 rate of 1e-4, r = 64, and $\alpha = 16$ following Ivison et al. (2023), and apply LoRA on all linear layers. We train all models for 15 epochs, For batch size, we used 16 for both WinoGrande and DialogSum, 817 and 32 for CAD, to fit all experiments on a single NVIDIA A100-40GB GPU. We use maximum 818 sequence lengths of 300, 128, and 512 tokens. For all experiments we use the same reference models 819 as the main models for fair comparison. 820

For *k*-Means clustering in S2L, Prototypicality, and SemDedup, we use 100 clusters on CAD and
DialogSum, and 200 clusters on WinoGrande, following the suggestions from Tirumala et al. (2023)
to set the number of clusters to around the square root of the number of instances. Moreover,
we compute gradients using the first five checkpoints for all experiments, and project them into a
1024-dimensional space using Park et al. (2023) (details see hyperparameter search).

Evaluation Metrics We evaluated CAD and DynaHate using the macro F1 score, WinoGrande by accuracy, and DialogSum by ROUGE-1, ROUGE-2, and ROUGE-L (from HuggingFace Evaluate), following the original studies (Vidgen et al., 2021; Sakaguchi et al., 2021; Chen et al., 2021).

Infrastructure All experiments were run on a single NVIDIA A100-40GB GPU using three random seeds. We used PyTorch 2.3, Transformers 4.42, and vLLM 0.5 for training and inference. Moreover, we use bfloat16 on all experiments to improve efficiency.

Hyperparameter Search We mainly searched for four hyperparameters: the number of training
 epochs, the number of clusters for *k*-Means clustering, the dimensionality of the projected gradients, and the checkpoints to use for gradient computation.

For the number of training epochs, we first perform a search of 3, 5, 7, and 10 epochs on all datasets and models, using three random seeds. We observe that models of different sizes share similar performance trends over epochs, with improvements continuing as the number of epochs increased. We therefore use the smaller models, i.e., DeBERTaV3_{Base} and OPT-125M, and extend this search over 15, 20, and 25 epochs. Across all datasets, the best performance is achieved with 15 epochs.

For the number of clusters, we search over 2, 5, 10, 20, 50, 100, and 200 clusters for each dataset and model, using three random seeds. The results are highly consistent across cluster numbers. Following Tirumala et al. (2023), we use the square root of the dataset size as a guideline, settling on 100 clusters for CAD and DialogSum, and 200 for WinoGrande.

For gradients, we use smaller models (DeBERTaV3_{Base} and OPT-125M) for hyper-parameter search, 847 and only one random seed (0) to avoid the high costs of computing and projecting gradients. We 848 compute the gradients for all 15 checkpoints, and project them into 1024, 2048, and 4096 dimensions. 849 First, we observe that different dimensionality compute similar results, and thus choose 1024 for 850 further experiments for efficiency. Second, we experimented with different strategies for selecting 851 checkpoints, including the first three, the last three, the first five, the last five, and evenly spaced three 852 and five checkpoints. Using the first checkpoints is the most consistent with using all checkpoints, 853 with the first five yielding a minimum Spearman's rank correlation of 0.96. We therefore use the first 854 five checkpoints for all experiments.

- 855 856
- 050
- 858
- 859
- 860
- 861
- 862
- 863

B OVERVIEW OF DATA PRUNING METHOD

	Representations			
Selection	Training dynamics	Hidden states	Gradients	
Max. diversity	SmallToLarge (Yang et al., 2024)	SemDedup (Abbas et al., 2023)		
Max. difficulty	Hard-to-learn (Swayamdipta et al., 2020; Jiang et al., 2021; İnce et al., 2023)	Prototypicality (Sorscher et al., 2022) SemDedup (Abbas et al., 2023)	Memorization (Feldman & Zhang, 2020)	
/al. contribution			LESS (Xia et al., 2024)	

Table 1: Pruning methods from §3.1, categorized by their representations and selection objectives.

918 С ADDITIONAL RESULTS 919

C.1 TOY EXAMPLE



Figure 6: Interactions between data representations and selection algorithms. We generate 600 data points from a 2D Gaussian mixture model and compare different methods to select 30% (180) of the data points. The color represents the ground-truth label, and the red Xs are the selected data points.

-0.75

-1.00

1.00

0.75

- 0.50

- 0.25

- 0.00

-0.25

-0.50

-0.75

-1.00

LESS -

Merry

CONSISTENCY BETWEEN DATA PRUNING METHODS C.2







(c) DeBERTaV3_{Base} on WinoGrande

0.04

0.04

(a) DeBERTaV3_{Base} on CAD





18

C.3 PERFORMANCE UNDER DIFFERENT DATA BUDGETS

934

935 936 937

938 939

940

920

921 922

923

945

946

947

948

949

950 951

952 953

954

955

956

957

958

959

960

961

962

963 964

965 966

967

968 969 970

971

Hard

Proto 0.08

Mem

LESS 0.45

SemDedup

0.05

Haro

0.05

0.05





(g) Pretrained vs. Fine-Tuned (FT) Hidden States: OPT-350M on Di- Hidden States: OPT-125M on Di- Hidden States: OPT-350M on DialogSum (Rouge-1)

(h) Pretrained vs. Fine-Tuned (FT) (i) Pretrained vs. Fine-Tuned (FT) alogSum (Rouge-2)

alogSum (Rouge-2)

Figure 9: Ablation studies on pretrained vs. fine-tuned hidden states.





Figure 11: The percentage of selected noisy instances in all noisy instances for different models, datasets, and noise rates.