DIVERSITY MEASUREMENT AND SUBSET SELECTION FOR INSTRUCTION TUNING DATASETS

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

023 024 Paper under double-blind review

Abstract

We aim to select data subsets for the fine-tuning of large language models to more effectively follow instructions. Prior work has emphasized the importance of diversity in dataset curation but relied on heuristics such as the number of tasks. In this paper, we use determinantal point processes to capture the diversity and quality of instruction tuning datasets for subset selection. We propose to measure dataset diversity with log determinant distance that is the distance between the dataset of interest and a maximally diverse reference dataset. Our experiments demonstrate that the proposed diversity measure in the normalized weight gradient space is correlated with instruction-following performance. Consequently, it can be used to inform when data selection is the most helpful and to analyze dataset curation strategies. We demonstrate the utility of our data selection approach on various instruction tuning datasets.

1 INTRODUCTION

Large language models (LLMs) are powerful but unwieldy for practical use. They often require demonstrations in context to elicit proper responses and even then may generate responses not intended by users. The base language model is typically "instruction tuned", i.e., finetuned to predict target responses given instructions. Instruction tuning enables the base language model to perform zero-shot tasks and follow users' intent more effectively, thus improving usability. Moreover, it is an indispensable step before additional preference learning to align the language model's output to human preference (Ouyang et al., 2022).

The number of instruction tuning datasets is rapidly growing, some with millions of data points (Ding et al., 2023; Zheng et al., 2024). This growth is facilitated by the ease of generating synthetic datasets by prompting LLMs (Wang et al., 2023b) and a growing effort to retain records of realworld user interactions with these models (Zhao et al., 2024; Zheng et al., 2024). Finetuning on ever-increasing data demands additional computational resources. As training on low quality data (e.g., incorrect responses) can lead to suboptimal models. Some data selection or pruning is required.

Practitioners in the field face an important challenge of selecting the optimal data subset for finetun-040 ing to maximize instruction following performance subject to a fixed computational budget. While 041 various solutions have been proposed for finding representative subsets in active learning (Sener 042 & Savarese, 2018), their applicability to natural language datasets remains underexplored. For in-043 stance, active learning methods that search for subsets with diverse weight gradients (Ash et al., 044 2019) were ineffective in our initial studies as they prioritized data points with short responses or those with large weight gradient norms. Most related methods aim to provide sufficient coverage of instruction tuning examples in the space of decoder-based language models' output token embed-046 dings (Bukharin & Zhao, 2023; Liu et al., 2024) that lacks semantic structure (Le & Mikolov, 2014). 047 Moreover, ensuring diversity in the embedding space of encoder-based masked language models is 048 limited by encoders' short context length. 049

Practitioners also grapple with a closely related question of estimating how much data allocated
 for model finetuning would achieve comparable performance with that of the entire dataset. One
 approach involves assigning a score to each dataset that indicates the extent to which a dataset
 can be reduced without compromising performance after model finetuning. While various scoring
 methods exist, here we focus on dataset diversity. Common measures of dataset diversity often rely

Distance (LDD) ے ق AlpacaEval (% Win) 0.50.4 AlpacaEval Win Rate (ρ_p =-0.88; ρ_s =-0.85) 0.3 0.2 LogDet I 0.0 .0 GP14-Alpaca · OASST? UltraChat Wittandhy ShareCPT Solf-Instruct FLAN Dolly Alpaca

062 063 064

065

066

067

068

054

056

058 059

060 061

Figure 1: The log determinant distance in Equation (4) of instruction tuning datasets is correlated with instruction following performance when the model is finetuned on these datasets, with a Pearson correlation of -0.88 and a Spearman's rank correlation of -0.85.

on intuitive heuristics, e.g., the number of tasks (Wei et al., 2022; Sanh et al., 2022), topics and user
 intents (Lu et al., 2024), or do not scale well with the dataset size (Friedman & Dieng, 2023).

We turn to determinantal point processes (DPPs) (Kulesza & Taskar, 2012) to identify diverse subsets of high quality instruction tuning data. We investigate several choices of data representations that capture data points' similarity and find that the radial basis kernel applied to the *normalized* weight gradients of the model is particularly effective when selecting from less diverse datasets.

In addition, we measure dataset diversity with *log determinant distance* that is the difference between the log determinant of kernel matrix of a maximally diverse dataset and that of the dataset under consideration, normalized by the dataset size. Log determinant distance is readily computable from the MAP inference algorithm that identifies the optimal subset. We demonstrate that log determinant distance is correlated with instruction following performance when using weight gradients as the data representation. As a result, the diversity measure can be used to evaluate the utility of instruction tuning datasets for finetuning and to predict, before any finetuning takes place, the extent to which we can prune data without sacrificing model performance. In addition, we investigate the implications of curation strategies on diversity.

084 085

2 RELATED WORK

087 088

089

2.1 INSTRUCTION TUNING DATASETS

Diversity and quality are recurring themes in the curation of instruction tuning datasets. Early in struction tuning datasets, e.g., Super-NaturalInstructions (Wang et al., 2022) and FLAN (Wei et al., 2022; Chung et al., 2022), are adapted from existing natural language processing benchmarks, with a particular focus on scaling the number of tasks to encourage task generalization.

094 Some instruction tuning datasets are curated using Self-Instruct and its variants (Honovich et al., 095 2023; Wang et al., 2023b) that prompt a LLM to generate a wide array of instructions and high-096 quality responses. These datasets, e.g., Alpaca (Taori et al., 2023), are typically distilled from 097 performant language models that underwent finetuning to generate user-preferred responses, e.g., 098 variants of InstructGPT (Ouyang et al., 2022), and are well-suited for the purposes of creating a chat 099 assistant. Moreover, they are distilled from increasingly powerful LLMs, e.g., GPT4-Alpaca (Peng et al., 2023), and contain more complex instructions, e.g., WizardLM (Xu et al., 2024), step-by-step 100 explanations in the responses, e.g., Orca (Mukherjee et al., 2023), or multi-turn conversations, e.g., 101 UltraChat (Ding et al., 2023). 102

Another family of instruction tuning datasets aims to better reflect LLMs' real-world use cases
that include significant human authorship. Some are manually curated from sources with helpful
responses such as Reddit, e.g., LIMA (Zhou et al., 2023), or from company employees, e.g., Dolly
(Conover et al., 2023). Alternatively, real-world user interactions are curated with state-of-the-art
LLMs from the internet, e.g., ShareGPT, RealChat-1M (Zheng et al., 2024), WildChat (Zhao et al., 2024). These datasets cover a wide range user intents, capturing real-world use scenarios.

We systematically study the relative diversity of aforementioned datasets and its impact on instruction following performance. Our experiments yield insights into the efficacy of the different curation approaches, e.g., distillation and human annotation.

111 112

112 2.2 DATA SELECTION

114 Analogous to data selection for finetuning, active learning selects informative examples to label from 115 a pool of unlabeled examples subject to a fixed labeling budget. Our approach is closely related to 116 research that formulates active learning as core-set selection, i.e., finding the representative data 117 subset (Tsang et al., 2005; Welling, 2009). Examples include searching for a covering of the full 118 dataset with the smallest cover radius by solving the k-center problem (Sener & Savarese, 2018) 119 and identifying subsets that are sufficiently spread out using k-means++ initialization (Ash et al., 2019). Related, data pruning methods remove redundant data points that are too close to each other 120 (Abbas et al., 2023) or to their respective cluster centroids (Sorscher et al., 2022). Similarly, our 121 work uses DPPs to model data subsets and relies on a greedy MAP algorithm (Chen et al., 2018) 122 to identify diverse subsets. The choice of distance metric and data representations is crucial. Prior 123 works have employed the ℓ -2 distance between neural network activations (Sener & Savarese, 2018; 124 Sorscher et al., 2022; Abbas et al., 2023) or between weight gradients of the log likelihood (Huang 125 et al., 2016; Ash et al., 2019). Here we investigate several data similarity measures on instruction 126 tuning datasets. 127

While choosing diverse subsets is driven by the notion that similar data points are redundant, an 128 alternative approach is motivated by the assumption that certain data points provide more value than 129 others. Specifically, many data selection algorithms define a quality score for each data point and 130 select the portion of the dataset with the highest scores. Various quality scoring functions have been 131 proposed for classification tasks, including the norm of the weight gradient (Settles, 2009; Huang 132 et al., 2016; Paul et al., 2021), the number of times an example transitions from correctly classified 133 to misclassified (i.e., "forgotten") during training (Toneva & Sordoni, 2019), the variability of the 134 ground-truth label likelihood over the course of training (Swayamdipta et al., 2020), and the average 135 ℓ -2 norm of the classification error vector (Paul et al., 2021). Our approach of modeling data subsets 136 with DPPs accommodates arbitrary scores and aims to strike a balance between choosing data points 137 with high quality scores and ensuring diversity within the selected subset.

138 Our work falls under a growing number of studies that select subsets of instruction tuning datasets 139 to finetune LLMs. Many use quality scores to rank and select data points including simple natural 140 language indicators like coherence (Cao et al., 2023) or perplexity (Li et al., 2023), and the LLM's 141 rating of data points based on metrics such as helpfulness (Chen et al., 2024; Liu et al., 2024). Others 142 select data subsets with sufficient coverage of topics and user intents (Lu et al., 2024). Our approach is closely related to methods that balance quality and diversity, e.g., by solving a variant of the facil-143 ity location problem (Bukharin & Zhao, 2023) or prioritize high quality data points while avoiding 144 duplicates (Liu et al., 2024). Different from prior approaches that solely focus on data selection, 145 we also aim to identify ways to characterize the diversity of instruction tuning datasets, beyond the 146 number of tasks (Wei et al., 2022), topics, and/or user intents (Lu et al., 2024). Specifically, we 147 model data subsets with DPPs that naturally emit a diversity metric over datasets that correlates 148 well with the downstream instruction following performance. This metric is useful for predicting 149 improvements in the instruction following performance and for comparing the dataset diversity.

150 151 152

153 154

155

3 Method

3.1 SUBSET SELECTION WITH DPPS

156 A point process on a set of N items is a probability distribution over all subsets of [N]. A DPP P 157 is a point process where the probability measure is parameterized by a positive semi-definite matrix 158 $L \in \mathbb{R}^{N \times N}$, i.e., $P(Y) \propto \det(L_Y)$ for any subset $Y \subset [N]$. $L_Y \equiv [L_{ij}]_{i,j \in Y}$ is a sub-matrix 159 of L indexed by Y in rows and columns. Intuitively, the diagonal elements of L are related to 160 the marginal probability of including the particular items, i.e., $P(\{i\}) \propto L_{ii}$. The off-diagonal 161 elements of L represents the similarity between items. Similar items are less likely to co-occur, i.e., $P(\{i,j\}) \propto L_{ii}L_{jj} - L_{ij}L_{ji}$.

183

189

190 191

196

197

205

206

209 210

Any positive semi-definite matrix L can be expressed as a Gram matrix VV^T for some matrix $V \in \mathbb{R}^{N \times D}$. Each row of V can be viewed as a feature vector for *i*-th item. The absolute value of the determinant of L_Y is the volume of the parallelepiped spanned by rows of V. Therefore, a high probability subset under V is a subset whose feature vectors span a large volume.

Given a dataset with N items $\{x_n\}_{n=1}^N$, we parameterize a determinantal point process P with a kernel matrix $K \in \mathbb{R}^{N \times N}$ that measures the similarity between data points and possibly a vector $q \in \mathbb{R}^N$ that indicates the quality of each data point. For instance, we can treat the cosine similarity between language models' output token embeddings as the similarity measure and the perplexity of the response conditioned on the instruction as data quality.

To select a moderately large subset of size M, the inner product kernel on features of dimension $D \ll M$ is unsuitable due to rank deficiency. Specifically, any subset Y with $|Y| > \operatorname{rank}(L)$ has zero probability mass $P(Y) \propto \det(L_Y) = 0$ and therefore the size of the most likely subset under P is upper bounded by $\operatorname{rank}(L)$. Instead, we use kernel functions that induce full rank Gram matrices, e.g., the radial basis function (RBF) kernel $K_{ij} = \exp\{-\gamma ||x_i - x_j||^2\}$, where a larger value of γ implies that the repulsive force between data points is more local. For data representations that are normalized to unit length, the radial basis function kernel reduces to $K_{ij} = \exp\{2\gamma x_i^T x_j\}$.

Following Kulesza & Taskar (2010), we define $L_{ij} = K_{ij}q_iq_j$. This is equivalent to scaling the kernel feature map by a scalar quality score. As long as K is positive semi-definite, so is L. This structure enables us to model similarity and quality independently while considering both components during inference. Moreover, the probability of any subset $Y \subset [N]$ factors, i.e.,

$$\log P(Y) \propto \sum_{i \in Y} \log q_i^2 + \log \det(K_Y).$$

The log likelihood is maximized for subsets with high quality (1st term) and diversity (2nd term). Similar to Chen et al. (2018), we introduce a hyperparameter $\lambda \in [0, 1]$ to control the relative importance of diversity and quality:

$$\log P(Y) \propto \lambda \sum_{i \in Y} q_i + (1 - \lambda) \log \det(K_Y), \tag{1}$$

that corresponds to a DPP parameterized by the kernel matrix $L = \text{diag}(e^{\beta q}) K \text{diag}(e^{\beta q})$ with $\beta = \lambda/(2(1-\lambda)).$

Given a data budget M, we pose subset selection as maximum a posteriori (MAP) inference under distribution P with a cardinality constraint:

V

$$^{*} = \underset{Y \subset [N]: |Y| = M}{\arg \max} \det(L_{Y}).$$

$$\tag{2}$$

Although this problem is NP-hard (Ko et al., 1995), the log probability in Equation (1) is submodular (Gillenwater et al., 2012) and therefore Equation (2) can be solved efficiently with a greedy algorithm (Nemhauser et al., 1978) with at least (1-1/e)-approximation guarantee, e.g., near optimal under certain assumptions (Sharma et al., 2015).

We use the greedy MAP inference algorithm (Chen et al., 2018) that grows the set of indices $S_1 \subset \cdots, S_N \subset [N]$ by adding

$$i^*(S) = \arg\max_{i \in [N] \setminus S} \left[\log \det(L_{S \cup \{i\}}) - \log \det(L_S) \right]$$
(3)

to the set at each iteration. We define $L_n \triangleq L_{S_n}$ as shorthand for the kernel matrix L indexed by the greedy solution S_n at the *n*-th iteration ($L_N \equiv L$). The marginal gains $\Delta_1(L) = \log \det(L_1)$ and

$$\Delta_n(L) = \log \det(L_n) - \log \det(L_{n-1})$$

for $n = 2, 3, \cdots$ approximate the rate of change in diversity of selected subsets $\{S_n\}$ over the iterations. Larger marginal gains means the selected item contributes more to the diversity of the already selected subset. The unnormalized probability for the whole dataset is the sum of the marginal gains:

214
215
$$\log \det(L) = \sum_{n=1}^{N} \Delta_n(L).$$

216 3.2 LOG DETERMINANT DISTANCE AS A MEASURE OF DIVERSITY

We propose a novel way to measure dataset diversity that is a byproduct of solving the MAP inference problem in Equation (2). The measure of diversity depends entirely on the kernel that defines the DPP. For a fixed kernel function, we can compare dataset diversity quantitatively.

While $\log \det(L)$ may seem like a natural choice, it is unsuitable for measuring dataset diversity for two reasons. First, $\log \det(L)$ is not invariant to scaling of the kernel matrix, leading to widely different values that complicate the interpretation of results. For instance, if a kernel matrix is scaled by a constant c > 0, $\log \det(cL) = N \log(c) + \log \det(L)$ changes by $N \log(c)$ for the same dataset. Second, $\log \det(L)$ depends heavily on the dataset size, particularly when there are significant marginal gains for each item selected.

227 To address the aforementioned challenges, we devise a metric that measures how far det(L) is 228 from the largest possible value than its absolute value. We operationally define a maximally diverse 229 dataset as a set of vectors on the hypersphere in \mathbb{R}^D that achieves the highest det(R). We use R to 230 denote the reference dataset's kernel matrix computed using the same kernel function $k(\cdot, \cdot)$. To our 231 best knowledge, there is not easy solution to obtain this set of vectors. In practice, we approximate such a set of vectors by sampling vectors randomly on the hypersphere. This yielded det(R) that is 232 in absolute value larger than det(L) derived from every real datasets we tried, and therefore served 233 our purpose of having a common reference to compare the diversity of different datasets to. We 234 define Log Determinant Distance as 235

$$LDD \triangleq \frac{1}{N} \log \frac{\det(R)}{\det(L)}.$$
(4)

The log determinant distance measures the average deviation of the volume of the parallelepiped spanned by the rows of the Gram matrix decomposition of L, i.e., $|\det(L)|$, from the largest possible volume, e.g., $|\det(R)|$. A smaller log determinant distance implies that the dataset is closer to the maximally diverse reference dataset, and therefore is more diverse. Alternatively, we can interpret the log determinant distance as the deficit in the average contribution of a data point to dataset diversity from optimum:

$$LDD \equiv \frac{1}{N} \sum_{n=1}^{N} \left(\Delta_n(R) - \Delta_n(L) \right)$$

The log determinant distance can be readily computed from the determinants of the kernel matrices det(L) and det(R) obtained by running the greedy MAP algorithm (Chen et al., 2018) twice.

Given our assumption that the reference dataset is maximally diverse, i.e., $|\det(R)| \ge |\det(L)|$ for any kernel matrix L, the non-negativity property holds: LDD ≥ 0 . Moreover, it is straightforward to show that the log determinant distance is invariant to scaling of kernels. The log determinant distance is also invariant to permutation of datasets since it is based on matrix determinants. In summary, the log determinant distance possesses favorable properties for measuring the dataset diversity.

255 256 3.3 WEIGHT GRADIENT AS DATA REPRESENTATION

236 237

244 245 246

We use the language model's weight gradient $\nabla_{\theta}\ell(x;\theta)$ of scalar-valued loss function ℓ as the data representation for data point x. As an example, ℓ can be the average log likelihood of the response conditioned on the instructions. For LLMs, the full weight gradient consists of billions of elements, rendering kernel computation infeasible. In this work, we apply Johnson-Lindenstrauss (JL) transforms (Johnson & Lindenstrauss, 1984) twice on weight gradients to reduce their dimensionality.

We first apply JL transform implicitly via Low-Rank Adaptation (LoRA) (Hu et al., 2022) to reduce memory as well as computation since most derivatives are neither stored nor computed. For weight matrix $W \in \mathbb{R}^{m \times n}$ in a fully connected layer, LoRA enforces a rank $r \ll min(m, n)$ update to the weight matrix that is a composition of two matrices: $B \in \mathbb{R}^{m \times r}$ and $A \in \mathbb{R}^{r \times n}$. For input activation $z \in \mathbb{R}^n$, the output activation $h \in \mathbb{R}^m$ after the update is

$$h = (W + \Delta W)z = Wz + BAz.$$

We initialize A to $\mathcal{N}(0, r^{-1})$ to construct a distance preserving random projection matrix and B to zero to preserve the forward pass activations. To obtain a lower-dimensional representation of the

full weight gradient $\nabla_W \ell$, we use LoRA at initialization to compute

 $\nabla_B \ell = \nabla_h \ell \cdot z^T A^T = \nabla_W \ell \cdot A^T.$ (5)

Here, we use $\ell \equiv \ell(x; \theta)$ for brevity. We also explored approaches that use LoRA to project $\nabla_W \ell$ onto a vector of size r, instead of m vectors of size r. For example, we can sum over the rows of $\nabla_B \ell$ or rows of $\nabla_B \ell$ after shifting the *i*-th row by *i* positions. We found these approaches induce a larger pairwise distance error.

Note that A is not applied to the entire weight gradient $\nabla_W \ell$. Instead, each row in $\nabla_W \ell$ of dimension n is projected to the corresponding row in $\nabla_B \ell$ of dimension r. The typical Johnson-Lindenstrauss Lemma also holds in this case.

281

282 283

295

296 297

298

Lemma 3.1. Let $\epsilon, \delta > 0$. If $r = O(\log(1/\delta)/\epsilon^2)$, then $\left| \|\operatorname{vec}(\nabla_{B}\ell)\|_{2}^{2} - \|\operatorname{vec}(\nabla_{W}\ell)\|_{2}^{2} \right| \leq \epsilon$

$$\left| \left\| \operatorname{vec}(\nabla_B \ell) \right\|_2^2 - \left\| \operatorname{vec}(\nabla_W \ell) \right\|_2^2 \right| \le$$

with probability at least $1 - \delta$.

286 The proof in Appendix A.2 involves simple application of the union bound.

We then apply the sparse JL transform to the concatenation of $\operatorname{vec}(\nabla_B \ell)$ for every fully connected layer in the neural network to further reduce storage and compute cost. Using a sparse projection matrix is necessary since concatenated $\operatorname{vec}(\nabla_B \ell)$ is still too costly to work with as it contains mlrentries where l is the number of fully connected layer in the network.

Theorem 3.1 can be extended trivially to include the second JL transform. It immediately follows that the two JL transforms together preserve the pairwise distance between weight gradients, in the same way that a single JL transform does on the entire weight gradient.

- 4 EXPERIMENTS
- 4.1 IMPLEMENTATION DETAILS

299 **Dataset** We employ a collection of instruction tuning datasets to understand the effect of data on 300 model's intruction following performance and to evaluate their relative diversity: FLAN (Wei et al., 301 2022), Self-Instruct (Wang et al., 2023b), Dolly (Conover et al., 2023), Alpaca (Taori et al., 2023), 302 GPT-4Alpaca (Peng et al., 2023), OASST2 (Köpf & Kilcher, 2023), Orca (Mukherjee et al., 2023), 303 UltraChat (Ding et al., 2023), WizardLM (Xu et al., 2024), and ShareGPT. We also evaluate the 304 diversity of preference datasets: OpenAI-Summarization (Stiennon et al., 2020), SHP (Ethayarajh 305 et al., 2022), UltraFeedback (Cui et al., 2024), and HH-RLHF (Bai et al., 2022). For each dataset, 306 we remove examples with sequence lengths greater than 2,048 to ensure the language model learns 307 to generate the end-of-sequence token properly. Except for Dolly and OASST2 that contain fewer examples, the aforementioned datasets are subsampled to 50,000 examples to control for the effect 308 of dataset size. 309

Model & Training In all experiments, we finetune Llama-7b (Touvron et al., 2023) for 3 epochs with learning rate of 2e-5 and a batch size of 128. We use AdamW optimizer with no weight decay and linearly decay the learning rate after warmup for 3% of the total number of training steps.

314 **Evaluations** We evaluate the performance of instruction-following models using a few bench-315 marks that measure model capabilities: factual knowledge across various subjects with Massive 316 Multitask Language Understanding (MMLU) (Hendrycks et al., 2020), reasoning on math problems 317 using Grade School Math (GSM) (Cobbe et al., 2021) and on general reasoning problems with the 318 Big-Bench Hard benchmark (BBH) (Suzgun et al., 2022), multilinguality with TydiQA (Clark et al., 319 2020), and coding skills with Codex-Eval (Chen et al., 2021). We use BENCHMARKS AVG to de-320 note the average performance across all aforementioned benchmarks. We use Alpaca-Eval (Dubois 321 et al., 2023) to evaluate instruction following. Specifically, We use length-controlled win rates, i.e., ALPACAEVAL LC-% WIN, to denote the proportion of times a model's generation is preferred by 322 GPT-4 over davinci-003's response, adjusting for for bias towards longer outputs (Dubois et al., 323 2024). We follow the evaluation procedure in (Wang et al., 2023a) closely.



Figure 2: Studies of the log determinant distance as a measure of diversity of instruction tuning and preference learning datasets. More diverse datasets yield a log determinant distance curve that is closer the LDD = 0 line. Distilling responses from capable large language models improves diversity (1st panel). The diversity of synthetic datasets generated using Self-Instruct (Wang et al., 2023b) increase with better teacher model (2nd panel). Using LLMs to re-write instructions to be more complex (Xu et al., 2024) also improves diversity (3rd panel). Curating instructions from diverse sources like ShareGPT & OASST2 yields consistently higher average marginal gains compared to those curated with less human involvement (4th panel). Preference datasets overall are a lot more diverse than instruction tuning datasets, some with no apparent drop off in average marginal gains (5th panel).

Kernel Function To compute the kernel matrix L, we fix the kernel function to radial basis kernel 349 and vary the data representations. We employ Llama-7b representing decoder-only language model 350 to compute the average output token embeddings (LLAMA EMB) & the weight gradients vectors 351 (LLAMA $\nabla_{\theta} \ell$), and MPNet (Song et al., 2020) representing encoder-only masked language model to 352 compute the average output token embeddings of instructions (MPNET EMB). We normalize these 353 data representations to unit length and use the abbreviation NOT NORM. to imply unnormalized 354 vectors. We will refer readers to Appendix B.1 for details on how γ is selected. 355

356 Log Determinant Distance To compute the log determinant distance of a dataset, we generate 357 a reference dataset by sampling vectors randomly on the surface of a D dimensional hypersphere; 358 D = 4096 for LLAMA EMB and LLAMA $\nabla_{\theta} \ell$, and D = 768 for MPNET EMB. We then use the 359 greedy MAP algorithm (Chen et al., 2018) to obtain the determinants of the kernel matrices det(L)360 and det(R), from which we compute the log determinant distance in Equation (4). 361

362 4.2 **RESULTS: DIVERSITY ASSESSMENT**

To assess the log determinant distance as a diversity measure, we compute the log determinant 364 distance using weight gradient vectors $\nabla_{\theta} \ell$ on 9 instruction tuning datasets and 4 preference learning 365 datasets detailed in Section 4.1. For instruction tuning datasets, ℓ is the average log likelihood of 366 tokens in the response conditioned on the instruction. For preference learning datasets, ℓ is the log 367 odds of the preferred response over an alternative worse response. Appendix C.1 demonstrates that 368 we can compute LDD in reasonable time. We also verified that the stochasticity from sampling 369 random vectors on the hypersphere to construct the reference dataset does not affect the estimate 370 significantly, e.g., average LDD of the Alpaca dataset over 5 runs has a standard deviation of 6e-7. 371

Figure 1 demonstrate that the log determinant distance of instruction tuning datasets is correlated 372 with instruction following performance of models finetuned on these datasets. Figure 4 compares the 373 log determinant distance of datasets computed across different data representations: MPNET EMB, 374 LLAMA EMB, and LLAMA $\nabla_{\theta} \ell$, and illustrates that LLAMA $\nabla_{\theta} \ell$ is the only data representation that 375 provides a useful predictor of instruction following performance. 376

Figure 2 compares the log determinant distance of instruction tuning datasets and preference learning 377 datasets. Using log determinant distance as a proxy for dataset diversity, there are a few takeaways:

340

341

342

343

344

345

346

347 348

324

325

326

327

328 329 330

331

333

334

Table 1: The performance of Llama-7b finetuned on 10k (20%) data subset obtained using our 379 DPP-based and alternative baseline data selection methods. We evaluate finetuned language models' 380 generic abilities BENCHMARK AVG and instruction following abilities ALPACAEVAL on Alpaca and 381 UltraChat. (\uparrow) indicates that data points with higher quality scores are selected. 382

302	1	U	1 7			
383	DATASETS	ALF	PACA	ULTRACHAT		
384	Methods	Benchmark Avg	AlpacaEval LC-% Win	BENCHMARK AVG	ALPACAEVAL LC-% WIN	
385	100% Data	23.2	28	22.9	40	
386	Random	21.6	27	22.6	36	
387	DPP (LLAMA $\nabla_{\theta} \ell$)	21.7	29	23.3	37	
388	DPP (LLAMA $\nabla_{\theta} \ell$) Not Norm.	22.7	15	22.7	36	
389	DPP (LLAMA EMB NOT NORM.) DPP (LLAMA EMB)	22.5	27	23.2	38	
390	DPP (MPNET EMB)	22.9	25 24	23.3	35	
391		22.8	24	23.0	39	
392	$\ \nabla_{\theta} \ell\ _2 (\downarrow)$ Alpagasus Rating (\uparrow)	22.7	31 27	- 23.3	- 38	
393	EL2N (\downarrow)	22.8	24	22.7	37	
394	IFD (\uparrow) Perplexity (\downarrow)	21.5 22.7	27	23.0	38 37	
395	#INPUT TOKENS (↑)	25.1	23	22.9	35	
396	#OUTPUT TOKENS (↑) #Total Tokens (↑)	23.4 20.4	31 30	22.7	38 38	
397	DPP (LLAMA $\nabla_{\theta} \ell$ + #OUTPUT TOKS)	22.8	32	22.6	41	
				1		

398

399 (1) the diversity of datasets improves from distilling responses or both instructions and responses 400 from a performant LLM, (2) distilling from better teacher models improves dataset diversity even 401 more, (3) rephrase instructions to be more complex also improves dataset diversity, (3) curating instructions from diverse sources, e.g., from real users on the internet or large-scale crowdsourc-402 ing, promotes dataset diversity, and (5) preference learning datasets are overall more diverse than 403 instruction tuning datasets. 404

405

406 4.3 RESULTS: DATA SELECTION WITH DPPS

407 In this section, we benchmark our DPP data selection approach on two instruction tuning datasets 408 of varying diversity: Alpaca (Taori et al., 2023) and UltraChat (Ding et al., 2023). The latter is more 409 diverse than the former (Figure 1). The data budget is 20% (10,000) of the total dataset size. 410

411 Baselines include random selection (RANDOM), set-cover based deduplication (DEDUP) (Abbas et al., 2023). We also include several rank-and-select approaches based on the norm of the weight 412 gradient ($\|\nabla_{\theta}\ell\|_2$), ChatGPT ratings of examples (ALPAGASUS RATING) (Chen et al., 2024), 413 (EL2N) (Paul et al., 2021), instruction following difficulty (IFD) (Li et al., 2023), the perplexity 414 of the response conditioned on the instruction (PERPLEXITY), and token counts (#INPUT TOKENS, 415 **#OUTPUT TOKENS, #TOTAL TOKENS).** 416

Table 1 reports the performance of models finetuned on data subsets selected using our method and 417 baselines on BENCHMARK AVG for generic abilities and ALPACAEVAL for instruction following. 418 We provide additional results demonstrating that our DPP-based data selection method works well 419 with different base language models (Table 5), scales to large dataset size (Table 6), handles varying 420 data budgets (Table 7). 421

422 We investigate the effect of data representation choice for our DPP-based approach. Using LLAMA 423 $\nabla_{\theta} \ell$ as the data representation yields the largest improvement in instruction following performance compared to alternative data representations on the Alpaca dataset. While simple data deduplica-424 tion DEDUP(MPNET EMB) yields the largest improvement on UltraChat. Appendix C.7 provides a 425 qualitative analysis of examples selected by the DPP-based approach across different data represen-426 tations. One interesting take-away is that examples that are close in the weight gradient space tend 427 to have similar answer structures (e.g., contain lists or long-form writing), even if they do not share 428 the same topics or keywords. 429

We also assess data selection methods based on quality scores. In general, data selection with EL2N 430 and #INPUT TOKENS results in subsets that perform worse than random subsets while using all other 431 quality scores improve instruction following performance. Retaining examples with small $\|\nabla_{\theta}\ell\|_{2}$,



Figure 3: Performance (ALPACAEVAL LC-% WIN) vs. diversity (LDD using $\nabla_{\theta} \ell$) for 20% subsets of the Alpaca dataset obtained from different data selection methods.

instead of large $\|\nabla_{\theta} \ell\|_2$ typically used in active learning (Park et al., 2022), leads to significant improvement. Similar gains on length-controlled win rates are observed when selecting examples with large #OUTPUT TOKENS, suggesting these gains aren't simply due to longer outputs.

To investigate how DPP-based approach balances diversity and quality, we balance most effective quality score #OUTPUT TOKENS with diversity in the normalized weight gradient space. Figure 5 illustrates the trade-off of our DPP-based selection approach and demonstrates that careful balancing of diversity and quality is necessary to achieve optimal performance. DPP (LLAMA $\nabla_{\theta}\ell$ + #OUT-PUT TOKS) leads to slightly improved instruction following performance compared to alternatives on both Alpaca and UltraChat, exceeding the performance of finetuning on the full dataset with 20% of the data.

Figure 3 shows the performance (ALPACAEVAL LC-% WIN) vs. diversity (LDD using $\nabla_{\theta} \ell$) for 20% subsets of the Alpaca dataset obtained from different data selection methods. AlpacaEval lengthcontrolled win rate and LDD is negatively correlated (e.g., $\rho_{pearson} = -0.82$) on Alpaca; While such correlation diminishes on the more diverse UltraChat dataset (e.g., $\rho_{pearson} = 0.44$). This finding has an intuitive explanation: it's worthwhile to enforce diversity only if there is room for improvement.

460 461

462

432

433 434

440 441

442

443 444

5 CONCLUSION

We introduced a DPP-based approach to select instruction tuning data subsets that provide a flexible 463 framework to integrate different notions of data similarity and quality. We explored several choices 464 of data representations and quality scores and determined their utility. We demonstrated that our 465 approach outperforms baselines in various ablation studies. More importantly, we proposed log 466 determinant distance (LDD) to quantify dataset diversity and demonstrated that LDD is correlated 467 with instruction following performance. We can use LDD to (1) gauge how much data should be 468 kept when selecting data subsets (2) whether we should care to implement algorithms to enforce 469 diversity, and (3) the impact of different data curation strategy on dataset diversity. 470

471 6 LIMITATIONS

Enforcing dataset diversity using DPP-based approach proves beneficial on less diverse datasets
(e.g., Alpaca). This benefit diminishes when applied to more diverse datasets (e.g., UltraChat).
Instruction tuning datasets that are based on crowdsourcing user interactions with strong LLMs
(e.g., WildChat) are pretty diverse to begin with. Enforcing diversity may provide limited gains
and random selection after basic text deduplication may be adequate. However, the log determinant
distance can still be used to understand the diversity of these datasets and to determine whether it is
worthwhile to implement more sophisticated ways to encourage diversity.

Our work suggests how to improve dataset diversity. We emphasize the importance of curating datasets with realistic instructions from diverse sources, e.g., internet user interactions with LLMs.
 If extensive human involvement is cost-prohibitive, an alternative approach is to distill the dataset entirely or re-write partially using the most capable LLMs. Surprisingly, preference learning datasets exhibit greater diversity compared to instruction tuning datasets, even if derived from the same source (e.g., UltraFeedback is curated from FLAN, UltraChat etc.). More work is required to better understand this phenomenon and its implications.

486 REFERENCES

- Amro Abbas, Kushal Tirumala, Dániel Simig, Surya Ganguli, and Ari S. Morcos. SemDeDup:
 Data-efficient learning at web-scale through semantic deduplication, March 2023.
- Jordan T. Ash, Chicheng Zhang, Akshay Krishnamurthy, John Langford, and Alekh Agarwal. Deep
 Batch Active Learning by Diverse, Uncertain Gradient Lower Bounds. In *International Confer*ence on Learning Representations, September 2019.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom Brown, Jack Clark, Sam McCandlish, Chris Olah, Ben Mann, and Jared Kaplan. Training a Helpful and Harmless Assistant with Reinforcement Learning from Human Feedback, April 2022.
- Gantavya Bhatt, Yifang Chen, Arnav M. Das, Jifan Zhang, Sang T. Truong, Stephen Mussmann,
 Yinglun Zhu, Jeffrey Bilmes, Simon S. Du, Kevin Jamieson, Jordan T. Ash, and Robert D.
 Nowak. An Experimental Design Framework for Label-Efficient Supervised Finetuning of Large
 Language Models, May 2024.
- Alexander Bukharin and Tuo Zhao. Data Diversity Matters for Robust Instruction Tuning, November
 2023.
- 508
 509 Yihan Cao, Yanbin Kang, and Lichao Sun. Instruction Mining: High-Quality Instruction Data
 510 Selection for Large Language Models, July 2023.
- Laming Chen, Guoxin Zhang, and Eric Zhou. Fast Greedy MAP Inference for Determinantal Point Process to Improve Recommendation Diversity. In *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.
- 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 *International Conference on Learning Representations*, May 2024.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared 519 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, 520 Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, 521 Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, 522 Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fo-523 tios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex 524 Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, 525 Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, 526 Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating 527 Large Language Models Trained on Code, July 2021. 528
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. Scaling Instruction-Finetuned Language Models, December 2022.
- 536

529

518

 Jonathan H. Clark, Eunsol Choi, Michael Collins, Dan Garrette, Tom Kwiatkowski, Vitaly Nikolaev, and Jennimaria Palomaki. TyDi QA: A Benchmark for Information-Seeking Question Answering in Typologically Diverse Languages. *Transactions of the Association for Computational Linguistics*, 8:454–470, 2020. doi: 10.1162/tacl_a_00317.

570

571

572

576

580

581

582

583

584

585

586

587

- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training Verifiers to Solve Math Word Problems, November 2021.
- 544 Mike Conover, Matt Hayes, Ankit Mathur, Jianwei Xie, Jun Wan, Sam Shah, Ali Ghodsi, Patrick
 545 Wendell, Matei Zaharia, and Reynold Xin. Free dolly: Introducing the world's first truly open
 546 instruction-tuned llm, 2023.
- Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Wei Zhu, Yuan Ni, Guotong Xie, Zhiyuan Liu, and Maosong Sun. UltraFeedback: Boosting Language Models with High-quality Feedback. In *International Conference on Learning Representations*, May 2024.
- Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Shengding Hu, Zhiyuan Liu, Maosong Sun, and
 Bowen Zhou. Enhancing Chat Language Models by Scaling High-quality Instructional Conversations. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 3029–3051, Singapore, December
 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.183.
- Yann Dubois, Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos
 Guestrin, Percy Liang, and Tatsunori Hashimoto. AlpacaFarm: A Simulation Framework for
 Methods that Learn from Human Feedback. In *Thirty-Seventh Conference on Neural Information Processing Systems*, November 2023.
- Yann Dubois, Balázs Galambosi, Percy Liang, and Tatsunori B. Hashimoto. Length-controlled
 alpacaeval: A simple way to debias automatic evaluators, April 2024.
- Kawin Ethayarajh, Yejin Choi, and Swabha Swayamdipta. Understanding Dataset Difficulty with
 \$\mathcal{V}\$-Usable Information. In *Proceedings of the 39th International Conference on Machine Learning*, pp. 5988–6008. PMLR, June 2022.
- Uriel Feige, Vahab S. Mirrokni, and Jan Vondrák. Maximizing Non-monotone Submodular Functions. *SIAM Journal on Computing*, 40(4):1133–1153, January 2011. ISSN 0097-5397, 1095-7111. doi: 10.1137/090779346.
 - Dan Friedman and Adji Bousso Dieng. The Vendi Score: A Diversity Evaluation Metric for Machine Learning. *Transactions on Machine Learning Research*, June 2023.
- Jennifer Gillenwater, Alex Kulesza, and Ben Taskar. Near-Optimal MAP Inference for Determinantal Point Processes. In *Advances in Neural Information Processing Systems*, volume 25. Curran Associates, Inc., 2012.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob
 Steinhardt. Measuring Massive Multitask Language Understanding. In *International Conference* on Learning Representations, October 2020.
 - Or Honovich, Thomas Scialom, Omer Levy, and Timo Schick. Unnatural Instructions: Tuning Language Models with (Almost) No Human Labor. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 14409–14428, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.806.
 - 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*, January 2022.
- Jiaji Huang, Rewon Child, Vinay Rao, Hairong Liu, Sanjeev Satheesh, and Adam Coates. Active
 Learning for Speech Recognition: The Power of Gradients. *NIPS Workshop*, 2016.
- 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 LM Adaptation with Tulu 2, November 2023.

594 595 596 597	William B. Johnson and Joram Lindenstrauss. Extensions of Lipschitz mappings into a Hilbert space. In Richard Beals, Anatole Beck, Alexandra Bellow, and Arshag Hajian (eds.), <i>Contemporary Mathematics</i> , volume 26, pp. 189–206. American Mathematical Society, Providence, Rhode Island, 1984. ISBN 978-0-8218-5030-5 978-0-8218-7611-4. doi: 10.1090/conm/026/737400.
598 599 600	Chun-Wa Ko, Jon Lee, and Maurice Queyranne. An Exact Algorithm for Maximum Entropy Sampling. <i>Operations Research</i> , 43(4):684–691, 1995. ISSN 0030-364X.
601 602 603	Andreas Köpf and Yannic Kilcher. OpenAssistant Conversations - Democratizing Large Language Model Alignment. In <i>Conference on Neural Information Processing Systems Datasets and Bench-</i> <i>marks Track</i> , November 2023.
604 605 606	Alex Kulesza and Ben Taskar. Structured Determinantal Point Processes. In Advances in Neural Information Processing Systems, volume 23. Curran Associates, Inc., 2010.
607 608 609	Alex Kulesza and Ben Taskar. Determinantal point processes for machine learning. <i>Foundations and Trends</i> ® <i>in Machine Learning</i> , 5(2-3):123–286, 2012. ISSN 1935-8237, 1935-8245. doi: 10.1561/2200000044.
610 611 612 613	Quoc Le and Tomas Mikolov. Distributed Representations of Sentences and Documents. In <i>Proceedings of the 31st International Conference on Machine Learning</i> , pp. 1188–1196. PMLR, June 2014.
614 615 616	Ming Li, Yong Zhang, Zhitao Li, Jiuhai Chen, Lichang Chen, Ning Cheng, Jianzong Wang, Tianyi Zhou, and Jing Xiao. From Quantity to Quality: Boosting LLM Performance with Self-Guided Data Selection for Instruction Tuning. September 2023.
617 618 619 620	Wei Liu, Weihao Zeng, Keqing He, Yong Jiang, and Junxian He. What Makes Good Data for Alignment? A Comprehensive Study of Automatic Data Selection in Instruction Tuning. In <i>International Conference on Learning Representations</i> , May 2024.
621 622 623	Keming Lu, Hongyi Yuan, Zheng Yuan, Runji Lin, Junyang Lin, Chuanqi Tan, Chang Zhou, and Jingren Zhou. #InsTag: Instruction Tagging for Analyzing Supervised Fine-tuning of Large Language Models. In <i>International Conference on Learning Representations</i> . arXiv, May 2024.
624 625 626	Subhabrata Mukherjee, Arindam Mitra, Ganesh Jawahar, Sahaj Agarwal, Hamid Palangi, and Ahmed Awadallah. Orca: Progressive Learning from Complex Explanation Traces of GPT-4, June 2023.
628 629 630	G. L. Nemhauser, L. A. Wolsey, and M. L. Fisher. An analysis of approximations for maximiz- ing submodular set functions—I. <i>Mathematical Programming</i> , 14(1):265–294, December 1978. ISSN 0025-5610, 1436-4646. doi: 10.1007/BF01588971.
631 632 633 634 635	Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. In <i>Advances in Neural Information Processing Systems</i> . Curran Associates, Inc., 2022.
636 637 638	Dongmin Park, Dimitris Papailiopoulos, and Kangwook Lee. Active Learning is a Strong Baseline for Data Subset Selection. <i>NeurIPS HITY Workshop</i> , 2022.
639 640 641	Mansheej Paul, Surya Ganguli, and Gintare Karolina Dziugaite. Deep Learning on a Data Diet: Finding Important Examples Early in Training. In <i>Advances in Neural Information Processing</i> <i>Systems</i> , November 2021.
642 643 644	Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. Instruction Tuning with GPT-4, April 2023.
645 646 647	Victor Sanh, Albert Webson, Colin Raffel, Stephen Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Arun Raja, Manan Dey, M. Saiful Bari, Canwen Xu, Urmish Thakker, Shanya Sharma Sharma, Eliza Szczechla, Taewoon Kim, Gunjan Chhablani, Nihal Nayak, Debajyoti Datta, Jonathan Chang, Mike Tian-Jian Jiang, Han Wang, Matteo Manica, Sheng Shen,

648 649 650 651	Zheng Xin Yong, Harshit Pandey, Rachel Bawden, Thomas Wang, Trishala Neeraj, Jos Rozen, Abheesht Sharma, Andrea Santilli, Thibault Fevry, Jason Alan Fries, Ryan Teehan, Teven Le Scao, Stella Biderman, Leo Gao, Thomas Wolf, and Alexander M. Rush. Multitask Prompted Training Enables Zero-Shot Task Generalization. In <i>International Conference on Learning Rep-</i>						
652	resentations, January 2022.						
653	Ozan Sanar and Silvio Savaraca Active Learning for Convolutional Neural Networks: A Core Set						
654	Approach In International Conference on Learning Representations February 2018						
655	rippiouent in International Conference on Learning Representations, Peeraaly 2010.						
657	Burr Settles. Active Learning Literature Survey. Technical Report, University of Wisconsin-						
658	Madison Department of Computer Sciences, 2009.						
659	Dravyansh Sharma, Ashish Kapoor, and Amit Deshpande. On Greedy Maximization of Entropy. In						
660	Proceedings of the 32nd International Conference on Machine Learning, pp. 1330–1338. PMLR,						
661	June 2015.						
662	John Shawe-Taylor and Nello Cristianini Kernel Methods for Pattern Analysis Cambridge Univer-						
663	sity Press, 2004.						
664							
665	Kaitao Song, Xu Tan, Tao Qin, Jianteng Lu, and Tie-Yan Liu. MPNet: Masked and Permuted Pre-						
666	volume 33 pp. 16857–16867 Curran Associates Inc. 2020						
667	volume 55, pp. 10057 10007. Curran Associates, me., 2020.						
660	Ben Sorscher, Robert Geirhos, Shashank Shekhar, Surya Ganguli, and Ari Morcos. Beyond neu-						
670	ral scaling laws: Beating power law scaling via data pruning. Advances in Neural Information						
671	Processing Systems, 55:19525–19550, December 2022.						
672	Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford,						
673	Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. In Advances						
674	in Neural Information Processing Systems, volume 33, pp. 3008–3021. Curran Associates, Inc.,						
675	2020.						
676	Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung,						
677	Aakanksha Chowdhery, Quoc V. Le, Ed H. Chi, Denny Zhou, and Jason Wei. Challenging BIG-						
678	Bench Tasks and Whether Chain-of-Thought Can Solve Them, October 2022.						
680	Swabha Swayamdipta, Roy Schwartz, Nicholas Lourie, Yizhong Wang, Hannaneh Hajishirzi,						
681	Noah A. Smith, and Yejin Choi. Dataset Cartography: Mapping and Diagnosing Datasets						
682	with Training Dynamics. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu (eds.),						
683	(EMNLP) pp. 0275, 0202, Opling, Neurophys. 2020, Association for Computational Linguistics						
684	doi: 10.18653/v1/2020.emnln-main.746.						
685							
686	Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy						
687	Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following liama model, 2023						
688	2025.						
600	Mariya Toneva and Alessandro Sordoni. An Empirical Study of Example Forgetting during Deep						
691	Neural Network Learning. International Conference on Learning Representations, 2019.						
692	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux. Timothée						
693	Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Ar-						
694	mand Joulin, Edouard Grave, and Guillaume Lample. LLaMA: Open and Efficient Foundation						
695	Language Models, February 2023.						
696	Ivor W. Tsang, James T. Kwok, and Pak-Ming Cheung. Core Vector Machines: Fast SVM Training						
697	on Very Large Data Sets. Journal of Machine Learning Research, 6(13):363-392, 2005. ISSN						
698	1533-7928.						
099 700	Yizhong Wang Swaroon Mishra Pegah Alipoormolahashi Yeganeh Kordi Amirreza Mirzaei						
701	Atharva Naik, Arjun Ashok, Arut Selvan Dhanasekaran, Anjana Arunkumar, David Stap, Eshaan Pathak, Giannis Karamanolakis, Haizhi Lai, Ishan Purohit, Ishani Mondal, Jacob Anderson, Kirby						

Kuznia, Krima Doshi, Kuntal Kumar Pal, Maitreya Patel, Mehrad Moradshahi, Mihir Parmar, Mirali Purohit, Neeraj Varshney, Phani Rohitha Kaza, Pulkit Verma, Ravsehaj Singh Puri, Rushang
Karia, Savan Doshi, Shailaja Keyur Sampat, Siddhartha Mishra, Sujan Reddy A, Sumanta Patro,
Tanay Dixit, and Xudong Shen. Super-NaturalInstructions: Generalization via Declarative Instructions on 1600+ NLP Tasks. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 5085–5109, Abu Dhabi, United Arab Emirates, December
2022. Association for Computational Linguistics.

- Yizhong Wang, Hamish Ivison, Pradeep Dasigi, Jack Hessel, Tushar Khot, Khyathi Chandu, David Wadden, Kelsey MacMillan, Noah A. Smith, Iz Beltagy, and Hannaneh Hajishirzi. How Far Can Camels Go? Exploring the State of Instruction Tuning on Open Resources. In *Thirty-Seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, November 2023a.
- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. Self-Instruct: Aligning Language Models with Self-Generated Instructions. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 13484– 13508, Toronto, Canada, July 2023b. Association for Computational Linguistics. doi: 10.18653/ v1/2023.acl-long.754.
- Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du,
 Andrew M. Dai, and Quoc V. Le. Finetuned Language Models are Zero-Shot Learners. In *International Conference on Learning Representations*, January 2022.
- Max Welling. Herding dynamical weights to learn. In *Proceedings of the 26th Annual International Conference on Machine Learning*, ICML '09, pp. 1121–1128, New York, NY, USA, June 2009. Association for Computing Machinery. ISBN 978-1-60558-516-1. doi: 10.1145/1553374.1553517.
- Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. WizardLM: Empowering Large Language Models to Follow Complex Instructions. In *International Conference on Learning Representations*, May 2024.
- Wenting Zhao, Xiang Ren, Jack Hessel, Claire Cardie, Yejin Choi, and Yuntian Deng.
 (InThe)WildChat: 570K ChatGPT Interaction Logs In The Wild. In *International Conference* on Learning Representations, May 2024.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Tianle Li, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zhuohan Li, Zi Lin, Eric Xing, Joseph Gonzalez, Ion Stoica, and Hao Zhang. RealChat-1M: A Large-Scale Real-World LLM Conversation Dataset. In *International Conference on Learning Representations*, May 2024.
- Chunting Zhou, Pengfei Liu, Puxin Xu, Srini Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat,
 Ping Yu, Lili Yu, Susan Zhang, Gargi Ghosh, Mike Lewis, Luke Zettlemoyer, and Omer Levy.
 LIMA: Less Is More for Alignment. In *Thirty-Seventh Conference on Neural Information Processing Systems*, November 2023.
- 744 745
- 746 747
- 748
- 749
- 750 751
- 752
- 753

754

756 A THEORY

758

759

760

775

782 783 784

A.1 APPROXIMATION GUARANTEES OF DPP-BASED DATA SELECTION THAT INCLUDES DIVERSITY AND QUALITY

For completeness sake, we will provide a (1-1/e)-approximation guarantee proof for the greedy maximization of the log determinant function. Specifically, we can make any log determinant function $f(X) = \log \det(K_X)$ monotone by scaling the matrix up such that the minimum eigenvalue is at least 1, as proved & discussed in Proposition 2 in (Sharma et al., 2015), thereby obtaining the (1-1/e) approximation guarantee for the greedy maximization of non-negative, monotone, and submodular functions. In practice, the greedy objective in Equation (3) is invariant to the scaling of the kernel matrix. Therefore, we would have obtained a solution with (1-1/e)-approximation guarantee without actually scale K.

Now, we argue that this approximation guarantee holds when we takes into account of both diversity and quality of data points. Specifically, as long as L is a valid kernel matrix, then the LogDet function $f(X) = \log \det(L_X)$ can be made monotone shown previously. We assumed that k is a valid kernel function, therefore $L_{ij} = K_{ij}e^{\beta q_i}e^{\beta q_j}$ is also a valid kernel function since (1) product of a real-valued function $g(x_i) = e^{\beta q_i}$ is a valid kernel and (2) kernels are closed under product (Refer to Proposition 3.22 in (Shawe-Taylor & Cristianini, 2004)).

776 A.2 PROOF OF THEOREM 3.1

⁷⁷⁷ Let ϵ, δ be given. For notation convenience, let $p \equiv \operatorname{vec}(\nabla_W \ell)$ and $q \equiv \operatorname{vec}(\nabla_B \ell)$. Let $p_1, \dots, p_m \in \mathbb{R}^n$ be rows of $\operatorname{vec}(\nabla_W \ell)$ and $q_1, \dots, q_m \in \mathbb{R}^r$ be rows of $\operatorname{vec}(\nabla_B \ell)$. Due to Equation (5), we have $q_k = Ap_k$ for $k = 1, \dots, m$. Provided $A \sim \mathcal{N}(0, \frac{1}{r})$ and $r = \mathcal{O}(\log(1/\delta)/\epsilon^2)$, the following holds due to Johnson-Lindenstrauss Lemma (Johnson & Lindenstrauss, 1984):

$$\mathbb{P}\left[\left|\left\|q_k\right\|_2^2 - \left\|p_k\right\|_2^2\right| < \frac{\epsilon}{m}\right] \ge 1 - \frac{\delta}{m}.$$
(6)

By union bound,

$$\mathbb{P}\left[\bigcup_{k=1}^{m}\left\{\left|\|q_{k}\|_{2}^{2}-\|p_{k}\|_{2}^{2}\right|>\frac{\epsilon}{m}\right\}\right]\leq\sum_{k=1}^{m}\mathbb{P}\left[\left|\|q_{k}\|_{2}^{2}-\|p_{k}\|_{2}^{2}\right|>\frac{\epsilon}{m}\right]\leq\sum_{k=1}^{m}\frac{\delta}{m}\leq\delta.$$
 (7)

If $\left| \|q_k\|_2^2 - \|p_k\|_2^2 \right| < \frac{\epsilon}{m}$ for all $k = 1, \dots, m$, then

$$\left| \|q\|_{2}^{2} - \|p\|_{2}^{2} \right| = \left| \sum_{k=1}^{m} \|q_{k}\|_{2}^{2} - \sum_{k=1}^{m} \|p_{k}\|_{2}^{2} \right| \le \sum_{k=1}^{m} \left| \|q_{k}\|_{2}^{2} - \|p_{k}\|_{2}^{2} \right| \le \sum_{k=1}^{m} \frac{\epsilon}{m} = \epsilon.$$
(8)

Therefore,

$$\mathbb{P}\left[\left|\|q\|_{2}^{2}-\|p\|_{2}^{2}\right| \leq \epsilon\right] \geq \mathbb{P}\left[\bigcap_{k=1}^{m}\left\{\left|\|p_{k}\|_{2}^{2}-\|q_{k}\|_{2}^{2}\right| < \frac{\epsilon}{m}\right\}\right] \geq 1-\delta \tag{9}$$

where the last inequality is by Equation (7).

800 801 802

804 805

799

796 797 798

B IMPLEMENTATION DETAILS

B.1 Details on how the kernel hyperparameter γ is selected

The goal is to find γ large enough such that the returned subset by f is close to full rank but not too large such that the kernel matrix becomes the identity matrix.

807 808

809

If γ is very small, the gain is also very small negative number. Depending on how the stopping criterionr's tolerance parameter ε is set, the algorithm terminates with a subset S of size |S| ≪ N. This would be problematic if we want to select a larger subset.

Table 2: The instruction following performance of Llama-7b finetuned on 20% data subset of Alpaca 811 with different values for the kernel hyperparamter γ . The subset selection algorithm pretty is robust 812 to the choice of γ . 813

814	γ	AlpacaEval LC-%Win
815	0.1	29.4
816	1.0	28.5
817	10.0	28.7
818		

819

810

820 821

822

834 835

836 837 838

839

• If γ is very large, the gain is close to 0, and therefore there is a risk of choosing examples close to arbitrary order.

823 We use Newton's method to get a rough idea on a γ large enough to cover a significant portion of the 824 dataset. Specifically, we treat the greedy algorithm as a black box function $f: \gamma \mapsto S$ that takes in 825 the kernel hyperparameter γ and returns a subset S. The greedy algorithm will return a subset with size at most equal to the numerical rank of the kernel matrix, and therefore the size of the subset 826 obtained depends on γ . We use Newton's method to find the root of $g: \gamma \mapsto |f(\gamma)| - N$ where N 827 is the dataset size. We pick the final γ to be larger than the γ obtained from the Newton's method, 828 usually a magnitude larger to ensure we can use the same γ for different datasets/data budget, and is 829 some exponent of 10, e.g., 1, 10, 1e-3. For example, when using normalized weight gradient as the 830 data representation, Newton's method returns $\gamma \approx 0.05$ and in the paper we just picked $\gamma = 1$. 831

We emphasize that the range of possible γ is pretty wide and as long as γ does not lie in the two 832 extreme cases, the algorithm is pretty robust to the choice of γ , as illustrated in Table 2. 833

С ADDITIONAL EXPERIMENTS AND RESULTS

Table 3: Time for a single run of greedy DPP MAP algorithm for various N and D on a Nvidia V100 40GB GPU.

N		D	
	256	1024	4096
10k	11 sec	10 sec	11 sec
50k	5.7 min	5.8 min	6.4 min
100k	17.1 min	16.3 min	17.7 min

845 846 847

848

C.1 RUNTIME COST OF COMPUTING THE LOG DETERMINANT DISTANCE

849 We use Chen et al. (2018)'s implementation with O(NMD) time and O(N(M+D)) memory com-850 plexity, respectively, where N is size of the original dataset, M is the subset size desired, and D is 851 dimension of data representation. The costly evaluation of kernel matrix entries at each iteration can 852 be parallelized on the GPU at the memory cost of O(ND). The algorithm is a feasible solution for selecting datasets at the scale of several hundred thousand examples. 853

854 To compute LDD, the time complexity is $O(N^2D)$. Table 3 reports the time of a single run of 855 the greedy DPP MAP algorithm for a few choices of N and D on a Nvidia V100 40GB GPU. 856 Note det(R) is computed only once and acts as a drop-in placement when computing LDD of each dataset. Therefore, the compute cost of obtaining det(R) can be amortized.

859 C.2 COMPARISONS OF DATA REPRESENTATIONS USED TO COMPUTE THE LOG 860 DETERMINANT DISTANCE

861

858

Figure 4 compares the log determinant distance (LDD) of datasets computed across different data 862 representations: MPNET EMB, LLAMA EMB, and LLAMA $\nabla_{\theta} \ell$. The log determinant distance based 863 on weight gradient provides the strongest correlation with instruction following performance. This



Figure 4: Comparison of the log determinant distance computed using different data representations: MPNET EMB (top row), LLAMA EMB (middle row), and LLAMA $\nabla_{\theta} \ell$ with respect to instruction tuning loss (bottom row). The log determinant distance based on weight gradient provides the strongest correlation with instruction following performance.



Figure 5: Vary λ interpolates between enforcing diversity and selecting for quality. The results indicate that the optimal performance for both datasets is achieved when $\lambda = 0.1$, highlighting the importance of accounting for both diversity and quality.

suggests that the weight gradient representation is the most informative data representation for measuring the diversity of instruction tuning datasets.

C.3 PERFORMANCE TRADE-OFF BETWEEN DIVERSITY AND QUALITY

Figure 5 shows the performance of the proposed selection approach DPP(LLAMA $\nabla_{\theta}\ell$ + #OUTPUT TOKENS) with varying λ . The λ parameter balances the contribution of examples to diversity in the normalized weight gradient space with the selection of examples with longer responses. The results indicate that the optimal performance for both datasets is achieved when $\lambda = 0.1$, highlighting the importance of accounting for both diversity and quality.

Table 4: Academic benchmark and instruction following evaluation of Llama-7b finetuned on 10k (20%) data subset of Alpaca (top block) and UltraChat (bottom block). This table compares random selection and full finetuning baseline as well as data selection methods that ensure diversity, quality, or both. (\uparrow) indicates that data points with higher quality score are selected.

023	Methods	ACADEMIC BENCHMARKS					AlpacaEval		
024		MMLU	GSM	BBH	TydiQA	CODEXEVAL	AVG	LC-% WIN	
924			AL	PACA					
920	100% Data	41.3	5.0	32.5	19.9	11.0	23.2	28	
920	RANDOM	32.3	4.8	33.4	22.4	8.5	21.6	27	
927	DPP (LLAMA $\nabla_{\theta} \ell$)	28.8	7.4	34.1	22.6	9.6	21.7	29	
928	DPP (Llama $\nabla_{\theta} \ell$) Not Norm.	41.3	5.3	26.3	25.0	8.5	22.7	15	
929	DPP (LLAMA EMB NOT NORM.)	36.9	5.1	31.6	21.5	8.7	22.1	27	
930	DPP (LLAMA EMB)	37.9	6.9 5 9	29.9	21.0	11.2 8 5	22.5	27	
001	DEDUP(MPNET EMB)	38.0	5.8	32.1	20.0	11.0	22.9	23	
931	$\ \nabla_{\mathbf{x}} \mathbf{x}\ $ (1)	36.0	6.0	33.1	20.0	8.5	22.7	31	
932	$\ \nabla \theta^{\mathcal{L}}\ _{2} (\psi)$ Alpagasus Rating (\uparrow)	36.2	5.2	31.3	19.8	9.1	21.6	27	
933	EL2N (\downarrow)	38.3	6.4	32.0	19.9	12.2	22.8	24	
934	IFD (†)	34.6	4.4	30.9	23.5	6.7	21.5	27	
025	Perplexity (\downarrow)	37.7	4.7	32.9	21.0	11.6	22.7	29	
930	#INPUT TOKENS (\uparrow)	43.8	5.3	34.0	24.8	10.4	25.1	23	
936	#OUTPUT TOKENS (\uparrow)	36.3	7.1	35.4	21.8	9.1	23.4	31	
937	#TOTAL TOKENS ()	25.1	1.2	55.5	21.7	8.5	20.4	50	
938	DPP (LLAMA $\nabla_{\theta} \ell$ + #OUTPUT TOKS)	34.1	6.5	34.2	22.2	11.6	22.8	32	
939			Ulth	RACHAT					
040	100% Data	37.8	7.6	31.9	20.4	10.4	22.9	40	
0.44	RANDOM	36.2	7.5	32.8	19.8	11.0	22.6	36	
941	DPP (LLAMA $\nabla_{\theta} \ell$)	36.2	7.6	34.5	20.1	12.8	23.3	37	
942	DPP (LLAMA $\nabla_{\theta} \ell$) Not Norm.	37.7	8.4	32.6	17.8	11.0	22.7	36	
943	DPP (LLAMA EMB NOT NORM.)	38.0	8.6	34.2	19.6	8.5	23.2	38	
944	DPP (LLAMA EMB) DPP (MoNet Emb)	38.1	8.9	32.8	18.9	12.2	23.3	35	
945	DEDUP(MPNET EMB)	37.4	8.7	31.5	21.2	9.1	23.0	39	
946	$\ \nabla_{\theta}\ell\ _{2}$ (1)	34.7	10.3	32.2	21.4	12.2	23.3	38	
540	$EL2N(\downarrow)$	38.6	7.5	32.4	18.2	11.0	22.7	37	
947	IFD (†)	34.5	10.2	33.4	21.4	7.9	23.0	38	
948	PERPLEXITY (\downarrow)	36.9	7.6	33.0	19.4	11.6	22.8	37	
949	#INPUT TOKENS ([†])	37.7	8.5	33.6	18.5	9.1	22.9	35	
050	#OUTPUT TOKENS (↑) #Total Tokens (↑)	34.5	10.3	31.9	20.8	9.1	22.7	38	
900		1 24.0	0.5	33.3	20.9	2.0	22.1	50	
951	DPP (LLAMA $\nabla_{\theta} \ell$ + #OUTPUT TOKS)	34.8	7.7	32.9	21.5	9.8	22.6	41	

THE PERFORMANCE OF DATA SELECTION APPROACHES ON DIFFERENT DATASETS AND C.4 BASE LANGUAGE MODELS

We provide additional experimental details on the performance of our proposed data selection method and baselines. Table 4 reports the performance of Llama-7b fintuned on 20% data subset of the Alpaca and UltraChat datasets. Table 5 reports the performance of Mistral-7b finetuned on 20% data subset of the Alpaca dataset. Together, these tables provide a comprehensive view of the performance characteristics of our proposed data selection method and baselines, and demonstrate that our claims hold across different datasets and base language models used.

C.5 THE PERFORMANCE OF DATA SELECTION APPROACHES ON LARGE-SCALE DATASET

We provide additional experiments on a data mix that consists of a balanced mixture of a subset of the datasets listed in Figure 1 of the paper, i.e., all but Self-Instruct and GPT4-Alpaca. The resulting mix has approximately 313k data points, with at most 50k data points from each sub-dataset. This is similar to the scale of the Tulu-v2 mixture (with 326k data points) that is a well-known data mixture for finetuning LLMs (Ivison et al., 2023). Table 6 reports the performance of Llama-7b fintuned on different 10k subsets (3% data budget) of the full dataset, demonstrating that our approach scales to a relatively large dataset.

973	Table 5: Academic benchmark and instruction following evaluation of Mistral-7b finetuned on 10k
974	(20%) data subset of Alpaca. This table compares random selection and full finetuning baseline as
975	well as data selection methods that ensure diversity, quality, or both. ([†]) indicates that data points
976	with higher quality score are selected.

Methods	ACADEMIC BENCHMARKS						ALPACAE
	MMLU	GSM	BBH	TydiQA	CODEXEVAL	AVG	LC-% W
100% Data	46.8	18.1	50.1	31.5	29.3	35.8	39
Random	61.1	21.4	49.4	31.2	30.5	39.6	39
DPP (Llama $\nabla_{\theta} \ell$)	60.3	21.2	49.1	33.5	34.8	40.3	44
DPP (LLAMA $\nabla_{\theta} \ell$) Not Norm.	60.1	19.6	49.9	30.0	31.7	39.0	40
DPP (LLAMA EMB NOT NORM.)	58.8	21.1	49.5	29.6	31.1	38.8	39
DPP (LLAMA EMB)	57.6	19.8	49.2	28.3	30.5	37.8	45
DPP (MPNET EMB)	58.3	20.8	48.1	29.1	29.3	38.0	41
DEDUP(MPNET EMB)	60.3	20.4	49.2	31.8	27.4	39.0	40
$\ \nabla_{\theta}\ell\ _{2}(\downarrow)$	59.8	26.3	49.4	21.9	30.5	38.4	45
ALPAGASUS RATING ([↑])	59.8	21.7	49.3	35.9	29.3	40.3	37
EL2N (\downarrow)	58.6	21.6	48.0	29.4	31.7	38.5	42
IFD (↑)	56.3	20.0	48.5	33.8	29.3	38.5	41
Perplexity (\downarrow)	59.3	21.4	48.9	27.5	34.8	38.8	41
#INPUT TOKENS (↑)	60.5	21.1	48.7	34.2	31.1	40.0	41
#Output Tokens (†)	60.7	26.7	50.1	22.4	32.3	39.1	47
DPP (LLAMA $\nabla_{\theta}\ell$ + #OUTPUT TOKS)	60.6	25.0	49.5	30.4	34.8	40.6	48

Table 6: Comparison of different data selection methods on a relatively large dataset (313k data points). Our proposed approach scales with the size of the dataset.

Method	BENCHMARK AVG	ALPACAEVAL LC-%WIN
Random	23.7	29.7
DPP(LLAMA $\nabla_{\theta} \ell$) DPP (LLAMA EMB)	24.6 24.1	28.4 27.2
Perplexity (\downarrow) $\ \nabla_{\theta}\ell\ _2 (\downarrow)$ #OUTPUT TOKENS (\uparrow)	24.7 23.4 22.6	35.8 38.2 39.7
DPP(LLAMA $\nabla_{\theta} \ell$ + #OUTPUT TOKS)	24.0	42.5

1006 C.6 The performance of data selection approaches with varying data budgets

Table 7: Performance of Llama-7b finetuned on subsets of Alpaca with varying data budgets. Ourproposed approach outperforms the baseline at low data budgets.

Method		Data Budget					
	10%	20%	40%	60%	100%		
Random	27.6	25.0	25.9	27.1	28.3		
#Output Tokens (\uparrow)	29.3	30.9	30.4	30.4	28.3		
DPP(LLAMA $\nabla_{\theta} \ell$ + #OUTPUT TOKS)	31.9	32.4	32.6	29.6	28.3		

Table 7 reports the performance of Llama-7b finetuned on subsets of Alpaca with varying data budgets. We select a strong quality-based baseline #OUTPUT TOKENS for comparison. We demonstrate that our proposed approach provides strong performance at different data budgets and outperforms the baselines at low data budgets.

C.7 QUALITATIVE ANALYSIS OF SELECTED EXAMPLES

1025 We visualize nearest neighbors of examples using different data representations to understand what data points would be removed (e.g., those that are similar) if we enforce diversity during selection.

The following data pair would be considered close by all data representations we investigated (MPNET EMB, LLAMA EMB, and LLAMA $\nabla_{\theta} \ell$) because they share keywords and belong to similar tasks:

1029 1030

1031

Instruction: As a personal trainer, create a nutrition plan for an athlete ...

Instruction: As a nutritionist, design a meal plan for a client with a rare genetic condition ...

1032 1033 Unlike other data representations, LLAMA $\nabla_{\theta} \ell$ tends to group examples with similar answer struc-1034 tures (e.g., lists or long-form writing) together, even if they do not share the same topics or keywords. 1035 The following pair is an example:

Instruction: Rename the following folder: Documents_1940. **Response**: 1940_Documents

1037 Instruction: Generate a username for Jamie that must contain the letter J and between 8-12 characters. Response: JaimeJr89

1038 1039

1040

D ANSWERS TO POTENTIAL QUESTIONS

1041
 1042
 1043
 D.1 WHY NOT USE DIVERSITY METRICS SUCH AS VENDI SCORE (FRIEDMAN & DIENG, 2023) ?

1044 Vendi Score $\exp(-\sum_{i=1}^{N} \lambda_i \log \lambda_i)$ depends on eigenvalues $\lambda_1, \dots, \lambda_N$ of the kernel matrix $\frac{1}{N}K$, 1045 where N is the number of examples. Computing the Vendi Score for arbitrary kernel function (e.g., 1046 radial basis kernel that we used in our work) takes $\mathcal{O}(N^3)$ time for eigendecomposition whereas 1047 computing LDD requires $\mathcal{O}(N^2)$ time. This implies that computing Vendi Score on moderately 1048 large datasets would be more prohibitive than LDD. Additionally, Vendi Score introduces additional 1049 constraints on the kernel, i.e., k(x, x) = 1, to ensure invariance to the scaling of the kernel matrix. In contrast, LDD maintains scale invariance without imposing this constraint on the kernel function.

In addition to proposing a diversity measure, we investigate different data representations (or the kernel functions) to determine their utility in the context of instruction tuning. Specifically, Figure 4 demonstrated that LDD based on normalized weight gradient provides the best predictor of instruction following performance. This is an orthogonal contribution that previous works have not tried to tackle. Therefore, even if we could feasibly compute the Vendi Score, we would still rely on the findings in our paper to pick the kernel function.

1057

D.2 WHY NOT USE FACILITY LOCATION FUNCTION FOR SUBSET SELECTION ?

1058 1059

1069

Table 8: Comparison of DPP and Facility Location based data selection approaches. We report ALPACAEVAL LC-% WIN of Llama-7b finetuned on different 20% data budget subset of the Alpaca dataset. The two algorithms are pretty comparable.

1060			
1003	Data Representation		Algorithm
1064	Data Representation	DPP	Facility Location
1065		DII	Taeinty Elocation
1066	Llama $ abla_{ heta}\ell$	28.5	28.2
1067	Llama Emb	26.6	26.3
1068	MPNET EMB	25.2	24.5

Facility location and DPPs share many commonalities for subset selection. For example, they both (1) can be solved with a (1-1/e)-greedy algorithm, (2) can incorporate quality scores, and (3) have the same time complexity $O(N^2)$ if using the greedy algorithm on the entire datasets. We consider our work concurrent to (Bukharin & Zhao, 2023; Bhatt et al., 2024) that uses facility location to select instruction tuning datasets.

We emphasize that we are not attached to the exact data selection framework, instead we care more about using some framework, be it DPP or facility location, to better understand how some properties of the dataset, e.g., diversity, quality or both, influence model performance. In addition to improving upon baselines, our work aims to answer questions that have practical implications, informing practitioners: (1) what diversity metric and quality scores are worth trying for selecting instruction tuning datasets, (2) alternative data representations (e.g., normalized weight gradient) that can be used to better quantify similarity between long-form text, (3) how diverse a new dataset
 is compared to existing ones, and if we should care to implement algorithms to enforce diversity, (4)
 the effect of different data curation strategies on dataset diversity.

For completeness sake, we will provide a comparison of DPP and facility location function based subset selection approaches. We use apricot's implementation of facility location with exact lazy greedy algorithm (Feige et al., 2011). We use the same kernel function for facility location function as used in DPP. Table 8 reports ALPACAEVAL LC-% WIN of Llama-7b finetuned on different 20% data budget subset of the Alpaca dataset. The two algorithms are pretty comparable.

D.3 IS ENFORCING DIVERSITY ALL WE NEED ?

While Figure 1 shows there exists a strong correlation between dataset diversity (LDD) and performance (ALPACAEVAL LC-% WIN) across datasets. We do not claim that diversity is correlated with performance all the time. For instance, diversity is strongly correlated with performance on 20% subsets of a less diverse dataset (e.g., Alpaca) while such correlation diminishes on that of a more diverse dataset (e.g., UltraChat), as shown in Figure 3.

1096 Diversity represents one axis of variation of a dataset, alongside with other factors such as quality. 1097 Knowing when such correlation exists for a specific dataset can help us determine if diversity matters 1098 and if we need to spend effort to enforce diversity. We can get an idea of where a dataset lie in this 1099 diversity spectrum with our proposed diversity measure: LDD.