## CLUES: Collaborative Private-domain High-quality Data Selection for LLMs via Training Dynamics

Wanru Zhao <sup>1\*</sup>, Hongxiang Fan<sup>2,1</sup>, Shell Xu Hu <sup>3</sup>, Bofan Chen <sup>1</sup>, Nicholas D. Lane<sup>1,4</sup>

<sup>1</sup> University of Cambridge, UK <sup>2</sup> Imperial College London, UK

<sup>3</sup> Samsung AI Center Cambridge <sup>4</sup> Flower Labs

## **Abstract**

Recent research has highlighted the importance of data quality in scaling large language models (LLMs). However, automated data quality control faces unique challenges in collaborative settings where sharing is not allowed directly between data silos. To tackle this issue, this paper proposes a novel data quality control technique based on the notion of data influence on the training dynamics of LLMs, that high quality data are more likely to have similar training dynamics to the anchor dataset. We then leverage the influence of the training dynamics to select high-quality data from different private domains, with centralized model updates on the server side in a collaborative training fashion by either model merging or federated learning. As for the data quality indicator, we compute the per-sample gradients with respect to the private data and the anchor dataset, and use the trace of the accumulated inner products as a measurement of data quality. In addition, we develop a quality control evaluation tailored for collaborative settings with heterogeneous medical domain data. Experiments show that training on the highquality data selected by our method can often outperform other data selection methods for collaborative fine-tuning of LLMs, across diverse private domain datasets, in medical, multilingual and financial settings. Our code is released at CLUES.

## 1 Introduction

Large language models (LLMs) training has predominantly relied on the accumulation of vast datasets. Recent observations suggest that even a modest quantity of high-quality diverse data can significantly enhance the instruction following capacity of LLMs. Previously, data quality control relied heavily on manual selection processes [37, 36]. This approach, while being commonly used, rendered scalability challenges due to the substantial labor costs. Recent advancements have seen automated low-quality data filters [3], such as perplexity filters [29] and de-duplication filters [22]. However, their effectiveness in data quality control in more complex environments remains to be explored, where data are spread across silos and locations in different formats and difficult to find.

Collaborative training techniques, such as model merging [11] and federated learning [20], are common paradigms for addressing data-sharing constraints and GDPR [28] compliance. However, data quality control for private data is even more challenging if users are in charge of manually filtering data. We summarize here the two unique challenges: (1) Quality Heterogeneity Some clients may possess a higher proportion of low-quality data compared to others, thus we should not select data from all clients with a fixed selection ratio. (2) Domain Heterogeneity Different data silos may come from different vertical domains, for example, in the multilingual setting, different languages have different quality standards that are never unified.

38th Conference on Neural Information Processing Systems (NeurIPS 2024).

<sup>\*</sup>Corresponding Author: Wanru Zhao (wz341@cam.ac.uk)

In this paper, we propose CLUES (collaborative learning under selection), an automated high-quality data selection method for collaborative fine-tuning of Large Language Models (LLMs), showcasing notable performance improvements in mixed-quality data environments from different private domain data. In these domains, private LLM vendors are supposed to build their specialized applications based on open-source LLMs using their own private data, which represent specialized domains with significant private (e.g., patient records) and public data (e.g., scientific papers). By tracing the training dynamics of each training sample, we leverage public dataset to define an anchor dataset and compute the influence of each training sample on the anchor dataset, and set a global threshold to provide effective collaborative quality controls compared with traditional local data quality selection methods in the following aspects: (1) General: Our method is a general pipeline to improve the generalization performance for LLM fine-tuning. It has an interpretation in terms of bi-level optimization with inner optimization in the client side and outer optimization in the server side to minimize the loss on the anchor dataset. (2) Collaborative: Our method is a collaborative fine-tuning paradigm that can be seamlessly integrated into existing model merging and federated learning frameworks, where the modification occurs on the server side only to incorporate data selection. (3) Scalable: We only employ an approximation to solve the bi-level optimization, which makes it scalable to LLMs.

We evaluate our proposed method on medical, multilingual and financial Question Answering (QA) datasets, demonstrating significant improvements of up to 67.3% on challenging medical and financial QA datasets, highlighting the effectiveness of our proposed method. Through extensive analyses, we demonstrate the significant impact of leveraging training dynamics on the collaborative data quality control of LLMs.

## 2 Problem formulation: Collaborative Data Quality

#### 2.1 Related Work

Collaborative LLM Fine-Tuning Paradigms: Model Merging and Federated Learning Collaborative fine-tuning exhibits certain advantageous properties as a distributed machine learning paradigm by shifting the traditional model training process towards sharing model parameters instead of raw data. Participating clients train models using their own private datasets locally, and the updated model parameters are aggregated on the server. This preserves the privacy of the underlying data while collectively benefiting from the knowledge gained during the training process [20]. We focus on merging fine-tuned models that are optimized from the same pre-trained backbone. Different fine-tuned models initialized from the same pre-trained model effectively share a part of the optimization trajectory and can often be merged without accounting for permutation symmetry [40, 11, 17]. Therefore, merging fine-tuned models can improve performance on a single target task [13, 6], improve out-of-domain generalization [2, 1], create multitask models from different tasks [23], and other settings [23, 4].

One of the most significant challenges plaguing model merging and federated learning methods in previous research is the concern that the model parameters might interfere with each other during weighted averaging or other merging operations. This undesirable interaction could potentially lead to a merged model that performs worse than the individual models before merging. We argue that it can be tackled from the perspective of data attribution.

**Model merging.** let  $f_{\theta} \in \mathcal{F}$  denote the language model and  $D_k \in \mathcal{D}$  denote the training dataset on client k. Given the training datasets  $D_k$ , we can define a model merging operator  $\mathcal{M}_K(\cdot; D_k, k \in K = \{1, \cdots, n\}) : \mathcal{F} \to \mathcal{F}$ . The model merging process can be expressed as

$$f_{merging} = \mathcal{M}_K(f)$$

**Federated Averaging.** Based on the notation of model merging, the federated averaging process can be expressed as

$$f_{fed} = (\prod_{t=1}^{T} \mathcal{M}_{S_t(K)})(f)$$

where T is the round number.  $S_t(K)$  is the index set of clients participated in the training in round t.

**Data Attribution and Selection for LLMs** The quality of the training data of a machine learning model can have a significant impact on its performance. One measure of data quality is the notion

of valuation, i.e., the degree to which a given training example affects the model and its predictive performance. Although data attribution is a well-known concept for researchers, the complexity behind large language models, coupled with their growing size, features, and datasets, has made quantification difficult. Recent methods include Perplexity Score, IFD [24], and DataInf [21], etc. More details are provided in Appendix E. However, those data attribution above have not been used in collaborative settings where each client has statistical heterogeneous and quality heterogeneous private-domain data. And previous data selection methods have not provide a way to determine the golden threshold to decide whether a training data sample should be kept or filter out.

**Training Dynamics** Previous works [38, 25, 35] that analyze training dynamics focus primarily on supervised learning and are largely model- and data-agnostic. Swayamdipta et al. [34] empirically demonstrated the influence of data by visually mapping individual training samples according to their impact on the correctness, confidence, and variability of a model.

#### 2.2 Assumption and Objective: Collaborative High-quality Data Selection for LLMs

**Definition 1.1** (Data Quality on Specific Domain k). Given a model architecture  $\theta$ , a training configuration (optimizer, etc.), and a validation set  $D_{val}$  in a specific domain k, the quality of training data z is defined as follows: for  $z_1, z_2 \in \mathcal{D}_{train}$ , if  $\mathcal{L}_{val}(\theta(z_1), D_{val}) < \mathcal{L}_{val}(\theta(z_2), D_{val})$ , then the quality of  $z_1$  is considered higher than that of  $z_2$ . Here,  $\mathcal{L}_{val}$  denotes the validation loss. In other words, the lower the validation loss, the higher the data quality.

**Definition 1.2** (Data Quality in Collaborative Private Domains). Given a model architecture  $\theta$ , a training configuration (optimizer, etc.), and a validation set  $D_{val} = \mathcal{D}_{val}^{(1)}, \mathcal{D}_{val}^{(2)}, \ldots, \mathcal{D}_{val}^{(K)}$  for all K tasks, the quality of training data z is defined based on the validation loss of the global model  $\theta_{merged}$  on  $D_{val}$ . Specifically, for  $z_1, z_2 \in \mathcal{D}_{train}^{(k)}$ , if  $\mathcal{L}_{val}(\theta_{merged}(z_1), D_{val}) < \mathcal{L}_{val}(\theta_{merged}(z_2), D_{val})$ , then the quality of  $z_1$  is considered higher than that of  $z_2$ . As in the single-domain case, lower validation loss indicates higher data quality.

Remarks 1 (Impact of Low Quality Data in Collaborative Private Domains). We manually construct low-quality data samples on each client. We change the proportion of low-quality data from 0% to 100%. Higher scores indicate better performance. From Fig. 1, a larger portion of low-quality data results in higher validation loss, and more unstable and less effective training loss curve. Fig. 2 shows the performance drop when we change the proportion of low-quality data from 0% to 60%.

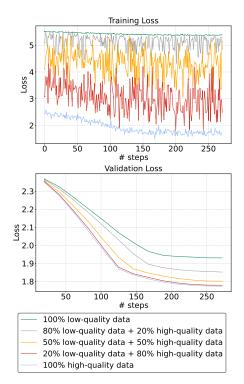
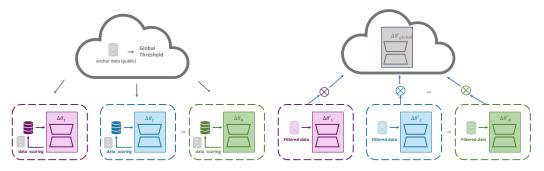


Figure 1: Validation loss and training loss.



Figure 2: Performance drop on the performance of collaborative fine-tuning of LLMs when we change the proportion of low-quality data from 0% to 60%. Higher scores indicate better performance.



Step One: Local Training for Data Quality Scoring Step Two: Collaborative Training with High-Quality

Figure 3: Overall workflow diagram consists of two phases: 1) Step One: client-side computes each sample's quality score with scoring functions using the public validation set and global model, then server-side calculates the score of a global threshold by anchor data 2) Step Two: clients filter data according to the global threshold and starts collaborative learning on selected high-quality data with adaptive weights on the model side.

Remarks 2 (Enhancing Data Quality on Collaborative Private Domains). In the collaborative learning framework, the ratio and distribution of low-quality data are unknown a priori. Only the server has access to the global distribution of both high-quality and low-quality data, while individual clients cannot infer the global distribution from their local distributions due to statistical heterogeneity. The server can infer the distribution of high-quality data from public anchor data. Our objective is to select data points that most significantly reduce the validation loss of the global model, rather than optimizing for each local model independently. It is important to note that the scope of this study does not consider new models joining during training or continual learning paradigms.

## 3 Methodology: CLUES

#### 3.1 Overview

In our workflow, each client performs local training using his own high-quality private data. We have a public validation set located on both the clients and the server, which consists of commonly recognized, high-quality public data. As illustrated in Figure 3, the overall workflow consists of two phases designed to achieve data quality control in the collaborative development of LLMs.

**Step One. Local Training for Data Quality Scoring** Local clients compute each sample's quality score via our training dynamics-based methods using the public validation set and their own fine-tuned model. The server determines a global threshold score, serving as a unified standard of data quality with only a very small amount of anchor data, and sends it to the clients.

Step Two. Collaborative Learning with High-Quality Data Each client then discards data samples that fall below the global threshold received, ensuring that only high-quality data verified by the unified standard are retained. The clients then utilize the high-quality filtered data sets  $\mathcal{D}_k'$  (where  $|\mathcal{D}_k'| \leq |\mathcal{D}_k|$ ) and the initial global model  $\theta^0$  for collaborative learning. After local fine-tuning with the selected high-quality curation data, clients send their local LoRA adapter to the server. The server then aggregates the LoRA parameters of the individual models.

#### 3.2 Step One: Training Dynamics-based Data Scoring

The idea behind our method is straightforward — trace the training process to capture changes in prediction as individual training examples are visited.

For each client, we have designed a data scoring step to calculate the score for each training data sample to measure its contribution to model prediction. Specifically, considering the training set of examples  $\mathcal{D}_k = \{z_1, \dots, z_K\}$  and a model  $\theta$ , we represent the validation set as  $\mathcal{D}'_k = \{z'_1, \dots, z'_K\}$ . We measure the performance of a model using a loss function  $\ell: \mathbb{R}^p \times Z \to \mathbb{R}$ . The loss of the model noted by  $\theta$  on an example z is given by  $\ell(\theta, z)$ . We fine-tune the model by finding parameters  $\theta$  that minimize the training loss  $\sum_{i=1}^K \ell\left(\theta, z_i\right)$ , through an iterative optimization procedure, such as

Stochastic Gradient Descent (SGD) or its variant, which utilizes one training example  $z_t$  in iteration t, updating the parameter from  $\theta_t$  to  $\theta_{t+1}$ :

$$\boldsymbol{\theta}^{t+1} - \boldsymbol{\theta}^t = -\eta_t \nabla \ell \left( \boldsymbol{z}; \boldsymbol{\theta}^t \right) \tag{1}$$

We trace the training process to capture changes in prediction as individual training examples are visited. The contribution of a particular training example z on a given test example z' is defined as the total reduction in loss on the test example z' that is induced by the training process whenever the training example z is utilized. We define the data quality of a particular training example z as the sum of the contribution of the whole validation dataset.

The simplified expression for data quality is as follows:

$$S(z) = \sum_{z'} \sum_{t=1}^{T} \bar{\eta}_i \nabla \ell(z', \theta_t) \cdot \nabla \ell(z, \theta_t)$$
 (2)

The per-sample gradients are calculated for each training sample from the checkpoint t saved during the model training. LLMs are generally tuned using AdamW, which has a more complicated update formula involving the moving averages of the gradient moments.

For Adam,

$$\boldsymbol{\theta}^{t+1} - \boldsymbol{\theta}^{t} = -\eta_{t} \mathcal{L}\left(\boldsymbol{z}, \boldsymbol{\theta}^{t}\right), \mathcal{L}\left(\boldsymbol{z}, \boldsymbol{\theta}^{t}\right) \triangleq \frac{\boldsymbol{m}^{t+1}}{\sqrt{\boldsymbol{v}^{t+1}} + \epsilon}$$
(3)

For AdamW,

$$\boldsymbol{\theta}^{t+1} - \boldsymbol{\theta}^{t} = -\eta_{t} \mathcal{L}\left(\boldsymbol{z}, \boldsymbol{\theta}^{t}\right), \mathcal{L}\left(\boldsymbol{z}, \boldsymbol{\theta}^{t}\right) \triangleq \frac{\boldsymbol{m}^{t+1}}{\sqrt{\boldsymbol{v}^{t+1}} + \epsilon} + \lambda \boldsymbol{\theta}^{t}$$
(4)

Therefore, the training data quality score for LLMs is calculated using the following formula:

$$S(z) = \sum_{z'} \sum_{t=1}^{T} \bar{\eta}_{i} \mathcal{L}(z', \boldsymbol{\theta}_{t}) \cdot \mathcal{L}(z, \boldsymbol{\theta}_{t})$$
(5)

The dot product of the loss gradients of the training example (z) and the test example (z') is weighted by the learning rate  $(\eta_i)$  at different checkpoints and summed up, where we implemented applying point-wise loss gradients to disentangle the relative contributions of each training example. We use the output of the checkpoints from the learning algorithm to capture the training process. The higher the score S(z), the higher the quality of the training sample z. We demonstrates an optimized training approach for collaborative learning of multiple models. By selecting high-quality training data for each local model, we select gradients that positively impact loss trajectories. These trimmed gradients accumulate, leading to an improved position in the weight space. Considering interference during our data selection (gradient selection) of  $\Delta\theta_1'$  and  $\Delta\theta_2'$ , we reduce the interference of weight updates from different models. After parameter aggregation, the merged model  $\Delta\theta_{merged}$  can be improved to an enhanced position in the weight space represented by  $\Delta\theta_{targeted}$ .

It is particularly well-suited for parameter-efficient fine-tuning techniques such as Low-Rank Adaptation (LoRA) [10], which involves freezing the pre-trained model weights and injecting trainable rank decomposition matrices into linear projects of the Transformer architecture. A neural network contains many dense layers that perform matrix multiplication. In the self-attention module, we denote the query projection matrices as  $W_q$ , the key projection matrices as  $W_k$ , the value projection matrices as  $W_v$ , the output project matrices as  $W_o$ . In principle, we can apply LoRA to any subset of weight matrices in a neural network to reduce the number of trainable parameters. In the Transformer architecture, there are four weight matrices in the self-attention module  $(W_q, W_k, W_v, W_o)$ . In our implementation, we apply LoRA only to  $W_q$  and  $W_v$  in most experiments for simplicity.

One straightforward solution is to calculate the quality scores on all weight parameters of LoRA, but may be computationally infeasible when larger models with several millions of parameters are used. To address the memory bottleneck of calculating and saving gradients, we take gradients with respect to a given layer. We propose operating on the first layer of the model, which contains the

least cancelation effect, since the early layers encode *unique logit*. Therefore, we develop the idea of LoRA-based training-data influence in the context of gradient descent. Our proposed influence score is scalable due to the sparse nature of low-rank gradients and contains both low-level and high-level information since the gradient to the low-rank layer can capture both high-level and low-level information about the input sentence.

Note that the above gradient computation process is based on one single checkpoint and there is no parameter update throughout the process. Hence, for each training data point, we can perform this process in parallel, which can facilitate the computation.

## 3.3 Step Two: Global Standard with Anchor Data Scoring

On the server, we use a small set of public data (10 samples in our paper) as our anchor data and calculate the average score of these 10 data points as the global threshold. This establishes a unified standard for division between low- and high-quality data for heterogeneous clients, allowing for the further filtering of local data.

Then we merge the parameters of individual models with adaptive weights on different models. For model merging techniques, we implemented  $Task\ Arithmetic\ [12]$  on task weights, the LoRA matrices are involved in weighted sum. In task arithmetic, one first computes the task weights which is the difference between fine-tuned and base model weights, then calculates a weighted sum of these task weights. Here, the delta weights considered are the individual matrices A and B instead of their product BA. Consider two LoRA adapters  $(A_1, B_1)$  and  $(A_2, B_2)$  along with weights  $w_1$  and  $w_2$  for the weighted merging of these two adapters, then the merging happens as follows:

$$A_{\text{merged}} = \sqrt{\mathbf{w}_1} A_1 + \sqrt{\mathbf{w}_2} A_2 \tag{6}$$

$$B_{\text{merged}} = \sqrt{\mathbf{w}_1} B_1 + \sqrt{\mathbf{w}_2} B_2 \tag{7}$$

We also implement a more efficient method for merging LoRA adapters by eliminating redundant parameters: *TrIm*, *Elect*, *and Merge (TIES) [42]*. First, redundant parameters are trimmed, then conflicting signs are resolved into an aggregated vector, and finally, the parameters whose signs are the same as the aggregate sign are averaged. This method takes into account that some values (redundant and sign disagreement) can degrade performance in the merged model.

## 4 Experiments

Unlike traditional data quality selection methods for pre-trained models or traditional fine-tuning, in our collaborative setting, the training data from vertical domains is very sensitive and subject to strict restrictions regarding sharing and privacy. Therefore, we propose a new experimental setting using medical domain data for downstream tasks and evaluation for open-ended medical QA tasks, considering both quality heterogeneity and domain heterogeneity.

## 4.1 Experimental Setup

**Tasks and Datasets** We conduct our evaluation on the open-ended question-answering (QA) tasks.

- (1) Medical QA: PMC-LLama [41] and Medalpaca-flashcards [7] cover medical question-answering, rationale for reasoning, and conversational dialogues, comprising a total of 202M tokens. We use 16k samples in total, with 8k samples randomly sampled from PMC-LLama and Medalpaca-flashcards each. We uniformly partition the total samples into 20 clients in this task to demonstrate the effectiveness of CLUES in terms of the scalability of the clients, where the clients are IID subsets of the original distribution. For low-quality data, 3.2k samples (40% total data) are polluted with cutting, deletion, or substitution. These 40% low-quality data, together with the rest of the high-quality data, composites the mix-quality data set.
- (2) Multilingual QA: MMedBench [33] is a medical muti-choice dataset of 6 different languages. It contains 45k samples for the trainset and 8,518 samples for the testset. Each question is accompanied by a right answer and high-quality rationale. We use 6312 samples randomly sampled from MMedBench and 1052 samples per language for each of the 6 clients. For the low-quality data, a certain ratio is either polluted with random noise.

(3) Financial QA: To demonstrate the generalizability of our proposed method across various domains, we also include FiQA [5], part of the training corpus of FinGPT [43], which consists 17.1k financial open Question-Answering instructions. We randomly sample 2000 data samples for each of the 4 clients from FiQA dataset, and pollute each of them with a low-quality data ratio of 80%, 20%, 10%, and 50% respectively.

Note that for all tasks, the anchor data and validation dataset used in our proposed method are selected as a held-out high-quality dataset from the same data source.

**Models** We use LLama2-7b [37] and Mistral-7b [14] as our pre-trained models, and fine-tune them with Low-Rank Adaptation (LoRA) [10] on each of the client side. As for the model merging technique, in our main experiments, we use TIES merging. We also compare it with Task Arithmetic in our ablation studies.

**Baselines** The *Oracle* shows the results that train only on the remaining high quality data in the mixed-quality dataset, serving as the theoretical upper bound. We implement the three baselines: existing methods mentioned in 2.1: Perplexity score, IFD [24], and DataInf [21] independently on each client.

**Evaluation metrics** The evaluations focus on two main aspects: (1) Question-Answering capabilities, assessed by GPT-4 [30] scoring within the test set splited from the same sources of the training dataset. In the medical QA and multilingual QA tasks, 200 samples are randomly selected from the medical dataset to serve as the test set. We evaluate the models that need to be compared on the test set to generate responses respectively. Then we use the OpenAI GPT-4 model API to assign scores to their responses. Each response of is rated by the judge on a scale from 0 to 1, reflecting how well the answer aligns with the ground truth. In our financial QA task, GPT-4 rate the responses of the fine-tuned model on our data set on a scale of 1 to 10, reflecting criteria including relevance, precision and fluency. To address potential positional bias, we send our response along with the benchmark output to GPT-4 twice, with different orders. We then calculate the average of these scores as the final performance score. (2) Knowledge acquisition, measured by average accuracy of responses to multiple-choice questions in the MMLU clinical topics [8, 9], MedMCQA [31], PubMedQA [16], and USMLE [15] datasets. Although the goal of private domain fine-tuning is not to increase knowledge, there shouldn't be too much knowledge forgetting during this process. (3) Data selection correctness Precision, Recall, F-1 Score, and Accuracy are widely-used evaluation metrics that provide complementary insights into the model's effectiveness from the data selection perspective. In our case, positive instances represent high-quality data, while negative instances represent low-quality data. Precision quantifies the proportion of correctly identified positive instances among all instances predicted as positive, while Recall measures the proportion of correctly identified positive instances among all actual positive instances in the dataset. The F1 Score offers a balanced measure of Precision and Recall, while Accuracy reflects the overall correctness of our data selection (based on our data scoring and threshold determining method) across all classes.

## 4.2 Main Results

Based on the low-quality dataset setup, we evaluate our data-quality control pipeline in collaborative LLM fine-tuning in both federated (communication round cr=300) and model merging (communication round cr=1) settings. Note that in federated learning, the server and clients need to intensively communicate the model updates during model training. We implement the three baseline methods described in the Section 2.1 to calculate scores for each training data, and set the unified scoring standard using corresponding scoring functions with anchor data.

We demonstrate the performance of data quality control methods in collaborative settings in the medical QA task (Tab. 1) and Multilingual QA task (Tab. 2).

**Federated Learning v.s. Model Merging** Firstly, for both pre-trained models and tasks, with other settings remaining the same, model merging performs better than federated learning. This indicates that loose communication between the local model and the server, compared to frequent communication, might lead to better generalization. Additionally, the performance boost with selected data in the federated setting is larger than in the model merging setting. This might be because during federated learning, we calculate the data score based on the global model (instead of the local model

Table 1: Data selection performance in federated setting on MedicalQA. We **bold** the highest performance and <u>underline</u> the second highest performance for each row.

	Mistral-7b		Llama2-7b	
Evaluation Metric	GPT-4 Scoring	Knowledge Avg	GPT-4 Scoring	Knowledge Avg
Mix-qual Data	0.085	0.194	0.0952	0.311
Oracle	<u>0.160</u>	0.233	0.099	0.440
PPL	0.079	0.346	0.045	0.362
IFD [19]	0.087	0.287	0.050	0.346
DataInf [24]	0.093	0.106	0.103	0.335
CLUES (ours)	<b>0.161</b> (100.6%)	<u>0.309</u> (132%)	<b>0.210</b> (212.1%)	0.356 (80.9%)

Table 2: Data selection performance on MMedBench. We **bold** the highest performance and <u>underline</u> the second highest performance for each row.

	Mist	ral-7b	Llam	na2-7b
Setting	Federated	Model Merging	Federated	Model Merging
Mix-qual Data	0.420	0.515	0.440	0.485
Oracle	0.451	0.530	<u>0.449</u>	0.490
CLUES (ours)	<u>0.435</u> (96.5%)	<u>0.525</u> (99.1%)	<b>0.477</b> (106.2%)	<u>0.487</u> (99.4%)

in the model merging setting) at each timestamp, which can better trace and regularize the training trajectory to the optimal location.

**Data Selection Performance** In both federated and model merging settings, our data selection can achieve over 96% and over 91% of the theoretical upper bound performance, respectively. Our method outperforms the other local data selection baselines under the GPT4 Scoring metrics. Compared to the other methods which cause severe forgetting during instruction tuning, the performance of our method on the Knowledge-based benchmark remains within an acceptable range. This shows that our methods are able to improve domain-specific tasks without forgetting knowledge injected during pretraining.

## 5 Analysis

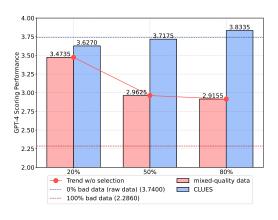
#### 5.1 Qualitative Analysis

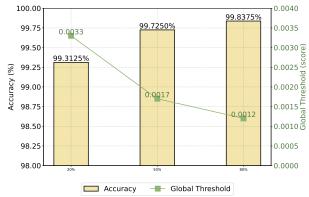
We performed a qualitative analysis by manually comparing the outputs generated by models finetuned on our selected high-quality data versus the original low-quality data. This comparison (Tab. 7 and Tab. 8) provides insights into the tangible improvements in model performance and output quality.

#### 5.2 Varying Levels of Low-Quality Data

To evaluate the robustness of our data selection method under different data quality conditions, we conducted a series of experiments with varying proportions of low-quality data. We maintained a consistent proportion of low-quality data across all clients for each experiment, ranging from 0% to 100%, including pollution levels 20%, 50%, and 80%.

Fig. 4 presents the performance of models trained with and without our data selection method across these different proportions. The results demonstrate that our method effectively enhances data quality across all scenarios with GPT-4 scoring. And in terms of accuracy of the data selection, our method consistently selected over 99% of the high-quality data across different proportions of low-quality data. Additionally, to understand the adaptability of our global threshold, we analyzed how the global threshold changes with different proportions of low-quality data. Fig. 4 illustrates that our global threshold adjusts across varying levels of data quality.





- (a) GPT-4 Scoring Performance in different proportions of low-quality data.
- (b) Selection accuracy and global threshold (score) in different proportion of low-quality data.

Figure 4: Experimental results for different levels of low-quality data

#### 5.3 Quality Heterogeneity

To provide a more comprehensive analysis, in addition to the experiments in the *domain heterogeneity* setting shown above, we conducted additional experiments in a *quality heterogeneity* setting using the FiQA dataset, which focuses on the answer of financial questions. Specifically, we randomly polluted 80%, 20%, 10%, and 50% of the training set for each of the four clients, respectively. The findings demonstrate that our method significantly enhances the quality of the data even when clients have different proportions of low-quality data.

**Varying Merging Techniques** Fig. 5 demonstrates that different weighted merging or aggregation techniques lead to varying performance. Notably, the performance of our data selection method with the *Linear Merging* technique does not even reach the performance of low-quality data with *TIES Merging* technique, highlighting the significant impact of weighted merging techniques on overall performance. Furthermore, we experimented with different merging techniques on the FiQA dataset, demonstrating the importance of weighted merging, shown in Fig. 5.

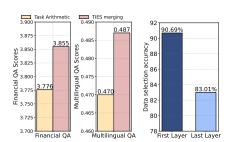


Figure 5: Left: Comparison of different merging techniques. Right: First layer v.s. last layer for low-rank tracing gradient.

# **Layer Selection for Low-Rank Tracing Gradient** In terms of layer selection, we evaluated both the last layer and the *token embeddings*.

We show that layer selection distorts the score (the inner product of two gradients). In our ablation study, we observe that since the activation connected to the last layer of weights contains *shared logic*, the data influenced calculated through the last layer weights are prone to a *cancellation effect*, where the data influence of different examples has a large magnitude that contradicts each other. The cancelation effect lowers the power of the influence score, and deleting influential examples according to this measure often does not change the model's behavior by much. From Fig. 5, we show that the first layer has a less severe cancelation effect than the last layer.

**Unified Scoring with Anchor Data** We conducted an ablation study on our global threshold to further validate our approach. Tab. 3 illustrates the advantage of using a global threshold determined by our anchor data for data selection in this heterogeneous setting, compared to selection based on average ratio or pre-determined scores. These results demonstrate that our approach successfully balances the identification of positive cases with the minimization of false positives, offering a robust and superior solution.

Table 3: Data selection performance on FiQA. We **bold** the highest performance for each row.

	Precision	Recall	F1 Score	Accuracy
Select by ratio	79.17%	79.17%	79.17%	75.00%
Select by a pre-determined score	92.77%	99.13%	95.84%	95.00%
Select by global threshold (Ours)	97.44%	99.38%	98.39%	97.91%

## 6 Discussion and Conclusion

Collaborative model development, including model merging and federated averaging, would benefit from different kinds of high-quality data, and for each of them, the definition of quality is slightly different. In this paper, we establish a data quality control pipeline for collaborative fine-tuning of LLMs, avoiding directly sharing any private data. Our experiments show that the selected high-quality data ensures an effective and reliable learning process, leading to improved model performance.

To the best of our knowledge, we are the first to propose a data selection method for large language models in a collaborative setting, while previous work has mainly focused on traditional centralized settings. We bring up the insights to view federated learning and model merging within the same framework, incorporate different experimental setups and unify federated learning and model merging methods, making it universally applicable. Additionally, our method performs well on generation datasets and takes into account scenarios with bad data, while previous work has not considered downstream domain-specific generation tasks for large language models. Our method does not require repeated training.

**Societal impact** Our work builds large language models that make it possible to create a collaborative instead of a monolithic ecosystem from open-source models while preserving the privacy of users' own data. The constant progress being made in machine learning needs to extend across borders if we are to democratize ML in developing countries. Adapting state-of-the-art (SOTA) methods to resource-constrained environments such as developing countries can be challenging in practice, pushing open source and inclusion.

**Limitations and future work** Our data quality control methods are based on the assumption that all the local models share the same model architectures. It is easy to achieve when our fine-tuning is based on the LoRA adapter. However, it may be worth extending it to adapt to different local model architectures, for example, different low ranks. Future work may explore the intrinsic relation between data selection and the model parameters and how our data selection methods can help reduce the interference of parameter vectors from different models.

#### Acknowledgement

The authors would like to thank Colin Raffel, Haokun Liu, Meghdad Kurmanji and Stefanos Laskaridis for useful discussions and feedback. This research was supported by the European Research Council via the REDIAL project and the Royal Academy of Engineering via the DANTE.

#### References

- [1] Devansh Arpit, Huan Wang, Yingbo Zhou, and Caiming Xiong. Ensemble of averages: Improving model selection and boosting performance in domain generalization. *arXiv* preprint *arXiv*:2110.10832, 2021.
- [2] Junbum Cha, Sanghyuk Chun, Kyungjae Lee, Han-Cheol Cho, Seunghyun Park, Yunsung Lee, and Sungrae Park. Swad: Domain generalization by seeking flat minima. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- [3] Together Computer. Redpajama: an open dataset for training large language models, 2023.
- [4] Shachar Don-Yehiya, Elad Venezian, Colin Raffel, Noam Slonim, Yoav Katz, and Leshem Choshen. Cold fusion: Collaborative descent for distributed multitask finetuning. *arXiv* preprint *arXiv*:2212.01378, 2022.
- [5] FinGPT. fingpt-fiqa\_qa. https://huggingface.co/datasets/FinGPT/fingpt-fiqa\_qa, 2023. Accessed: [Insert Access Date Here].
- [6] Vipul Gupta, Santiago Akle Serrano, and Dennis DeCoste. Stochastic weight averaging in parallel: Large-batch training that generalizes well. *arXiv preprint arXiv:2001.02312*, 2020.
- [7] Tianyu Han, Lisa C Adams, Jens-Michalis Papaioannou, Paul Grundmann, Tom Oberhauser, Alexander Löser, Daniel Truhn, and Keno K Bressem. Medalpaca–an open-source collection of medical conversational ai models and training data. *arXiv preprint arXiv:2304.08247*, 2023.
- [8] Dan Hendrycks, Collin Burns, Steven Basart, Andrew Critch, Jerry Li, Dawn Song, and Jacob Steinhardt. Aligning ai with shared human values. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021.
- [9] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021.
- [10] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *CoRR*, abs/2106.09685, 2021.
- [11] Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Suchin Gururangan, Ludwig Schmidt, Hannaneh Hajishirzi, and Ali Farhadi. Editing models with task arithmetic. *arXiv preprint arXiv:2212.04089*, 2022.
- [12] Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Ludwig Schmidt, Hannaneh Hajishirzi, and Ali Farhadi. Editing models with task arithmetic. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum? id=6t0Kwf8-jrj.
- [13] Pavel Izmailov, Dmitrii Podoprikhin, Timur Garipov, Dmitry Vetrov, and Andrew Gordon Wilson. Averaging weights leads to wider optima and better generalization. arXiv preprint arXiv:1803.05407, 2018.
- [14] Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- [15] Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. What disease does this patient have? a large-scale open domain question answering dataset from medical exams. *Applied Sciences*, 11(14):6421, 2021.
- [16] Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William W Cohen, and Xinghua Lu. Pubmedqa: A dataset for biomedical research question answering. *arXiv* preprint arXiv:1909.06146, 2019.
- [17] Xisen Jin, Xiang Ren, Daniel Preotiuc-Pietro, and Pengxiang Cheng. Dataless knowledge fusion by merging weights of language models. *arXiv preprint arXiv:2212.09849*, 2022.

- [18] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua Bengio and Yann LeCun, editors, 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings, 2015.
- [19] Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In *International conference on machine learning*, pages 1885–1894. PMLR, 2017.
- [20] Jakub Konečný, H Brendan McMahan, Felix X Yu, Peter Richtárik, Ananda Theertha Suresh, and Dave Bacon. Federated learning: Strategies for improving communication efficiency. *arXiv* preprint arXiv:1610.05492, 2016.
- [21] Yongchan Kwon, Eric Wu, Kevin Wu, and James Zou. Datainf: Efficiently estimating data influence in lora-tuned llms and diffusion models. arXiv preprint arXiv:2310.00902, 2023.
- [22] Katherine Lee, Daphne Ippolito, Andrew Nystrom, Chiyuan Zhang, Douglas Eck, Chris Callison-Burch, and Nicholas Carlini. Deduplicating training data makes language models better. *arXiv preprint arXiv:2107.06499*, 2021.
- [23] Margaret Li, Suchin Gururangan, Tim Dettmers, Mike Lewis, Tim Althoff, Noah A Smith, and Luke Zettlemoyer. Branch-train-merge: Embarrassingly parallel training of expert language models. *arXiv preprint arXiv:2208.03306*, 2022.
- [24] 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. *ArXiv*, abs/2308.12032, 2023.
- [25] Yuanzhi Li and Yang Yuan. Convergence analysis of two-layer neural networks with relu activation. *Advances in neural information processing systems*, 30, 2017.
- [26] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019. OpenReview.net, 2019.
- [27] Sourab Mangrulkar, Sylvain Gugger, Lysandre Debut, Younes Belkada, Sayak Paul, and Benjamin Bossan. Peft: State-of-the-art parameter-efficient fine-tuning methods, 2022.
- [28] Christopher F. Mondschein and Cosimo Monda. *The EU's General Data Protection Regulation (GDPR) in a Research Context*, pages 55–71. Springer International Publishing, Cham, 2019. ISBN 978-3-319-99713-1. doi: 10.1007/978-3-319-99713-1\_5.
- [29] Niklas Muennighoff, Alexander M Rush, Boaz Barak, Teven Le Scao, Aleksandra Piktus, Nouamane Tazi, Sampo Pyysalo, Thomas Wolf, and Colin Raffel. Scaling data-constrained language models. *arXiv preprint arXiv:2305.16264*, 2023.
- [30] OpenAI, :, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mo Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan

Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O'Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. Gpt-4 technical report, 2023.

- [31] Ankit Pal, Logesh Kumar Umapathi, and Malaikannan Sankarasubbu. Medmcqa: A large-scale multi-subject multi-choice dataset for medical domain question answering. In *Conference on Health, Inference, and Learning*, pages 248–260. PMLR, 2022.
- [32] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, Advances in Neural Information Processing Systems, volume 32. Curran Associates, Inc., 2019.
- [33] Pengcheng Qiu, Chaoyi Wu, Xiaoman Zhang, Weixiong Lin, Haicheng Wang, Ya Zhang, Yanfeng Wang, and Weidi Xie. Towards building multilingual language model for medicine, 2024.
- [34] Swabha Swayamdipta, Roy Schwartz, Nicholas Lourie, Yizhong Wang, Hannaneh Hajishirzi, Noah A. Smith, and Yejin Choi. Dataset cartography: Mapping and diagnosing datasets with training dynamics. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9275–9293, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.746.
- [35] Yuandong Tian, Tina Jiang, Qucheng Gong, and Ari Morcos. Luck matters: Understanding training dynamics of deep relu networks. *arXiv preprint arXiv:1905.13405*, 2019.
- [36] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models, 2023.
- [37] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas

Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models, 2023.

- [38] Charles L. Wilson, James L. Blue, and Omid M. Omidvar. Training dynamics and neural network performance. *Neural Networks*, 10(5):907–923, 1997. ISSN 0893-6080. doi: https://doi.org/10.1016/S0893-6080(96)00119-0.
- [39] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online, October 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-demos.6.
- [40] Mitchell Wortsman, Gabriel Ilharco, Samir Ya Gadre, Rebecca Roelofs, Raphael Gontijo-Lopes, Ari S Morcos, Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon Kornblith, et al. Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time. In *International conference on machine learning*, pages 23965–23998. PMLR, 2022.
- [41] Chaoyi Wu, Xiaoman Zhang, Ya Zhang, Yanfeng Wang, and Weidi Xie. Pmc-llama: Further finetuning llama on medical papers. *arXiv preprint arXiv:2304.14454*, 2023.
- [42] Prateek Yadav, Derek Tam, Leshem Choshen, Colin Raffel, and Mohit Bansal. TIES-merging: Resolving interference when merging models. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [43] Hongyang Yang, Xiao-Yang Liu, and Christina Dan Wang. Fingpt: Open-source financial large language models. *arXiv preprint arXiv:2306.06031*, 2023.

#### A Data Attribution

Perplexity (PPL) serves as a fundamental metric in language modeling to measure the model's ability to predict text sequences accurately. Mathematically, it is defined as the exponential of the average negative log-likelihood: PPL =  $\exp(-\frac{1}{N}\sum_{i=1}^{N}\log P(w_i|w_1,\ldots,w_{i-1}))$ , where N is the number of tokens in the sequence and  $P(w_i|w_1,\ldots,w_{i-1})$  represents the probability the model assigns to token  $w_i$  given its preceding context. A lower perplexity score indicates better model performance, as it suggests that the model assigns higher probabilities to the correct tokens in the sequence, effectively measuring how "surprised" the model is by new text.

Instruction Following Difficulty (IFD) [24] provides a quantitative metric for evaluating the difficulty of instruction-following tasks in language models. This score is calculated as the ratio between two key measurements: the Conditioned Answer Score  $s_{\theta}(A|Q)$  and the Direct Answer Score  $s_{\theta}(A)$ : IFD $_{\theta}(Q,A) = \frac{s_{\theta}(A|Q)}{s_{\theta}(A)}$ , where  $s_{\theta}(A|Q)$  measures the model's ability to generate responses with instructional context, and  $s_{\theta}(A)$  evaluates the model's capability to generate answers in isolation. The Direct Answer Score is computed as  $s_{\theta}(A) = -\frac{1}{N} \sum_{i=1}^{N} \log P(w_i^A|w_1^A, \dots, w_{i-1}^A; \theta)$ . This metric quantifies the extent to which instructions aid in response generation, where a higher IFD score suggests that the given instruction provides limited useful context for the model's response generation, indicating greater difficulty in following the instruction.

Influence Functions [19]: DataInf represents an efficient algorithm for computing influence functions, distinguished by its closed-form expression that reduces computational and memory complexity compared to existing methods. The algorithm approximates the inverse Hessian calculation  $(G_l(\theta^*) + \lambda_l I_{d_l})^{-1}$  through the key transformation:  $\frac{1}{n} \sum_{i=1}^n \left( \nabla_{\theta_l} \ell_i \nabla_{\theta_l} \ell_i^T + \lambda_l I_{d_l} \right)^{-1} \approx \frac{1}{n} \sum_{i=1}^n \left( I_{d_l} - \frac{\nabla_{\theta_l} \ell_i \nabla_{\theta_l} \ell_i^T}{\lambda_l + \nabla_{\theta_l} \ell_i^T \nabla_{\theta_l} \ell_i} \right)$ , where the Sherman-Morrison formula enables a closed-form solution. The influence function is computed as  $\mathcal{I}_{\text{DataInf}}(x_k, y_k) = \sum_{l=1}^L \frac{1}{\lambda_l} \left( \frac{1}{n} \sum_{i=1}^n \frac{L_{l,i}}{\lambda_l + L_{l,ii}} L_{l,ik} - L_{l,k} \right)$ .

## **B** Priliminary Results on Low-quality Data

Table 4: Preliminary results on MedicalQA.

	Mistral-7b		tral-7b Llama2-7b	
Evaluation Metric	GPT-4 Scoring	KnowledgeAvg	GPT-4 Scoring	Knowledge Avg
Raw Data Mix-qual Data	0.165 0.085 (\dag48.5%)	0.343 0.194 (↓43.4%)	0.265 0.0925 (\dot 65.1%)	0.424 0.311 (\\26.7%)

Table 5: Preliminary results on MMedBench.

	Mistral-7b		Llam	a2-7b
Setting	Federated	Model Merging	Federated	Model Merging
Raw Data Mix-qual Data	0.455 0.420 (\pm,7.69%)	0.540 0.515 (\dagger{6.48%})	0.450 0.440 (\pm2.22%)	0.505 0.485 (\dagger3.96%)

## C Detailed Method Description

**Stage 1 (On each client)** Local fine-tuning with low-quality data; save model checkpoints.

**Stage 2 (On each client)** Calculate gradients, compute scores for each training sample, send scores to the server.

**Stage 3 (On the server)** Calculate gradients, compute scores of anchor data, determine the global threshold using anchor data scores and client scores.

Stage 4 (On each client) Select data with scores not lower than the global threshold.

Stage 5 (On each client) Local fine-tuning with high-quality data, then send model parameters to the server.

**Stage 6 (On the server)** Merge client model parameters to obtain the final global model.

#### **Algorithm 1** Our data selection method for collaborative fine-tuning

```
\begin{array}{ll} \textit{Initialization} & \text{Initial global model:} & \theta^0; & \text{Training datasets (private):} D_{train} = \\ \left\{ \mathcal{D}_{train}^{(1)}, \mathcal{D}_{train}^{(2)}, \dots, \mathcal{D}_{train}^{(K)} \right\}, \mathcal{D}_{train}^{(k)} = \{z_1^{(k)}, \dots, z_n^{(k)}\}; \\ & \text{Validation datasets (public):} & D_{val} = \left\{ D_{val}^{(1)}, \mathcal{D}_{val}^{(2)}, \dots, \mathcal{D}_{val}^{(K)} \right\}, \mathcal{D}_{val}^{(k)} = \{z_1^{\prime(k)}, \dots, z_m^{\prime(k)}\}; \\ & \text{Anchor data (public):} & D_{anc} = \left\{ \mathcal{D}_{anc}^{(1)}, \mathcal{D}_{anc}^{(2)}, \dots, \mathcal{D}_{anc}^{(K)} \right\}, \mathcal{D}_{anc}^{(k)} = \{z_1^{*(k)}, \dots, z_v^{*(k)}\}; \\ \end{array}
```

```
On the Client k

Train model \theta_k on \mathcal{D}^{(k)}_{train}

For each training sample z_i \in \mathcal{D}^{(k)}_{train}

Calculate \nabla \ell \left( \boldsymbol{z}', \boldsymbol{\theta}_t \right)

for each checkpoint t in T do

for each validation sample z_i' \in \mathcal{D}^{(k)}_{val} do

Calculate \nabla \ell \left( \boldsymbol{z}'_i, \boldsymbol{\theta}_t \right)

end for

end for

S(z) = \sum_{z'} \sum_{t=1}^{T} \bar{\eta}_i \nabla \ell \left( \boldsymbol{z}', \boldsymbol{\theta}_t \right) \cdot \nabla \ell \left( \boldsymbol{z}, \boldsymbol{\theta}_t \right)

plate Scoring end for

\mathcal{D}'^{(k)}_{train} = \left\{ z_i \in \mathcal{D}^{(k)}_{train}, z_i \geq \tau \right\}

Select training data with scores above the threshold

Train model \theta_k' on \mathcal{D}'^{(k)}_{train}

Fraining with High-Quality Data Send updated \theta_k' to server
```

```
On the Server for each anchor data z_i \in \mathcal{D}^{(k)}_{anc} do Calculate \nabla \ell\left(\boldsymbol{z}, \boldsymbol{\theta}_t\right) for each checkpoint t in T do for each validation sample z_i' \in \mathcal{D}^{(k)}_{val} do Calculate \nabla \ell\left(\boldsymbol{z}_i', \boldsymbol{\theta}_t\right) end for end for end for S(z) = \sum_{t=1}^T \bar{\eta}_i \nabla \ell\left(\boldsymbol{z}', \boldsymbol{\theta}_t\right) \cdot \nabla \ell\left(\boldsymbol{z}, \boldsymbol{\theta}_t\right) \Rightarrow Anchor data score end for Determine the global threshold \tau with anchor data D_{anc} Send \tau to each clients for training with High-quality data Aggregate client updates: \theta' = \sum_{k=1}^K \frac{|\mathcal{D}_k|}{\sum_{k=1}^K |\mathcal{D}_k|} \theta_k' \Rightarrow Model Merging or Aggregation
```

## D Complexity of Data Scoring

The overall compute complexity, where N is number of checkpoints, and d is gradient dimension.

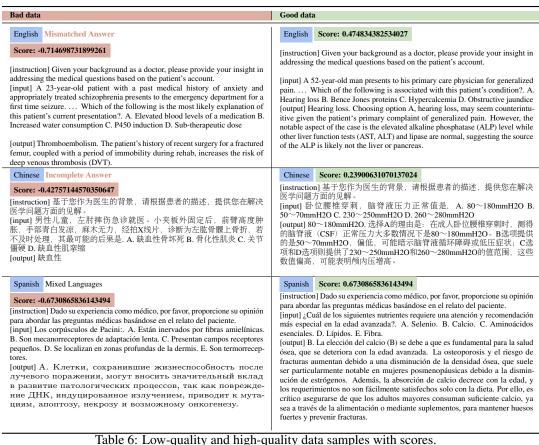
$$\mathcal{O}\left(N\cdot|\mathcal{D}|\cdot|\mathcal{D}_{\text{val}}|\cdot d\right)$$

The overall storage complexity:

$$\mathcal{O}(|\mathcal{D}| \cdot N \cdot d + |D_{val}| \cdot N \cdot d)$$

## **E** Implementation Details

We use Hugging Face's transformers library [39] and PEFT library [27] for loading pre-trained models and LoRA configurations. For our training and evaluation, we use PyTorch as the underlying auto-differentiation framework [32]. We use the AdamW optimizer [26, 18] for all experiments. All the experiments are conducted on the machines with the same hardware configuration using



NVIDIA A40 and A6000. In the experiments of Medical QA, we consider 300 communication rounds, n=20 clients, 10 local steps of model training. In the experiments of Multilingual Medical QA, we distrubute 6 languages on 6 clients, with each clients 1052 dataset. In all experiments, we use 8 bit quantization, set  $batch\_size = 16$ , and LoRA rank to 16 (all the participating LoRA adapters have same rank), and we do 3 epochs for local training.

#### F **Additional Experimental Results**

Questions	Ground Truth Response	Model Response (w/o data selection)	Model Response (w/ our data selection)
初规时察23岁,小观宫开,给100观宫持时,23岁,小观宫持时,100观宫持时,1400处密,1400处密,1400处密,1400处密,1400处密,1400处。14	严密观察产程进展.选的原名《严密观察产程进展》的原因是因为初产妇程宫正在第一个一个一个一个一个一个一个一个一个一个一个一个一个一个一个一个一个一个一个	选手缩前律以展产选产产但产静是法妇的,所有的的间的代况观一法妇项素疗表的的间的代况观一法妇项素疗表的的的间的代况观一法妇项素疗表剖的的的间的代况观一法妇项素疗表剖的的的间的代况观一法妇项素疗表剖的的间的代况观一法妇项素疗表	A.R.一,进产种程产部生征取因观适 为A.R.一,进产种程产部生征取因观适 外区,是一个的成医合采。密合 为A.R.一,进产种程产部生征取因观适 外层,是一产的成医合采。密合 型产宫前示可征观产要帮赈,措产最 型产宫前示可征观产要帮赈,措产最 型产宫前示可征观产要帮赈,措产最 以至现货,并产量 以至现货,并产量。
What type of injury to the arm/elbow most often leads to supracondylar fractures?	Supracondylar fractures most often occur after hyperextension injuries of the arm/elbow.	1000 Hg is the pressure required to rupture an average aortic aortic	Supracondylar fractures are a type of fracture that occurs above the elbow joint, and they are most commonly caused by a fall onto an outstretched arm.

Table 7: Samples of the output of merged models on multilingual MMedBench dataset.

Questions	Model Response (w/o data selection)	Model Response (w/ our data selection)	
How does a 2 year treasury note work?	It is a 2 year treasury note.	Treasury notes are a short-term government security that matures in two years. They are issued by the U.S. Treasury Department to help inance the nation's deicit.	
What causes a stock to drop in price?	Investors are afraid to buy stocks because they are afraid the stock will drop in price.	1	

Table 8: Samples of the output of merged models on FiQA dataset.

## **NeurIPS Paper Checklist**

The checklist is designed to encourage best practices for responsible machine learning research, addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove the checklist: **The papers not including the checklist will be desk rejected.** The checklist should follow the references and follow the (optional) supplemental material. The checklist does NOT count towards the page limit.

Please read the checklist guidelines carefully for information on how to answer these questions. For each question in the checklist:

- You should answer [Yes], [No], or [NA].
- [NA] means either that the question is Not Applicable for that particular paper or the relevant information is Not Available.
- Please provide a short (1–2 sentence) justification right after your answer (even for NA).

The checklist answers are an integral part of your paper submission. They are visible to the reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it

(after eventual revisions) with the final version of your paper, and its final version will be published with the paper.

The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation. While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided a proper justification is given (e.g., "error bars are not reported because it would be too computationally expensive" or "we were unable to find the license for the dataset we used"). In general, answering "[No]" or "[NA]" is not grounds for rejection. While the questions are phrased in a binary way, we acknowledge that the true answer is often more nuanced, so please just use your best judgment and write a justification to elaborate. All supporting evidence can appear either in the main paper or the supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification please point to the section(s) where related material for the question can be found. IMPORTANT, please:

- · Delete this instruction block, but keep the section heading "NeurIPS paper checklist",
- Keep the checklist subsection headings, questions/answers and guidelines below.
- Do not modify the questions and only use the provided macros for your answers.

#### 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: See Section 3 and Section 4.

#### Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

#### 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: See Appendix 6.

#### Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was
  only tested on a few datasets or with a few runs. In general, empirical results often
  depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.

- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

## 3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: The paper does not have any theoretical results.

#### Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and crossreferenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

## 4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: See Section 4 and Appendix 6.

#### Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
- (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.

- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

## 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: See the submitted artifacts.

#### Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/ public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https: //nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

#### 6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: See Section 4 and Appendix 6.

## Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

## 7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: See the bar plots in Section 4.

#### Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error
  of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

## 8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: See Section 4 and Appendix 6.

#### Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

#### 9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: We followed the NeurIPS Code of Ethics.

#### Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

## 10. Broader Impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: We have provided the statement of societal impact of the work performed. Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

## 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: This paper proposes an approach for data selection, and does not pose high risks for misuse.

#### Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with
  necessary safeguards to allow for controlled use of the model, for example by requiring
  that users adhere to usage guidelines or restrictions to access the model or implementing
  safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do
  not require this, but we encourage authors to take this into account and make a best
  faith effort.

#### 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: This paper is licensed under CC-BY-NC-SA 4.0. All other codes, datasets, and references are properly cited.

#### Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.

- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the
  package should be provided. For popular datasets, paperswithcode.com/datasets
  has curated licenses for some datasets. Their licensing guide can help determine the
  license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

#### 13. New Assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: We did not have any released new assets.

#### Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

## 14. Crowdsourcing and Research with Human Subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: We did not involve crowdsourcing nor research with human subjects.

#### Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

# 15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: We did not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.