EFFICIENT ENSEMBLES IMPROVE TRAINING DATA AT TRIBUTION

Anonymous authors

Paper under double-blind review

ABSTRACT

Training data attribution (TDA) methods aim to quantify the influence of individual training data points on the model predictions, with broad applications in data-centric AI, such as mislabel detection, data selection, and copyright compensation. However, existing methods in this field, which can be categorized as *retraining-based* and *gradient-based*, have struggled with the trade-off between computational efficiency and attribution efficacy. Retraining-based methods can accurately attribute complex non-convex models but are computationally prohibitive, while gradientbased methods are efficient but often fail for non-convex models. Recent research has shown that augmenting gradient-based methods with ensembles of multiple independently trained models can achieve significantly better attribution efficacy. However, this approach remains impractical for very large-scale applications.

In this work, we discover that expensive, fully independent training is unnecessary for ensembling the gradient-based methods, and we propose two efficient ensemble strategies, DROPOUT ENSEMBLE and LORA ENSEMBLE, alternative to naive independent ensemble. These strategies significantly reduce training time (up to 80%), serving time (up to 60%), and space cost (up to 80%) while maintaining similar attribution efficacy to the naive independent ensemble. Our extensive experimental results demonstrate that the proposed strategies are effective across multiple TDA methods on diverse datasets and models, including generative settings, significantly advancing the Pareto frontier of TDA methods with better computational efficiency and attribution efficacy. We conduct a theoretical analysis that provides insights into the success of our empirical findings.

031 032 033

034

003 004

010 011

012

013

014

015

016

017

018

019

021

025

026

027

028

029

1 INTRODUCTION

Training data plays an increasingly crucial role in modern artificial intelligence (AI) models (Kaplan et al., 2020). Consequently, data-centric AI emerges as a vital paradigm, emphasizing the collection, curation, and understanding of training data. *Training Data Attribution* (TDA) is a family of methods that assess the influence of each training sample on a model's output. Numerous TDA methods have been developed and applied to a wide range of data-centric AI applications, such as mislabel detection (Koh & Liang, 2017), data selection (Engstrom et al., 2024), and copyright compensation (Deng & Ma, 2023), thereby gaining increasing popularity.

However, accurately attributing the training data influence on very large-scale AI applications remains 043 an open challenge. Existing TDA methods can be generally categorized into two groups: retraining-044 based methods and gradient-based methods (Hammoudeh & Lowd, 2024). Retraining-based methods involve the systematic retraining of the model with and without specific training samples to observe 046 changes in the model's output (Ghorbani & Zou, 2019; Jia et al., 2019; Feldman & Zhang, 2020; 047 Ilyas et al., 2022; Wang & Jia, 2023). Such methods often require thousands of model retraining, 048 sometimes even growing with the size of the training data, to achieve satisfactory performance, which makes them infeasible for moderately large models. Gradient-based methods, on the other hand, estimate the influence by tracking the gradients of data samples (Koh & Liang, 2017; Yeh et al., 2018; 051 Pruthi et al., 2020). These methods are typically computationally more efficient as they do not require the training of multiple models. But empirically they can be brittle in handling complex non-convex 052 models (Basu et al., 2020; Bae et al., 2022) and be sensitive to the randomness associated with model initialization and training dynamics (Søgaard et al., 2021).

Recent studies have shown that gradient-based TDA methods can be significantly improved by
ensembling tens of models independently trained with different random seeds, leading to state-of-theart attribution efficacy on modern neural network models (Søgaard et al., 2021; Park et al., 2023).
Aggregating these independently trained models helps to mitigate the randomness introduced by
training dynamics, such as random initializations and stochastic optimization. However, despite
its impressive performance, scaling such a naive independent ensemble approach further for very
large models remains challenging due to the significant computational costs associated with training
multiple models.

062 In this work, we hypothesize that training models independently is unnecessary for the purpose of 063 TDA ensemble, and we propose two efficient ensemble strategies alternative to the naive independent 064 ensemble approach. Our first strategy, DROPOUT ENSEMBLE, is motivated by dropout (Srivastava et al., 2014), a common deep learning module initially designed to efficiently approximate ensemble. 065 DROPOUT ENSEMBLE reduces the computational costs by replacing the independently trained models 066 in the naive ensemble approach with multiple dropout-masked models using the same original model. 067 Our second strategy, LORA ENSEMBLE, is motivated by an efficient fine-tuning technique, LORA (Hu 068 et al., 2021). Similar to DROPOUT ENSEMBLE, LORA ENSEMBLE reduces the computational costs 069 by replacing the independently trained models with LoRA fine-tuned models from the same original model, which is particularly suitable for generative Transformer models (Vaswani et al., 2017). 071

We evaluate the proposed DROPOUT ENSEMBLE and LORA ENSEMBLE with extensive experiments. 072 Our experiments span various datasets (MNIST (LeCun et al., 1998), CIFAR (Krizhevsky et al., 073 2009), and MAESTRO (Hawthorne et al., 2018)), model architectures (Multi-Layer Perceptrons 074 (MLP), Residual Neural Networks (ResNet) (He et al., 2016), and Transformers (Vaswani et al., 075 2017)), and TDA methods (TRAK (Park et al., 2023), Influence Function (Koh & Liang, 2017), and 076 Grad-Dot/Grad-Cos (Charpiat et al., 2019)). Compared to the naive independent ensemble approach, 077 DROPOUT ENSEMBLE and LORA ENSEMBLE can significantly reduce the training time, serving 078 time, and space costs by respectively up to 80%, 60%, and 80% while maintaining similar TDA 079 efficacy.

Additionally, we present a theoretical analysis to demonstrate the effectiveness of our proposed efficient ensemble strategies over the naive independent ensembling method. Specifically, we first introduce a general ensemble scheme: *two-step ensemble estimator*, which includes both DROPOUT ENSEMBLE and LORA ENSEMBLE. *Two-step ensemble estimator* utilizes a number of base estimators and generates several variants on top of each of these base estimators. Based on mild assumptions, our theory shows that compared to the naive ensemble estimator, *two-step ensemble estimator* can achieve smaller squared error for estimating the optimal attribution score obtained from any given TDA method.

We summarize the contributions of this work as follows:

- We show that fully independent training is not a strict requirement for effective TDA with ensembles, which opens up a promising direction for developing more efficient and effective TDA methods.
- We provide theoretical analysis that compares the TDA method with naive independent ensemble and the two novel efficient ensemble strategies: DROPOUT ENSEMBLE and LORA ENSEMBLE.
- Our proposed DROPOUT ENSEMBLE and LORA ENSEMBLE achieve significant efficiency improvement across a diverse range of machine learning models, datasets, and TDA methods, advancing the state-of-the-art of TDA.

2 RELATED WORK

090

091

092

TDA methods quantify the influence of each training sample on the model predictions by assigning
a TDA score to the sample. These methods can be categorized into retraining-based ones and
gradient-based ones (Hammoudeh & Lowd, 2024), and our work focuses on the latter. Please refer to
Appendix A for more detailed related works.

Ensembling for gradient-based TDA methods. Recent studies have shown the effectiveness of
 ensembling for improving TDA scores computed with gradient-based methods (Søgaard et al., 2021;
 Park et al., 2023), which mitigates their typical issues associated with non-convexity and sensitivity
 to randomness. Ensembling normally applies the TDA method to many independently trained models.
 Either averaging the final TDA scores (Søgaard et al., 2021) or aggregating some intermediate terms
 for score calculation (Park et al., 2023). Besides independently trained models, Park et al. (2023)
 suggest that model checkpoints at different stages of a single training can be used for ensembling as

well. All these ensembling methods, though effective, also require a non-trivial amount of ensembles to perform well. Empirical studies show that their TDA performance suffers significantly when ensemble size is limited (Park et al., 2023). Consequently, the ensemble size, and thus the cost associated with each ensemble, poses a significant barrier to the effective use of ensembling in current TDA methods.

3 Method

Following our hypothesis that independently trained models are unnecessary for ensembling TDA methods, we propose efficient ensemble strategies alternative to the naive independent ensemble.

118 119 120

121

122

161

113 114

115 116

117

3.1 PRELIMINARIES

We start by formalizing the TDA problem and the naive independent ensemble method for TDA.

123 **The TDA problem.** We have a training set $S = \{x_1, \ldots, x_n\}$ where each $x_i \in \mathcal{X}_{\text{train}}$, a test 124 set $\mathcal{T} = \{x_1, \ldots, x_m\}$ where each $x_i \in \mathcal{X}_{\text{test}}$, and a trained model output function f_{Θ} that is 125 parameterized by Θ . Here $\mathcal{X}_{\text{train}}$ and $\mathcal{X}_{\text{test}}$ are the space of the training and test data respectively. We 126 consider both supervised learning and generative modeling settings. For the supervised learning 127 setting, $\mathcal{X}_{\text{train}}$ and $\mathcal{X}_{\text{test}}$ are typically identical, and each element of them corresponds to a data point with both the feature and label. For the generative modeling setting, \mathcal{X}_{train} refers to the space of 128 training data (e.g., text sequence segments for autoregressive language models) while \mathcal{X}_{test} refers to the 129 space of model generation. For the model output function f_{Θ} , typically we will use the model learned 130 from the training set S (i.e., the model parameters will be chosen as $\Theta^* = \operatorname{argmin}_{\Theta} \sum_{x_i \in S} \mathcal{L}(x_i; f_{\Theta})$ for some loss function \mathcal{L}). However, this is not always the case in practice (Pruthi et al., 2020). For a 131 132 training set S and any test sample $x \in \mathcal{T}$, a TDA method τ derives the TDA scores $\tau(x, S; f_{\Theta}) \in \mathbb{R}^n$ 133 to quantify the influence of each training data point in S on the model learned from S. More 134 specifically, the *i*-th element, $\tau(x, S; f_{\Theta})_i \in \mathbb{R}$, is a real-valued score indicating the importance of 135 x_i on the model output on the test sample x. In the supervised classification setting, the "model" 136 output on x" usually refers to the loss, (log-)likelihood, or logit for the model predicting the correct 137 class of x (Koh & Liang, 2017; Park et al., 2023); in the generative setting, it typically refers to the (log-)likelihood for the model generating x (Deng & Ma, 2023). 138

The naive independent ensemble method for TDA. To mitigate the randomness led by the training dynamics of non-convex deep learning models, the naive independent ensemble method first trains a set of I independent models, $\{\Theta^{(i)}\}_{i=1}^{I}$, and then derives ensembled TDA scores $\tau_{\text{ens}}(x, S; \{f_{\Theta^{(i)}}\}_{i=1}^{I}) \in \mathbb{R}^n$ based on the set of models.

The ensembled TDA scores could be simply obtained by averaging over the TDA scores for individual models, i.e.,

$$\tau_{\text{ens}}(x, \mathcal{S}; \{f_{\Theta^{(i)}}\}_{i=1}^{I}) = \frac{1}{I} \sum_{i=1}^{I} \tau(x, \mathcal{S}; f_{\Theta^{(i)}}).$$
(1)

They could also come from more sophisticated aggregation over the individual models. For example, the TRAK method (Park et al., 2023) works as the following:

$$\tau_{\mathrm{TRAK}}(x,\mathcal{S};\{f_{\Theta^{(i)}}\}_{i=1}^{I}) = \left(\frac{1}{I}\sum_{i=1}^{I}\mathbf{Q}_{f_{\Theta^{(i)}}}\right) \left(\frac{1}{I}\sum_{i=1}^{I}\phi_{f_{\Theta^{(i)}}}\left(\Phi_{f_{\Theta^{(i)}}}^{\top}\Phi_{f_{\Theta^{(i)}}}\right)^{-1}\Phi_{f_{\Theta^{(i)}}}^{\top}\right), \quad (2)$$

where $\Theta^{(i)}$ are parameters of models independently trained on the training data; $\mathbf{Q}_{f_{\Theta^{(i)}}}$ is a diagonal matrix with each diagonal element corresponding to the "one minus correct-class probability" of a training data point under model $\Theta^{(i)}$; $\phi_{f_{\Theta^{(i)}}}$ is the vector gradient of $f_{\Theta^{(i)}}(x)$ with respect to $\Theta^{(i)}$; and $\Phi_{f_{\Theta^{(i)}}}$ is the matrix of the vector gradients of $f_{\Theta^{(i)}}(x_j)$ stacked over the training samples $x_j \in S$.¹

¹Some details of the exact TRAK implementation are omitted here.



Figure 1: DROPOUT ENSEMBLE consists of two steps: (1) train I models ($\{\Theta^{(i)}\}_{i=1}^{I}$) independently; (2) get D dropout-masked models ($\{f_{\Theta^{(i)}}^{(d)}\}_{d=1}^{D}$) for each i = 1, ..., I.

178 179

188

189

190

191

192

193

199

162 163

164

166

167

169 170 171

172

173 174

175 176

177

3.2 Computational costs of TDA methods

While recent research has shown that ensembling gradient-based TDA methods can achieve decent attribution efficacy with tens of independently trained models (as opposed to hundreds or even thousands of model retraining in retraining-based TDA methods) (Park et al., 2023), it is still impractical for very large-scale or edge-device applications due to the increased time and space costs. To better quantify the time and space costs associated with TDA methods, we categorize the costs into three parts: training time cost, serving time cost, and storage cost. We provide a detailed explanation for these costs as follows:

- **Training time cost**: This refers to the computational time incurred by processing and training models to be used by TDA ensembles. This is the major computational challenge for ensembling TDA methods.
- Serving time cost: This refers to the remaining computational time to deploy TDA ensembles *after model training*. Typically, ensembling gradient-based TDA methods derives different TDA scores using multiple trained models and then aggregates them. The serving time cost will include all computational costs incurred during this process.
- Space cost: This refers to the parameters of the trained models that are stored by TDA ensembles. These model parameters are required to run forward and backward passes on top of the models for TDA score calculation. The space costs are generally proportional to the model size.

3.3 DROPOUT ENSEMBLE

General methodology. Dropout, a common module used in many modern deep learning models, is
 initially designed as an efficient approximation of ensemble (Srivastava et al., 2014). Motivated by
 this idea, we propose a simple, efficient ensemble strategy, DROPOUT ENSEMBLE, for TDA methods.
 Instead of performing the TDA ensemble over a large number of independently trained models, the
 proposed method utilizes multiple dropout masks on the same model to perform the TDA ensemble.

In practice, DROPOUT ENSEMBLE consists of two steps as illustrated in Figure 1. In the first step, we train *I* independent models, $\{\Theta^{(i)}\}_{i=1}^{I}$, as shown by the green arrows in Figure 1. Next, as shown by the blue arrows in Figure 1, for each model $\Theta^{(i)}$, $i \in \{1, ..., I\}$, we obtain *D* variants of the model with different dropout masks, which are denoted as $\{f_{\Theta^{(i)}}^{(d)}\}_{d=1}^{D}$. For any TDA method, DROPOUT ENSEMBLE calculates the ensembled TDA scores through $\tau_{ens}(x, S; \{f_{\Theta^{(i)}}^{(d)}\}_{1 \le i \le I, 1 \le d \le D})$.

In comparison to the naive independent ensemble method, DROPOUT ENSEMBLE can significantly
 reduce the training time and space costs. As we will show in Section 4.2, this method allows
 us to replace the expensive independently trained models with dropout-masked models that incur
 no additional training cost or model parameters. There is some serving time cost overhead caused
 by DROPOUT ENSEMBLEAs we will show in Appendix M, the overhead is small and DROPOUT
 ENSEMBLE could achieve better efficacy-efficiency trade-off.



Figure 2: LORA ENSEMBLE consists of two steps: (1) train I models ($\{\Theta^{(i)}\}_{i=1}^{I}$) independently; (2) get L LORA fine-tuned models ($\{\Theta^{(i,l)}_{LORA}\}_{l=1}^{L}$) for each i = 1, ..., I.

TDA-method-specific optimization. It is possible to further optimize the proposed ensemble strategy when applying it to a particular TDA method. For example, we develop a variant of DROPOUT ENSEMBLE tailored for the TRAK method to reduce the serving time cost of DROPOUT ENSEMBLE. Recall that in Eq. (2), the quantities Q's only involve the forward pass of the models on the training data, while ϕ 's and Φ 's require the backward pass on the training data. We find that, perhaps a bit surprisingly, if we only use the dropout-masked models to calculate Q's while using the original models to calculate ϕ 's and Φ 's, we can get similar attribution efficacy. Concretely, directly applying DROPOUT ENSEMBLE on TRAK leads to

$$\tau_{\mathrm{ens}}\left(x, \mathcal{S}; \{f_{\Theta^{(i)}}^{(d)}\}_{\substack{1 \le i \le I \\ 1 \le d \le D}}\right) = \left(\frac{1}{I \cdot D} \sum_{i=1}^{I} \sum_{d=1}^{D} \mathbf{Q}_{f_{\Theta^{(i)}}^{(d)}}\right) \left(\frac{1}{I \cdot D} \sum_{i=1}^{I} \sum_{d=1}^{D} \phi_{f_{\Theta^{(i)}}^{(d)}} \left(\Phi_{f_{\Theta^{(i)}}^{\top}}^{\top} \Phi_{f_{\Theta^{(i)}}^{(d)}}\right)^{-1} \Phi_{f_{\Theta^{(i)}}^{\top}}^{\top}\right)$$

while in the optimized version, we have

$$\tau_{\mathrm{ens}}\left(x, \mathcal{S}; \{f_{\Theta^{(i)}}^{(d)}\}_{\substack{1 \leq i \leq I\\ 1 \leq d \leq D}}, \{f_{\Theta^{(i)}}\}_{1 \leq i \leq I}\right) = \left(\frac{1}{I \cdot D} \sum_{i=1}^{I} \sum_{d=1}^{D} \mathbf{Q}_{f_{\Theta^{(i)}}^{(d)}}\right) \left(\frac{1}{I} \sum_{i=1}^{I} \phi_{f_{\Theta^{(i)}}} \left(\Phi_{f_{\Theta^{(i)}}}^{\top} \Phi_{f_{\Theta^{(i)}}}\right)^{-1} \Phi_{f_{\Theta^{(i)}}}^{\top}\right).$$

In this case, we only need to evaluate the forward pass on the dropout-masked models and reduce the serving time cost by avoiding the backward pass. We call this variant of our method as DROPOUT ENSEMBLE (forward-only).

3.4 LORA ENSEMBLE

For large-scale generative models, especially Transformer models (Vaswani et al., 2017), LoRA adapters have shown great success for efficiently fine-tuning the models (Hu et al., 2021). Our second efficient ensemble strategy, LORA ENSEMBLE, is motivated by the LoRA techniques. In particular, we propose to replace the independently trained models in the naive ensemble method with LoRA fine-tuned models.

Similar to DROPOUT ENSEMBLE, LORA ENSEMBLE also consists of two steps as illustrated in Figure 2. The first step trains I independent models, $\{\Theta^{(i)}\}_{i=1}^{I}$ (green arrows in Figure 2), while the second step further trains L LoRA fine-tuned models, $\{\Theta^{(i,l)}_{LORA}\}_{l=1}^{L}$, for each $\Theta^{(i)}$, $i \in \{1, ..., I\}$ (brown arrows in Figure 2).

In comparison to DROPOUT ENSEMBLE, LORA ENSEMBLE will introduce a small amount of
 additional training time cost and model parameters for LoRA fine-tuning. However, since we only
 use the parameters in LoRA adapters to compute the TDA scores, LORA ENSEMBLE can reduce
 the serving time cost compared to either dropout-masked models in DROPOUT ENSEMBLE or
 independent models in the naive method. This leads to a unique advantage for LORA ENSEMBLE if
 the serving time cost is critical among the computational costs.



Figure 3: The LDS of naive independent ensemble and DROPOUT ENSEMBLE with different numbers of dropout-masked passes (D) and independently trained models (I). We apply the ensemble methods to the TDA method, TRAK. There are four experiment settings: MLP classifiers trained on MNIST and CIFAR-2 (top row); ResNet9 trained on CIFAR-2 (bottom-left); and Music Transformer trained on MAESTRO (bottom-right). The x-axis indicates the training time cost measured by the number of independently trained models (I). The y-axis indicates the attribution efficacy measured by LDS.

4 EXPERIMENTS

In this section, we empirically evaluate the efficiency and efficacy of the proposed DROPOUT ENSEMBLE and LORA ENSEMBLE.

4.1 EXPERIMENTAL SETUP

We conduct extensive experiments on a wide range of settings.

TDA methods. We apply the proposed methods to various gradient-based TDA methods, including
 influence function (IF) (Koh & Liang, 2017), Grad-Dot/Grad-Cos (Charpiat et al., 2019), and
 TRAK (Park et al., 2023). The detailed descriptions and implementations of these algorithms are
 provided in Appendix B.

Datasets and models. We consider models with different architectures trained on diverse datasets:
(1) a three-layer MLP classifier trained on the MNIST-10 dataset (LeCun et al., 1998), (2) a three-layer
MLP classifier and a ResNet-9 classifier (He et al., 2016) trained on the CIFAR-2 dataset (a two-class
subset of the CIFAR-10 dataset (Krizhevsky et al., 2009)), and (3) a Music Transformer (Huang
et al., 2018) trained on the MAESTRO dataset (Hawthorne et al., 2018). Notably, the former two
are supervised classification settings, while the last one is a generative modeling setting. We sample
5000 training samples and 500 test samples from MNIST-10 and CIFAR-2 datasets. For MAESTRO
dataset, we sample 5000 training samples and 178 generated samples. More detailed setups for



Figure 4: The LDS of naive independent ensemble and DROPOUT ENSEMBLE on more TDA methods, IF, Grad-Dot, and Grad-Cos. The experiments are performed on MLP classifiers trained on MNIST. The plot setup is similar as Figure 3.

each dataset and model are listed in Appendix C. MNIST-10 dataset holds CC BY-SA 3.0 license.
CIFAR-10 dataset holds CC-BY 4.0 license. MAESTRO dataset holds CC BY-NC-SA 4.0 license.

Evaluation metric for TDA efficacy. We utilize the linear datamodeling score (LDS) (Park et al., 2023) to evaluate the efficacy of the TDA methods augmented by different ensembling methods. Intuitively, LDS measures the rank correlation between the TDA scores among training samples and the change of model outputs by removing a subset of training samples and retraining the model from scratch. A higher LDS value corresponds to a better alignment between the TDA scores and the influence of the training samples on the model outputs, thus better TDA quality. We refer the reader to Appendix E for the exact definition of LDS.

348 Evaluation metrics for TDA efficiency. As introduced in Section 3.2, we measure the TDA 349 efficiency in terms of the training time cost, serving time cost, and space cost. For both training 350 and serving time costs, we measure the wall-clock time on a single A40 GPU. For the space cost, 351 although ideally one would measure it by memory usage, this can be challenging because memory usage is sensitive to the specific implementation and parallelization. Therefore, we measure the space 352 cost by the total parameter count instead. For DROPOUT ENSEMBLE, we will also use the number of 353 independently trained models (I) as a measure of the training time cost and space cost, as both of the 354 two costs are proportional to I in this case. More details about the measurements are provided in 355 Appendix F. 356

357 358

359

335

336

337

338

4.2 DROPOUT ENSEMBLE

360 **Improvement over the training time cost.** To illustrate the improvement over training time cost 361 comparing DROPOUT ENSEMBLE to the naive ensemble, we first report the results on the TRAK 362 method across four experiment settings. As can be seen in Figure 3, across all four experiment 363 settings, increasing the number of dropout-masked passes (D) results in significant LDS improvement 364 for a fixed number of independently trained models (I). Recall that DROPOUT ENSEMBLE does not incur any additional training time cost for a fixed I. This implies that DROPOUT ENSEMBLE 366 can significantly reduce the training time cost for achieving the same level of attribution efficacy 367 as measured by LDS. For example, DROPOUT ENSEMBLE with I = 5 can reach higher LDS than 368 naive independent ensemble with I = 25 for MLP on MNIST and MLP/ResNet9 on CIFAR-2, which achieves a 80% reduction on training time cost. For Music Transformer on MAESTRO, DROPOUT 369 ENSEMBLE with I = 1 can reach higher LDS than naive independent ensemble with I = 10, which 370 achieves a 90% reduction. 371

We find that DROPOUT ENSEMBLE works similarly well for other TDA methods. In Figure 4, we report the results on IF, Grad-Dot, and Grad-Cos. We only experiment on the setting of MLP classifiers trained on MNIST, as these TDA methods do not have meaningfully good efficacy on more complex settings. Similar to the TRAK experiments, we observe that DROPOUT ENSEMBLE with I = 1 could match the LDS of naive independent ensemble with I = 10 for IF and Grad-Dot, and match that with I = 5 for Grad-Cos. Therefore, we expect that the proposed DROPOUT ENSEMBLE can be generalized to various gradient-based TDA methods.



In this section, we present a theoretical analysis to compare TDA method with the naive ensemble (regular independent ensemble) and our two-step ensemble (a general case containing DROPOUT ENSEMBLE and LORA ENSEMBLE). With notation defined in Section 3.1, for any test sample $x \in \mathcal{T}$ and S as the training set, we assume a TDA method τ finds an "optimal" attribution score $\tau(\Theta^*)$



Figure 6: The LDS of naive ensemble and LORA ENSEMBLE with respect to different cost measurements. Here, we apply the ensemble methods to TRAK on Music Transformer trained on the
MAESTRO dataset. The x-axis of Figure (6a, 6b, 6c) indicates training/serving time costs (running
time on a single A40), while that of Figure 6d specifies the space cost (total parameter count). The
y-axis of all figures is the attribution efficacy measured by LDS. LORA ENSEMBLE has significantly
fewer costs in all aspects than naive ensemble for achieving similar LDS.

based on the optimal parameter Θ^* that some ensemble estimators try to approximate (omit f_{Θ} here for notational convenience). Then, we define the two types of ensemble estimators for $\tau(\Theta^*)$ and through Lemma 5.1 below to show that DROPOUT ENSEMBLE and LORA ENSEMBLE outperform the regular ensemble in terms of the approximation error.

469 We start with the definition of the two types of ensemble estimators. For a TDA method τ , the 470 regular ensemble estimator τ_{ens} is defined as the average of $\tau(\Theta^{(i)})$, where $\Theta^{(i)}$ for $i = 1, \dots, I$ 471 are I identically distributed (i.d.) individual estimators. The two-step ensemble estimator $\tau_{2-\text{step}}$ is defined as the average of $\tau(\Theta^{(k,d)})$, where each of $\Theta^{(k)}$ for $k = 1, \dots, K$ is an individual estimator 472 (just as $\Theta^{(i)}$ in regular ensemble estimator) and by using $\Theta^{(k)}$ as a base estimator, D variants of 473 it, $\Theta^{(k,d)}$, are generated for the second-step ensemble. This is a general definition and includes 474 475 DROPOUT ENSEMBLE and LORA ENSEMBLE as special cases. Our goal is to compare the squared 476 error of τ_{ens} and τ_{2-step} with respect to $\tau(\Theta^*)$, which we formalize as the following lemma (proof in Appendix J.1). 477

Lemma 5.1. Define τ_{ens} and τ_{2-step} as the regular and two-step ensemble estimators for $\tau(\Theta^*)$:

485

$$\tau_{ens} = \frac{1}{I} \sum_{i=1}^{I} \tau(\Theta^{(i)}) \quad and \quad \tau_{2\text{-step}} = \frac{1}{KD} \sum_{k=1}^{K} \sum_{d=1}^{D} \tau(\Theta^{(k,d)})$$

Assume each individual estimate $\tau(\Theta^{(i)})$ and $\tau(\Theta^{(k,d)})$ have the same distribution. Then, the difference in squared error $\Delta = \|\tau_{ens} - \tau(\Theta^*)\|^2 - \|\tau_{2-step} - \tau(\Theta^*)\|^2$ is given by:

$$\Delta = \frac{1}{I} \left(\Sigma_{1,1} + (I-1)\Sigma_{i,j} \right) - \frac{1}{KD} \left(\Sigma_{(k,m),(k,m)} + (D-1)\Sigma_{(k,m),(k,n)} + D(K-1)\Sigma_{(k,m),(l,n)} \right)$$

486 *where:*

488

489

496

500

501

- $\Sigma_{i,j} = Cov(\tau(\Theta^{(i)}), \tau(\Theta^{(j)}))$ is the covariance between individual estimators in the regular ensemble. Given the i.d. assumption, $\Sigma_{i,j}$ are the same for all $i \neq j$, and $\Sigma_{i,j} = \Sigma_{1,1}$ for i = j.
- $\Sigma_{(k,m),(l,n)} = Cov(\tau(\Theta^{(k,m)}), \tau(\Theta^{(l,n)}))$ is the covariance between variants of base individual estimators in the two-step ensemble. We further define $\Sigma_{(k,m),(k,n)}$ as the within-group covariance and $\Sigma_{(k,m),(l,n)}$ with $k \neq l$ as the between-group covariance.

⁴⁹³ Now, we analyze Δ and show that Δ will stay positive under reasonable assumptions, which implies ⁴⁹⁴ that $\tau_{2\text{-step}}$ outperforms τ_{ens} . We especially derive Δ for two interesting cases: K = I and KD = I⁴⁹⁵ (see assumptions and derivation in Appendix J.2).

497 **Case 1:** K = I means the number of individual estimators K in the two-step ensemble is equal 498 to the number of individual estimators I in the regular ensemble, but the two-step ensemble has D499 variants for each base estimator, resulting in a total of $K \cdot D = I \cdot D$ estimates. In this case, we drive

$$\Delta = \frac{D-1}{ID} \left(\Sigma_{1,1} - \Sigma_{(k,m),(k,n)} \right). \tag{3}$$

As long as the within-group covariance $\Sigma_{(k,m),(k,n)}$ is smaller than variance of each individual 502 estimator $\Sigma_{1,1}$, $\tau_{2-\text{step}}$ will outperform the regular ensemble τ_{ens} . This is likely to be the case for any 503 pair of different variants of the same base estimator through dropout or LoRA fine-tuning. For a fixed 504 small I, when D is small, $\frac{D-1}{D}$ has bigger impact on Δ and increase D can quickly increase Δ means 505 $\tau_{2\text{-step}}$ outperforms τ_{ens} more. In contrast, when D is large, $\frac{D-1}{D}$ is close to 1 and increase D will not change Δ much, meaning the performance gain of $\tau_{2\text{-step}}$ over τ_{ens} will saturate, which matches 506 507 our empirical results in Figure 3 and Figure 6. On the other hand, for a large I, changing D will 508 not change Δ much as Δ is already small, and the performance gain of $\tau_{2-\text{step}}$ over τ_{ens} will be less 509 significant, which also aligns with our empirical observation in Figure 3 and Figure 6. In summary, 510 when K = I, $\tau_{2\text{-step}}$ clearly outperforms τ_{ens} when the within-group covariance is small benefiting 511 from averaging over more total estimators (ID versus I), and the performance gain saturates as the 512 within-group covariance increases.

Case 2: KD = I means the total number of estimators in the two-step ensemble $K \cdot D$ is equal to the number of estimators I in the regular ensemble. Thus, both ensembles have the same number of estimates, but the two-step ensemble introduces a base-variant relationship. In this case, we derive:

517 518

529

513

$$\mathbf{A} = \frac{D-1}{I} \left(\Sigma_{i,j} - \Sigma_{(k,m),(k,n)} \right). \tag{4}$$

This case is about comparing the within-group covariance $\Sigma_{(k,m),(k,n)}$ and the individual estimator 519 covariance $\Sigma_{i,j}$. In general, $\Sigma_{(k,m),(k,n)}$ is expected to be larger than $\Sigma_{i,j}$, because $\Sigma_{(k,m),(k,n)}$ is the 520 covariance between variants of the same base estimator, while $\Sigma_{i,j}$ is the covariance between different 521 estimators. This can lead to negative Δ , indicating that the two-step ensemble is outperformed by the 522 regular ensemble in this case. However, since KD = I, $\frac{D-1}{I}$ is close to $\frac{1}{K}$, meaning $\tau_{2-\text{step}}$ is not 523 significantly outperformed by τ_{ens} when a large K is set. In summary, when KD = I, τ_{2-step} can still 524 outperform τ_{ens} if the within-group covariance is small. However, if this covariance is large, then 525 $\tau_{2-\text{step}}$'s advantage diminishes, and τ_{ens} may perform similarly or better due to the independent nature of its estimators. Even for the latter, two-step ensemble still enjoy drastic efficiency gain as generating 527 more variants for each base estimator is much cheaper, e.g., via dropout or LoRA fine-tuning, than 528 generating more individual estimators from re-training.

530 6 CONCLUSION

531 We present DROPOUT ENSEMBLE and LORA ENSEMBLE as efficient alternatives to the naive 532 independent ensemble approach for improving gradient-based TDA methods. The proposed strategies 533 significantly reduce training time (up to 80%), serving time (up to 60%), and space cost (up to 80%), 534 while maintaining similar attribution efficacy in comparison to the naive ensemble. Empirical results from our extensive experiments show that the proposed efficient ensembles can remarkably advance 536 the Pareto frontier of TDA methods with better computational efficiency and TDA efficacy. Notably, we demonstrate the proposed methods work well on a generative modeling setting. Our theoretical analysis shows our ensemble method has a smaller squared error for estimating the optimal scores 538 compared to naive ensembles. In the future, we will utilize our methods in more real-world generative settings that were blocked by high computational costs or low effectiveness of TDA methods.

540 REFERENCES

542 543 544	Juhan Bae, Nathan Ng, Alston Lo, Marzyeh Ghassemi, and Roger B Grosse. If influence functions are the answer, then what is the question? <i>Advances in Neural Information Processing Systems</i> , 35:17953–17967, 2022.
545 546 547	Elnaz Barshan, Marc-Etienne Brunet, and Gintare Karolina Dziugaite. Relatif: Identifying explanatory training samples via relative influence. In <i>International Conference on Artificial Intelligence and Statistics</i> , pp. 1899–1909, PMLR, 2020.
548 549 550	Samyadeep Basu, Philip Pope, and Soheil Feizi. Influence functions in deep learning are fragile. <i>arXiv preprint arXiv:2006.14651</i> , 2020.
551 552	Guillaume Charpiat, Nicolas Girard, Loris Felardos, and Yuliya Tarabalka. Input similarity from the neural network perspective. <i>Advances in Neural Information Processing Systems</i> , 32, 2019.
553 554 555	Junwei Deng and Jiaqi Ma. Computational copyright: Towards a royalty model for ai music generation platforms. <i>arXiv preprint arXiv:2312.06646</i> , 2023.
556 557	Logan Engstrom, Axel Feldmann, and Aleksander Madry. Dsdm: Model-aware dataset selection with datamodels. <i>arXiv preprint arXiv:2401.12926</i> , 2024.
558 559 560 561	Vitaly Feldman and Chiyuan Zhang. What neural networks memorize and why: Discovering the long tail via influence estimation. <i>Advances in Neural Information Processing Systems</i> , 33:2881–2891, 2020.
562 563	Amirata Ghorbani and James Zou. Data shapley: Equitable valuation of data for machine learning. In <i>International conference on machine learning</i> , pp. 2242–2251. PMLR, 2019.
564 565 566 567	Han Guo, Nazneen Rajani, Peter Hase, Mohit Bansal, and Caiming Xiong. FastIF: Scalable influence functions for efficient model interpretation and debugging. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> . Association for Computational Linguistics, 2021. doi: 10.18653/v1/2021.emnlp-main.808.
569 570	Zayd Hammoudeh and Daniel Lowd. Training data influence analysis and estimation: A survey. <i>Machine Learning</i> , pp. 1–53, 2024.
571 572 573	Curtis Hawthorne, Andriy Stasyuk, Adam Roberts, Ian Simon, Cheng-Zhi Anna Huang, Sander Dieleman, Erich Elsen, Jesse Engel, and Douglas Eck. Enabling factorized piano music modeling and generation with the maestro dataset. <i>arXiv preprint arXiv:1810.12247</i> , 2018.
575 576 577	Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pp. 770–778, 2016.
578 579 580	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. <i>arXiv preprint arXiv:2106.09685</i> , 2021.
581 582 583 584	Cheng-Zhi Anna Huang, Ashish Vaswani, Jakob Uszkoreit, Noam Shazeer, Curtis Hawthorne, Andrew M Dai, Matthew D Hoffman, and Douglas Eck. Music transformer: Generating music with long-term structure. <i>arXiv preprint arXiv:1809.04281</i> , 2018.
585 586	Andrew Ilyas, Sung Min Park, Logan Engstrom, Guillaume Leclerc, and Aleksander Madry. Data- models: Predicting predictions from training data. <i>arXiv preprint arXiv:2202.00622</i> , 2022.
587 588 589 590 591	Ruoxi Jia, David Dao, Boxin Wang, Frances Ann Hubis, Nick Hynes, Nezihe Merve Gürel, Bo Li, Ce Zhang, Dawn Song, and Costas J Spanos. Towards efficient data valuation based on the shapley value. In <i>The 22nd International Conference on Artificial Intelligence and Statistics</i> , pp. 1167–1176. PMLR, 2019.
592 593	Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models.

arXiv preprint arXiv:2001.08361, 2020.

594 595 596	Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In <i>International conference on machine learning</i> , pp. 1885–1894. PMLR, 2017.
597	Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
598 599 600	Yongchan Kwon and James Zou. Beta shapley: a unified and noise-reduced data valuation framework for machine learning. In <i>International Conference on Artificial Intelligence and Statistics</i> , pp. 8780–8802. PMLR, 2022.
602 603	Yongchan Kwon, Eric Wu, Kevin Wu, and James Zou. Datainf: Efficiently estimating data influence in lora-tuned llms and diffusion models. <i>arXiv preprint arXiv:2310.00902</i> , 2023.
604 605	Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. <i>Proceedings of the IEEE</i> , 86(11):2278–2324, 1998.
607 608	Sung Min Park, Kristian Georgiev, Andrew Ilyas, Guillaume Leclerc, and Aleksander Madry. Trak: Attributing model behavior at scale. <i>arXiv preprint arXiv:2303.14186</i> , 2023.
609 610 611	Garima Pruthi, Frederick Liu, Satyen Kale, and Mukund Sundararajan. Estimating training data influence by tracing gradient descent. <i>Advances in Neural Information Processing Systems</i> , 33: 19920–19930, 2020.
612 613 614	Andrea Schioppa, Polina Zablotskaia, David Vilar, and Artem Sokolov. Scaling up influence functions. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 36, pp. 8179–8186, 2022.
615 616	Anders Søgaard et al. Revisiting methods for finding influential examples. <i>arXiv preprint</i> arXiv:2111.04683, 2021.
618 619	C. Spearman. The proof and measurement of association between two things. <i>The American Journal of Psychology</i> , 15(1):72–101, 1904.
620 621 622	Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. <i>The journal of machine learning research</i> , 15(1):1929–1958, 2014.
623 624	John Tukey. Bias and confidence in not quite large samples. Ann. Math. Statist., 29:614, 1958.
625 626 627	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. <i>Advances in neural information processing systems</i> , 30, 2017.
628 629 630 631	Jiachen T Wang and Ruoxi Jia. Data banzhaf: A robust data valuation framework for machine learning. In <i>International Conference on Artificial Intelligence and Statistics</i> , pp. 6388–6421. PMLR, 2023.
632 633 634	Chih-Kuan Yeh, Joon Kim, Ian En-Hsu Yen, and Pradeep K Ravikumar. Representer point selection for explaining deep neural networks. <i>Advances in neural information processing systems</i> , 31, 2018.
635 636 637	
638 639	
640 641	
642 643	
644 645	
646 647	

648 A MORE RELATED WORK

650 **Retraining-based TDA methods.** Retraining-based methods compute TDA scores by systematically 651 retraining the model with and without specific training samples to quantify their influence on the 652 model's output. These methods are computationally expensive due to the requirement of the large 653 amount of model retrainings. For example, Leave-One-Out (LOO) influence Tukey (1958) measures the prediction difference between the model trained on the full training dataset and models trained 654 on subsets with only one specific sample dropped. Data Shapley (Ghorbani & Zou, 2019; Jia et al., 655 2019), Beta-Shapley (Kwon & Zou, 2022), and Data Banzhaf (Wang & Jia, 2023) extends the LOO 656 idea to consider data interactions for more equitable TDA scores, but they require retraining models 657 on all possible subsets of the training dataset. Similarly, DataModels (Ilyas et al., 2022) tries to learn 658 the model predictions when the model is trained on each subset of the training dataset, which requires 659 a nontrivial amount of model retraining. While, in practice, sampling or approximation will be used 660 to reduce the number of model retrainings, these methods typically still require thousands of or more 661 retrainings to achieve satisfactory attribution efficacy. In summary, the high demand for retraining 662 makes these TDA methods computationally inefficient and limits their practical applicability to even 663 moderately large models.

664 Gradient-based TDA methods. Gradient-based methods are another group of TDA methods 665 that usually provide closed-form TDA scores using gradients. Since the seminal work of influence 666 function by Koh & Liang (2017), gradient-based methods have become increasingly popular due 667 to their scalability. The influence function (Koh & Liang, 2017) and subsequent studies (Guo 668 et al., 2021; Barshan et al., 2020; Schioppa et al., 2022; Kwon et al., 2023) obtain TDA scores 669 by approximating the effect of upweighting a training sample on the loss function. Moreover, 670 Representer Point Selection (Yeh et al., 2018) decomposes the pre-activation of a neural network 671 as a linear combination of training samples. TracIn (Pruthi et al., 2020) traces the loss changes on the test points during the training process. TRAK (Park et al., 2023) uses the neural tangent kernel 672 with random projection to assess influence. These gradient-based methods significantly reduced the 673 computational cost compared to retraining-based methods. The downside is that they typically rely 674 on the convexity assumption and Taylor approximation to calculate TDA scores. These requirements 675 lead to performance degradation on non-convex neural networks and sensitivity to the randomness 676 inherent in model initialization and training. 677

B TDA METHODS

678

679 680

681

682

683 684

685

686 687

688 689

690

691

692 693

694

695 696 697 In this section, we provide details on the TDA methods we perform efficient ensembling on. Additionally, we also introduce the ensembling aggregation methods, i.e., the way to aggregate $\tau_{\text{ens}}(x, S; \{f_{\Theta^{(i)}}\}_{i=1}^{I})$, for each method. The notation will follow Section 3.1.

Tracing with the Randomly-projected After Kernel (TRAK). This is a state-of-the-art gradientbased TDA method provided by Park et al. (2023). It natively introduces ensembling in its definition.

$$\tau_{\mathrm{TRAK}}(x,\mathcal{S};\{f_{\Theta^{(i)}}\}_{i=1}^{I}) = \left(\frac{1}{I}\sum_{i=1}^{I}\mathbf{Q}_{f_{\Theta^{(i)}}}\right) \left(\frac{1}{I}\sum_{i=1}^{I}\phi_{f_{\Theta^{(i)}}}\left(\Phi_{f_{\Theta^{(i)}}}^{\top}\Phi_{f_{\Theta^{(i)}}}\right)^{-1}\Phi_{f_{\Theta^{(i)}}}^{\top}\right)$$

Detailed notation description is described in Section 3.1. We use 2048 as the random projection dimension for TRAK. For each independently trained model, we train the model on half sampled training set following Park et al. (2023).

Influence functions based on the conjugate gradients (IF). First proposed by Koh & Liang (2017), the definition of IF is

$$\tau_{\mathrm{IF}}(x,\mathcal{S};\{f_{\Theta^{(i)}}\}_{i=1}^{I}) = \left[\frac{1}{I}\sum_{i=1}^{I}g_{f_{\Theta^{(i)}}}(x_{j})^{\top}H_{f_{\Theta^{(i)}}}^{-1}g_{f_{\Theta^{(i)}}}(x):x_{j}\in\mathcal{S}\right],$$

where $g_{f_{\Theta}(i)}(x)$ is the vector gradient of model output $f(x; \Theta^{(i)})$ with respect to the parameters $\Theta^{(i)}$, $H_{f_{\Theta}(i)}^{-1}$ is the inverse hessian matrix with respect to the training set, and $g_{f_{\Theta}(i)}(x_j)^{\top}$ is the vector gradient to the training sample. The product of the first two terms are an inverse-hessian-vectorproduct problem, we implement conjugate gradients approach to solve it. **Grad-Dot.** Grad-Dot is proposed by Charpiat et al. (2019). We simply calculate Grad-Dot multiple times on trained parameters $\Theta^{(i)}, i \in \{1, \dots, K\}$ and take the average value.

$$\tau_{\text{Grad-Dot}}(x, \mathcal{S}; \{f_{\Theta^{(i)}}\}_{i=1}^{I}) = \left[\frac{1}{I}\sum_{i=1}^{I}g_{f_{\Theta^{(i)}}}(x)^{\top}g_{f_{\Theta^{(i)}}}(x_j) : x_j \in \mathcal{S}\right],$$

where $g_{f_{\Theta(i)}}(x)^{\top}$ is the vector gradient of model output $f(x; \Theta^{(i)})$ with respect to the parameters $\Theta^{(i)}$, and $g_{f_{\Theta(i)}}(x_j)^{\top}$ is the gradient to the training sample.

Grad-Cos. Similar to Grad-Dot,

$$\tau_{\operatorname{Grad-Cos}}(x, \mathcal{S}; \{f_{\Theta^{(i)}}\}_{i=1}^{I}) = \left[\frac{1}{I} \sum_{i=1}^{I} \frac{g_{f_{\Theta^{(i)}}}(x)^{\top}}{\|g_{f_{\Theta^{(i)}}}(x)\|} \frac{g_{f_{\Theta^{(i)}}}(x_j)}{\|g_{f_{\Theta^{(i)}}}(x_j)\|} : x_j \in \mathcal{S}\right],$$

where $g_{f_{\Theta}(i)}(x)^{\top}$ is the vector gradient of model output $f(x; \Theta^{(i)})$ with respect to the parameters $\Theta^{(i)}$, and $g_{f_{\Theta}(i)}(x_j)^{\top}$ is the gradient to the training sample.

720 721 722

723 724

719

710

711 712

C DETAILED EXPERIMENT SETUP

MLP on MNIST-10. For the MNIST-10 (LeCun et al., 1998) experiment, we sampled 5000 training
 samples and 500 testing samples. We used a 3-layer MLP with hidden layer sizes equal to 128 and 64
 and placed dropout layers after the first two linear layers with a rate of 0.1, which resulted in a total
 of about 0.11M parameters. We employed an SGD optimizer with learning rate 0.01, momentum 0.9,
 and batch size 64 to train this MLP classifier for 100 epochs.

730

731 MLP/ResNet-9 on CIFAR-2. For CIFAR-2 experiment, we construct the CIFAR-2 dataset by 732 sampling from the CIFAR-10 dataset that only include the "cat" and "dog" classes (same as the 733 setting in Park et al. (2023)). To incorporate a diverse set of model architectures, we train an MLP 734 and a CNN-based model on this dataset. We also only consider a subset of the CIFAR-2 dataset with a training size equal to 5000 and a testing size equal to 500. Firstly, we use a 3-layer MLP 735 with hidden layer sizes equal to 120 and 84 and place dropout layers after the first two linear layers 736 with a rate of 0.1, which results in a total of about 0.38M parameters. We employ an SGD optimizer 737 with a learning rate of 0.01, momentum of 0.9, and batch size 64 to train this MLP classifier for 50 738 epochs. Additionally, we consider a standard ResNet-9 model (He et al., 2016) with dropout layers 739 being placed after all convolution layers, which has roughly 4.83M trainable parameters. We train 740 this model for 50 epochs.

741 742

MusicTransformer on MAESTRO. For the MAESTRO experiment, we use the MIDI and Audio 743 Edited for Synchronous TRacks and Organization (MAESTRO) dataset (v2.0.0) Hawthorne et al. 744 (2018) and construct a Music Transformer following the original setting in (Huang et al., 2018). 745 Specifically, the number of layers equals to 6, the number of independent heads equal to 8, the input 746 feature size is 512 and the dimension of the feedforward network is 1024. For data processing, we 747 follow the basic experiment setup used by (Deng & Ma, 2023). To be more specific, we define 748 a vocabulary set of size equal to 388, which includes "NOTE ON" and "NOTE OFF" events for 749 128 different pitches, 100 "TIME SHIFT" events, and 32 "VELOCITY" events. The raw data is 750 pre-processed as sequences of about 90K events. Due to computational constraints, we train a Music 751 Transformer model on a subset of the official training set in the MAESTRO dataset with a size of 752 5000. For training, the batch size has been set to 64, and the model is trained by a classic seq2seq loss function. We employ an Adam optimizer with a learning rate equal to 1e-4, $\beta_1 = 0.9$ and $\beta_2 = 0.98$. 753 We apply zero warm-up steps since the dataset size is comparably small, and we train the model 754 for 20 epochs. For music event generation, we use 178 samples from the official testing dataset as 755 prompts to generate music with a single event, which is used to evaluate the TDA methods.

⁷⁵⁶ D EFFICIENT ENSEMBLE SETUP

758 D.1 DROPOUT ENSEMBLE

760 Dropout is applied to the same layers as the model training stated in Appendix C. Dropout rate is set to 0.1 for all experiment settings.

763 D.2 LORA ENSEMBLE

This method is only applied to MusicTransformer trained on the MAESTRO dataset. According to 765 Section 3.4, we first train I independent models on the full training dataset with 10 epochs using 766 different model initialization, and all the other remaining settings follow Appendix C. The LoRA 767 adapters used in this experiment all have rank r = 8, alpha $\alpha = 8$, and trainable biases. No dropouts 768 are applied for LoRA parameters. We augment LoRA adapters to only the W_q and W_v matrices in 769 the self-attention modules within all the layers of the MusicTransformer. Then, we define fine-tuning 770 datasets for each LoRA adapter as random training data subsets with a size equal to 2500 (i.e., half of 771 the original training dataset size). We further fine-tune each of the aforementioned I independent 772 models using L different LoRA adapters for an additional 10 epochs on these random training data 773 subsets, which results in a total of $I \cdot L$ LoRA fine-tuned models.

774 775

776

783 784 785

762

764

E LINEAR DATAMODELING SCORE (LDS)

Park et al. (2023) proposed the *linear datamodeling score* (LDS), aiming at probing the TDA method's ability to make counterfactual predictions based on the attribution score derived from the learned model output function f_{Θ} and the corresponding dataset to train f_{Θ} . Because most TDA methods are assumed to be *additive*², the TDA scores can be used to predict the model output function learned from a subset of training data in a summation form. Formally, the *attribution-based output predictions* of the model output function $f_{\Theta,S'}$ is defined as follows:

$$g_{\tau}(x, \mathcal{S}'; \mathcal{S}) \triangleq \sum_{i: x_i \in \mathcal{S}'} \tau(x, \mathcal{S}; f_{\Theta})_i,$$
(5)

where S is the training set, $S' \subseteq S$ is a subset of S and $f_{\Theta_{S'}}$ is the model output function with $\Theta_{S'}$ learned from S'. Intuitively, $g_{\tau}(x, S'; S)$ computes the overall attribution of the subset S' on example x, which should be a powerful indicator of the model prediction on x (i.e., $f_{\Theta_{S'}}(x)$) if the TDA method works well. The *linear datamodeling score* (LDS) is defined to measure the predictive power of $g_{\tau}(x, S'; S)$ and can be formalized as follows:

Definition E.1 (Linear datamodeling score). *Given a training set* S, *a model output function* f_{Θ} , *and a corresponding TDA method* τ . *Let* $\{S_1, \ldots, S_m : S_j \subseteq S\}$ *be m randomly sampled subsets of* S, *each of size* $\alpha \times n$ *for some fixed* $\alpha \in (0, 1)$. *The linear datamodeling score (LDS) of* τ *for a specific example* x *is defined as*

$$LDS(\tau, x) \triangleq \boldsymbol{\rho}(\{f_{\Theta_{\mathcal{S}_i}}(x) : j \in [m]\}, \{g_{\tau}(x, \mathcal{S}_j; \mathcal{S}) : j \in [m]\}),$$

where ρ is the Spearman rank correlation (Spearman, 1904), $f_{\Theta_{S_j}}$ is the model output function with Θ_{S_j} learned from S_j and $g_{\tau}(x, S_j; S)$ is defined in Eq (5).

To compute LDS for our experiment settings, we use 50 models that are independently trained on random subsets with size half of the full dataset (i.e., we set m = 50 and $\alpha = 0.5$ in Definition E.1).

800 801 802

808

809

796

797

798

799

F WALL-CLOCK TIME MEASUREMENTS

In this paper, we use the wall-clock time on a single A40 GPU to measure the computational costs if the number of independently trained models (i.e., the value of I) can not precisely demonstrate the costs, e.g., the training time cost and serving time cost of LORA ENSEMBLE.

Here, we define several components that dominate the wall-clock time of different ensemble methods.

 $^{^{2}}$ If a TDA method is additive, then it defines an attribution score that the overall influence of a group is the sum of the individual influence in the group.

- 810 • T_{Train}: The time to train a model from scratch. 811
 - T_{Train, Base}: The time to train a base model from scratch for LoRA tuning. Normally speaking, T_{Train, Base} < T_{Train}.
 - T_{Train, LoRA}: The time to fine-tune for one LoRA adapter.
- 814 • T_{Serving}: The time to calculate the TDA scores for one trained model *after model training*.
- T_{Serving, Forward-only}: The time to calculate the TDA scores for one trained models after model training 815 with shared gradients. (will only be used by DROPOUT ENSEMBLE(forward only). Normally 816 speaking, $T_{Serving, Forward-only} < T_{Serving, Forward-only}$. 817
- T_{Serving, LoRA}: The time to calculate the TDA scores after model training for one LoRA adapter. 818 Normally speaking, $T_{Serving, LoRA} < T_{Serving}$. 819
- The total computational cost is approximated and summarized in t: 820
 - Training time cost (naive independent ensemble/DROPOUT ENSEMBLE/DROPOUT ENSEMBLE (forward only)): $I \times T_{\text{Train}}$.
 - Training time cost (LORA ENSEMBLE): $I \times T_{\text{Train, Base}} + I \times L \times T_{\text{Train, LoRA}}$.
 - Serving time cost (naive independent ensemble): $I \times T_{\text{Serving}}$.
 - Serving time cost (DROPOUT ENSEMBLE): $I \times D \times T_{Serving}$.
- 826 • Serving time cost (DROPOUT ENSEMBLE (forward only)): $I \times T_{Serving} + I \times (D-1) \times T_{Serving}$ T_{Serving, Forward-only}.
 - Serving time cost (LORA ENSEMBLE): $I \times L \times T_{\text{Serving, LORA}}$.

G ADDITIONAL EXPERIMENT FOR DROPOUT ENSEMBLE

G.1 ABLATION EXPERIMENT TO RANDOM PROJECTION

Here, we examine the root of the DROPOUT ENSEMBLE's performance through an ablation experi-834 ment. There are no random factors other than dropout for IF, Grad-Dot, or Grad-Cos, while there is 835 another random factor, i.e., random projection, in TRAK. We perform D dropout-masked passes with 836 dropout enabled, i.e., "DROPOUT ENSEMBLE' and D dropout-masked passes with dropout disabled, 837 i.e., "Only Random Projection", in Table 1 and calculate the LDS. Random projection does contribute 838 to the accuracy improvement, but the improvement will saturate when D is large. 839

		1	3	5
Naive Independent Ensemble	N/A	0.122	0.217	0.275
DROPOUT ENSEMBLE Only Random Projection	10	0.249 0.210	0.399 0.335	0.457 0.398
DROPOUT ENSEMBLE Only Random Projection	25	0.316 0.217	0.458 0.351	0.502 0.413

Table 1: Ablating the contribution of test-time dropout. The experiment is carried out on MNIST+MLP. Note that D is the number of dropout-masked passes, and I is the number of independently trained models.

852 853 854

855

856

858

850

851

812

813

821

822

823

824

825

827

828 829 830

831 832

833

G.2 INTERMEDIATE CHECKPOINTS

Here, we show that DROPOUT ENSEMBLE ensemble can also be applied to intermediate checkpoints and improve the accuracy in Table 2. Park et al. (2023) states that intermediate checkpoints with fewer epochs can be used for ensembling to reduce the training time cost. We save multiple checkpoints from different epochs of the same training process (I = 1).

- 859 860
- 861 862

MEMORY COSTS OF DROPOUT ENSEMBLE AND LORA ENSEMBLE Η

Here we record the peak memory usage at serving time for DROPOUT ENSEMBLE (Table 3) and 863 LORA ENSEMBLE (Table 4) applied on TRAK algorithm. The memory of vanilla DROPOUT

Ensembling methods	<i>#ckpts</i>	1	3	5	10
Naive Independent Ensemble	N/A	0.122	0.170	0.188	0.201
DROPOUT ENSEMBLE	10	0.249	0.318	0.342	0.362
	25	0.316	0.359	0.373	0.389

Table 2: The TDA efficacy of DROPOUT ENSEMBLE on intermediate checkpoints from different epochs of the same training process (I = 1). Note that D is the number of dropout-masked passes, and #ckpts is the number of intermediate checkpoints used.

877 ENSEMBLE is the same as naive independent ensemble. The meomory of DROPOUT ENSEMBLE
878 (Forward Only) is larger than vanilla DROPOUT ENSEMBLE because some of the cached terms. The
879 memory usage of LORA ENSEMBLE is slightly lower than naive independent ensemble because of
880 the reduction in parameter size.

Datasets and Models	D method variants	3 10 25
MNIST+MLP	DROPOUT ENSEMBLE DROPOUT ENSEMBLE (Forward Only)	342M 636M 1956M 4787M
CIFAR2+MLP	DROPOUT ENSEMBLE DROPOUT ENSEMBLE (Forward Only)	344M 638M 1966M 4797M
CIFAR2+ResNet9	DROPOUT ENSEMBLE DROPOUT ENSEMBLE (Forward Only)	477M 785M 2087M 4877M
MAESTRO+MusicTransformer	DROPOUT ENSEMBLE DROPOUT ENSEMBLE (Forward Only)	538M 634M 1529M 3421M

Table 3: The peak memory usage of DROPOUT ENSEMBLE and its alternative on TRAK with different numbers of dropout-masked passes (D) and the number of independently trained models fixed to 5 (I = 5).

Ensembling Methods		1	3	5	10
Naive Independent Ensemble	N/A	431M	538M	538M	538M
LORA ENSEMBLE LORA ENSEMBLE LORA ENSEMBLE	3 10 25	404M 404M 404M	404M 404M 404M	404M 404M 404M	404M 404M 404M

Table 4: The peak memory usage of LORA ENSEMBLE and naive independent ensemble applied on TRAK for MusicTransformer trained on MAESTRO dataset. Note that *I* represents the number of independently trained models, and *L* is the number of LoRA adapters augmented on each model.

I SPACE COSTS OF DROPOUT ENSEMBLE AND LORA ENSEMBLE

917 Here we record the parameter count to present the space cost for DROPOUT ENSEMBLE (Table 5) and LORA ENSEMBLE (Table 6).

settings I(w/ any D)	MNIST +MLP	CIFAR-2 +MLP	CIFAR-2 +ResNet-9	MAESTRO +MusicTransformer
1	0.11M	0.38M	4.83M	13.11M
3	0.33M	1.14M	14.48M	39.35M
5	0.55M	1.89M	24.13M	65.58M
10	1.09M	3.79M	48.25M	131.16M
25	2.73M	9.48M	120.63M	327.89M

Table 5: The space cost (total parameter count) of different experiment settings under different numbers of independently trained models (I). Note that DROPOUT ENSEMBLE will **not** incur any additional storage cost with more dropout-masked passes (i.e., the space cost is fixed with respect to D).

	Naive Independent ensemble $(L = 0)$	3	10	25
1	13.11M 30.35M	13.41M	14.09M	15.57M
5	65.58M	40.25M 67.05M	42.29M 70.49M	77.87M
10	131.76M	134.11M	140.99M	155.73M

939Table 6: The space cost (total parameter count) of MusicTransformer under different numbers of940independently trained models (I) and different numbers of LoRA adapters (L). LORA ENSEMBLE941only add marginal space cost (at most a 18.7% increment for D = 25 across all I) compared to naive942independent ensembling.

J PROOFS

J.1 PROOF OF LEMMA 5.1

Proof. We proceed by analyzing the bias and variance contributions for both the naive ensemble and the two-step ensemble. We first analyze the bias of the two estimators. Let $\mu_i = \mathbb{E}[\tau(\Theta^{(i)})]$ and $\mu_{(k,d)} = \mathbb{E}[\tau(\Theta^{(k,d)})]$ represent the expected values of the individual and variant estimators, respectively. Since we assume that each individual estimate $\tau(\Theta^{(i)})$ and $\tau(\Theta^{(k,d)})$ have the same and identical distribution, we have $\mu_i = \mu_{(k,d)}$, and thus

$$\operatorname{Bias}^{2}(\tau_{\operatorname{ens}}) = \left(\frac{1}{I}\sum_{i=1}^{I}\mu_{i} - \tau(\Theta^{*})\right)^{2} = \left(\frac{1}{KD}\sum_{k=1}^{K}\sum_{d=1}^{D}\mu_{(k,d)} - \tau(\Theta^{*})\right)^{2} = \operatorname{Bias}^{2}(\tau_{2}\operatorname{-step}) \quad (6)$$

Therefore, the difference in squared error Δ arises solely from the variance terms.

We then analyze the variance of the two estimators. For the naive ensemble estimator τ_{ens} , the variance is given by:

$$\operatorname{Var}(\tau_{\operatorname{ens}}) = \frac{1}{I^2} \sum_{i=1}^{I} \operatorname{Var}(\tau(\Theta^{(i)})) + \frac{2}{I^2} \sum_{i < j} \operatorname{Cov}(\tau(\Theta^{(i)}), \tau(\Theta^{(j)}))$$
(7)

968 Using the simplified notation, where $\Sigma_{1,1} = \text{Var}(\tau(\Theta^{(i)}))$ and $\Sigma_{i,j} = \text{Cov}(\tau(\Theta^{(i)}), \tau(\Theta^{(j)}))$ for 969 $i \neq j$, this can be written as:

$$\operatorname{Var}(\tau_{\operatorname{ens}}) = \frac{1}{I} \left(\Sigma_{1,1} + (I-1)\Sigma_{i,j} \right)$$
(8)

For the two-step ensemble estimator $\tau_{2-\text{step}}$, the variance is given by:

$$\operatorname{Var}(\tau_{2\operatorname{-step}}) = \frac{1}{(KD)^2} \sum_{k=1}^{K} \sum_{d=1}^{D} \operatorname{Var}(\tau(\Theta^{(k,d)})) + \frac{2}{(KD)^2} \sum_{(k,d) < (k',d')} \operatorname{Cov}(\tau(\Theta^{(k,d)}), \tau(\Theta^{(k',d')}))$$
(9)

This variance can be decomposed into within-group and between-group covariances. Let $\Sigma_{(k,m),(k,n)} = \operatorname{Cov}(\tau(\Theta^{(k,m)}), \tau(\Theta^{(k,n)}))$ represent the within-group covariance (i.e., between variants of the same base estimator), and let $\Sigma_{(k,m),(l,n)} = \text{Cov}(\tau(\Theta^{(k,m)}), \tau(\Theta^{(l,n)}))$ represents the between-group covariance (i.e., between variants of different base estimators). The variance simplifies to:

$$\operatorname{Var}(\tau_{2-\operatorname{step}}) = \frac{1}{KD} \left(\Sigma_{(k,m),(k,m)} + (D-1) \Sigma_{(k,m),(k,n)} + D(K-1) \Sigma_{(k,m),(l,n)} \right)$$
(10)

Finally, the difference in variance between the two estimators, denoted as Δ , is:

$$\Delta = \operatorname{Var}(\tau_{\text{ens}}) - \operatorname{Var}(\tau_{2\text{-step}}) \tag{11}$$

Substituting the variance expressions for τ_{ens} and τ_{2-step} , we have:

$$\Delta = \frac{1}{I} \left(\Sigma_{1,1} + (I-1)\Sigma_{i,j} \right) - \frac{1}{KD} \left(\Sigma_{(k,m),(k,m)} + (D-1)\Sigma_{(k,m),(k,n)} + D(K-1)\Sigma_{(k,m),(l,n)} \right)$$
(12)

This expression shows that the difference in error depends on the number of estimators (I in the naive ensemble and KD in the two-step ensemble) and the covariance structure among the estimators. The two-step ensemble will outperform the naive ensemble if the within-group and between-group covariance are sufficiently small compared to the overall covariance in the naive ensemble.

J.2 DERIVATION OF THE ERROR DIFFERENCE FOR K = I and KD = I

For both derivations, we use the following assumptions:

$$\Sigma_{1,1} = \Sigma_{(k,m),(k,m)}$$
 and $\Sigma_{i,j} = \Sigma_{(k,m),(l,n)}$ for $i \neq j, k \neq l$ (13)

This first equality implies that the variance of each estimator in both ensemble approaches is the same. The second equality implies that the covariance between different estimators in the naive ensemble is the same as the covariance between different estimator variants in the two-step ensemble. These are reasonable assumptions following the i.d. assumption on the individual estimators.

Case 1: K = I. Substitute K = I into the expression for Δ and use the assumption $\Sigma_{1,1} =$ $\Sigma_{(k,m),(k,m)}$ and $\Sigma_{i,j} = \Sigma_{(k,m),(l,n)}$, this simplifies Δ to:

$$\Delta = \frac{1}{I} \left(\Sigma_{1,1} + (I-1)\Sigma_{i,j} \right) - \frac{1}{ID} \left(\Sigma_{1,1} + (D-1)\Sigma_{(k,m),(k,n)} + D(I-1)\Sigma_{i,j} \right)$$
(14)

$$=\frac{1}{I}\Sigma_{1,1} + \frac{I-1}{I}\Sigma_{i,j} - \frac{1}{ID}\Sigma_{1,1} - \frac{D-1}{ID}\Sigma_{(k,m),(k,n)} - \frac{I-1}{I}\Sigma_{i,j}$$
(15)

$$= \frac{1}{2} \sum_{1,1} - \frac{1}{2} \sum_{1,1} - \frac{D-1}{2} \sum_{(k,m)} \frac{(k,m)}{(k,m)}$$
(16)

$$= \frac{1}{I} \sum_{i,1} - \frac{1}{ID} \sum_{i,1} - \frac{1}{ID} \sum_{(k,m),(k,n)}$$
(16)
$$D - 1 \qquad D - 1$$

1023
1024
$$= \frac{D-1}{ID} \Sigma_{1,1} - \frac{D-1}{ID} \Sigma_{(k,m),(k,n)}$$
1025
$$D-1 = (-1) \sum_{k=1}^{n} \sum_{k=1}^$$

$$= \frac{D-1}{ID} \left(\Sigma_{1,1} - \Sigma_{(k,m),(k,n)} \right)$$
(18)

Case 2: KD = I. Substitute KD = I into the expression for Δ and use the assumption $\Sigma_{1,1} = \Sigma_{(k,m),(k,m)}$ and $\Sigma_{i,j} = \Sigma_{(k,m),(l,n)}$, this simplifies Δ to:

$$\Delta = \frac{1}{I} \left(\Sigma_{1,1} + (I-1)\Sigma_{i,j} \right) - \frac{1}{I} \left(\Sigma_{1,1} + (D-1)\Sigma_{(k,m),(k,n)} + (I-D)\Sigma_{i,j} \right)$$
(19)

$$= \frac{1}{I} \Sigma_{1,1} + \frac{I-1}{I} \Sigma_{i,j} - \frac{1}{I} \Sigma_{1,1} - \frac{D-1}{I} \Sigma_{(k,m),(k,n)} - \frac{I-D}{I} \Sigma_{i,j}$$
(20)

$$= \frac{I-1}{I} \Sigma_{i,j} - \frac{D-1}{I} \Sigma_{(k,m),(k,n)} - \frac{I-D}{I} \Sigma_{i,j}$$
(21)

$$= \frac{D-1}{I} \sum_{i,j} - \frac{D-1}{I} \sum_{(k,m),(k,n)}$$
(22)

$$= \frac{D-1}{I} \left(\Sigma_{i,j} - \Sigma_{(k,m),(k,n)} \right)$$
(23)

1040 K ADDITIONAL EXPERIMENTS

1039

1071

Here we report the results of additional experiments, including a new experiment setting, GPT on
 the Shakespeare dataset for language modeling task, and two existing settings with a larger scale,
 MNIST with 60K training images and MAESTRO with 15K training music sequences. As can be
 seen in Figure 7 and Figure 8 (respectively for DROPOUT ENSEMBLE and LORA ENSEMBLE), both
 methods still show clear improvements in these new experiments.



Figure 8: LDS of Naive Ensemble and LORA ENSEMBLE. Please refer to Figure 6 for the plot format. LORA ENSEMBLE is applied to Transformer models only therefore the ResNet + MNIST setting is omitted. One point in the MT on MAESTRO (15K) plot is missing as we did not finish the experiment on time.

1072 L COMPUTATION OVERHEAD

In Table 7, we list the computational overhead by different ensemble methods measured on Music-Transformer with TRAK. The overhead of naive ensemble is proportional to the number of models used in the ensemble. In comparison, the DROPOUT ENSEMBLE does not incur any additional cost in terms of training time cost (no additional training) and space cost (the dropout-masked variants do not need to be stored), while incurring additional serving time cost. The LORA ENSEMBLE incurs additional costs in terms of all three types of costs, but they are much smaller than naive ensemble. Furthermore, the serving time cost is significantly reduced in comparison to DROPOUT ENSEMBLE since we only need to calculate the gradients with respect to the LoRA parameters.

1080	Cost Type	No Ensemble	Naive (3)	Naive (10)	D=3	D=10	L=3	L=10
1082	Training Time	1	3	10	1	1	2.00	5.50
1083	Serving Time	1	3	10	3	10	1.85	6.07
1084	Space	1	3	10	1	1	1.03	1.10

Table 7: The relative computation overhead for Naive Ensemble (varying number of ensembles),
DROPOUT ENSEMBLE (one base model and varying D), and LORA ENSEMBLE (one base model and varying L). The "No Ensemble" column refers to the cost with only one model, where the entries are normalized to one for easier comparison. The relative costs of Naive Ensemble and DROPOUT ENSEMBLE are based on simple counting while the relative costs of LORA ENSEMBLE are based on experiments on MusicTransformer with TRAK.

1092 1093

1094

M COMPUTATIONAL COST OF EACH EFFICIENT ENSEMBLE METHODS

Here we report the computational cost of each efficient ensemble method. The results of the ratio of training over serving cost (i.e. training time cost serving time cost) shown in Table 8 demonstrate that training cost dominates the overall computational cost of TDA method.

We also record the training and serving time of each ensemble methods and plot them in Figure 9.
 The result shows that DROPOUT ENSEMBLE (and the forward-only variant) and LORA ENSEMBLE
 can achieve better efficacy-efficiency trade-off. In other words, both ensembles could improve the
 efficacy with the same computational cost or reduce the computational cost to reach the same efficacy.

 $40 \times$

 $60 \times$

 $25 \times$

 $4.2 \times$

Table 8: The training cost / serving cost ratio of each experiment settings.

Training/Serving | Training/Serving(forward-only)

 $200 \times$

 $450 \times$

 $208 \times$

 $42 \times$

- 1103
- 1104 1105

1106

- 1107
- 1108
- 1109 1110
- 1111
- 1112

1114

1113 N QUALITATIVE RESULTS

Here we provide a qualitative results to show how DROPOUT ENSEMBLE improve the TDA result
without training more models. We select two random test images from MNIST-10 and the corresponding training samples that are most helpful (with highest score) and most detracting (with
lowest score). The result calculated by DROPOUT ENSEMBLE (D=50) is more reasonable by manual
inspection. We can observe that the training images identified by DROPOUT ENSEMBLE (D=50)
have "shapes" similar to the corresponding test images while helpful (detracting) examples are of the
same (different) class as the test images themselves.

1122 1123

1124

O NOISY LABEL DETECTION RESULTS

Experiment Setting

MNIST-10 + MLP

CIFAR-2 + MLP

CIFAR-2 + ResNet9

MAESTRO + MusicTransformer

The improvement of efficient ensembles could also be demonstrated in downstream tasks of TDA, such as noisy label detection. Here we present the AUC of noisy label detection on MNIST-10 with 20% randomly flipped training samples. The performance of TRAK with DROPOUT ENSEMBLE outperforms the naive independent ensemble without training more models.

```
1130 P DISCUSSION OF DROPOUT ENSEMBLE(FORWARD-ONLY)
```

- 1131
- DROPOUT ENSEMBLE(forward-only) performs well in the experiments demonstrated in Figure 5 and Figure 9, especially considering the competitive efficacy-efficiency trade-off. A potential explanation to this phenomenon is that the term (Q) affects the performance "substantially" as shown





in TRAK (Park et al., 2023)'s ablation study. Focusing on mitigating the effect of randomness on
this term and cache other terms could improve the TDA performance effectively and efficiently. It is
worth noting that this phenomenon may be closely related to some open questions existing in the
TRAK algorithm. For instance, empirical evidence indicates that a particular approach to ensembling
models, known as "term-wise" ensemble in the TRAK paper, plays a critical role in TRAK's superior
performance. However, the underlying reasons for this remain insufficiently understood.

1188 1189	Ensembling methods	#ckpts D	1	5
1190 1191	Naive Independent Ensemble	N/A	0.778	0.810
1192		10	0.802	0.849
1193	DROPOUT ENSEMBLE	25	0.810	0.849
1194				

Table 9: The AUC (Area-Under-Curve) of noisy label detection on MNIST-10 with 20% randomly
flipped training samples. TRAK with efficient ensembles can outperform naive TRAK under all
settings with drastically fewer computational cost.

1200 Q SERVING COST OVERHEAD

Here we also record the serving time of each ensemble methods and plot them in Figure 11. The results shows that DROPOUT ENSEMBLE(forward-only) and LORA ENSEMBLE could still outperform the naive ensemble if only serving cost is considered. DROPOUT ENSEMBLE has a larger serving time cost overhead. It is worth noting that the training cost and the serving cost are both for the TDA process, which are defined respectively in Section 3.2. The computational cost of TDA should consider both costs together.



