# ADVERSARIAL ATTACKS ON DATA ATTRIBUTION

Anonymous authors

Paper under double-blind review

# ABSTRACT

Data attribution aims to quantify the contribution of individual training data points to the outputs of an AI model, which has been used to measure the value of training data and compensate data providers. Given the impact on financial decisions and compensation mechanisms, a critical question arises concerning the adversarial robustness of data attribution methods. However, there has been little to no systematic research addressing this issue. In this work, we aim to bridge this gap by detailing a threat model with clear assumptions about the adversary's goal and capabilities and proposing principled adversarial attack methods on data attribution. We present two methods, *Shadow Attack* and *Outlier Attack*, which generate manipulated datasets to inflate the compensation adversarially. The Shadow Attack leverages knowledge about the data distribution in the AI applications, and derives adversarial perturbations through "shadow training", a technique commonly used in membership inference attacks. In contrast, the Outlier Attack does not assume any knowledge about the data distribution and relies solely on black-box queries to the target model's predictions. It exploits an inductive bias present in many data attribution methods-outlier data points are more likely to be influentialand employs adversarial examples to generate manipulated datasets. Empirically, in image classification and text generation tasks, the Shadow Attack can inflate the data-attribution-based compensation by at least 200%, while the Outlier Attack achieves compensation inflation ranging from 185% to as much as 643%.

# 1 INTRODUCTION

031

000

001 002 003

004

006 007

008 009

010

011

012

013

014

015

016

017

018

019

021

024

025

026

027 028 029

Data attribution aims to quantify the contribution of individual training data points to the outputs
 of an Artificial Intelligence (AI) model (Koh & Liang, 2017). A key application of data attribution
 is to measure the value of training data in AI systems, enabling appropriate compensation for data
 providers (Ghorbani & Zou, 2019; Jia et al., 2019). With the rapid advancement of generative AI,
 these methods have gained increased relevance, particularly in addressing copyright concerns. Re cent studies have explored economic frameworks using data attribution for copyright compensation,
 showing promising preliminary results (Deng & Ma, 2023; Wang et al., 2024).

Given the significant potential of data attribution methods for data valuation and compensation, an important question arises regarding their adversarial robustness. As these methods influence financial decisions and compensation mechanisms, they may attract malicious actors seeking to manipulate the system for personal gain. This underscores the need to investigate whether data attribution methods can be manipulated and exploited, as their vulnerabilities could lead to unfair compensation, undermining the trustworthiness of the solutions built on top of these methods.

However, adversarial attacks on data attribution methods have received little to no exploration in
 prior literature. Along with the proposal of a data-attribution-based economic solution for copyright
 compensation, Deng & Ma (2023) briefly experimented with a few heuristic approaches (e.g., duplicating data samples) for attacking data attribution methods. A systematic study that clearly defines
 the threat model and develops principled adversarial attack methods has yet to be conducted.

This work presents the first comprehensive study to fill this gap. We first outline the threat model
by detailing the data compensation workflow and specifying the assumptions we made. One key assumption is that the data contribution is periodic, and there is certain persistence across consecutive
iterations of data contributions, which is the source of knowledge that the adversary could exploit.
We also assume that the adversary may either have access to the distribution of the data used by the

target model of the AI system or can get black-box queries of the target model's predictions, both are commonly seen in the AI security literature (Shokri et al., 2017; Chen et al., 2017).

Subsequently, we propose two adversarial attack strategies, Shadow Attack and Outlier Attack, re-057 spectively relying on different assumptions about the adversary's capabilities. The Shadow Attack relies on the access to data distribution and employs the "shadow training" technique commonly used in membership inference attacks (Shokri et al., 2017) to train "shadow models" that imitate the 060 target model. The adversary can then directly perturb their dataset to achieve a higher compensation 061 on these shadow models. The Outlier Attack, instead, does not assume knowledge about the data dis-062 tribution but only relies on black-box queries of the target model's predictions. The key idea behind 063 this method lies in an inductive bias of many data attribution methods—outlier data points are more 064 likely to be more influential. The proposed Outlier Attack utilizes adversarial examples (Goodfellow et al., 2015; Chen et al., 2017) to generate realistic outliers in a black-box fashion. 065

We conduct extensive experiments, including both image classification and text generation settings, to demonstrate the effectiveness of the proposed attack methods. Our results show that by only adding imperceptible perturbations to real-world data features, the Shadow Attack can inflate the adversary's compensation to at least 200% and up to 456%, while the Outlier Attack can inflate the adversary's compensation to at least 185% and up to 643%.

Overall, our study reveals a critical practical challenge—adversarial vulnerability—in deploying data attribution methods for data valuation and compensation. Moreover, the design of the proposed attack methods, especially the Outlier Attack that exploits a common inductive bias of data attribution methods, offers deeper insights into these vulnerabilities. These findings provide valuable directions for future research to enhance the robustness of data attribution methods.

076 077

078

# 2 RELATED WORK

079 Data Attribution for Data Valuation and Compensation. Data attribution methods have been widely used for quantifying the value of training data in AI applications and compensating data 081 providers (Ghorbani & Zou, 2019; Jia et al., 2019; Yoon et al., 2020; Kwon & Zou, 2022; Feldman & Zhang, 2020; Xu et al., 2021; Lin et al., 2022; Kwon & Zou, 2023; Just et al., 2023; Wang 083 & Jia, 2023; Deng & Ma, 2023; Wang et al., 2024). With the rapid advancement of generative 084 AI, these methods have gained increasing relevance due to the growing concerns around copyright. 085 Recent studies have proposed economic solutions using data attribution for copyright compensation, yielding promising preliminary results (Deng & Ma, 2023; Wang et al., 2024). Given the significant potential of data attribution methods for data valuation and compensation, a critical question arises 087 regarding their adversarial robustness. Aside from one earlier exploration using heuristic approaches 880 to attack data attribution methods (Deng & Ma, 2023), no systematic study has addressed this issue. 089 This work presents the first comprehensive study that outlines a detailed threat model and proposes 090 principled and effective approaches for adversarial attacks on data attribution methods. 091

**Membership Inference Attack.** Membership inference attacks aim to infer whether a specific 092 data point was used during the training of a machine learning model, typically without knowing the actual training dataset or having white-box access to the model (Shokri et al., 2017). The general 094 strategy involves leveraging information such as model architecture, training data distribution, or 095 black-box model predictions (Shokri et al., 2017; Song & Mittal, 2021). We refer the readers to Hu 096 et al. (2022) for a detailed survey on this topic. One of the proposed attack methods, the Shadow 097 Attack, is inspired by the "shadow training" technique (Shokri et al., 2017) commonly used in 098 membership inference attacks, where the adversary draws "shadow samples" following the same 099 distribution as the actual training dataset used by the target model to be attacked, and trains "shadow models" based on the shadow samples to imitate the target model. 100

Adversarial Example. Adversarial example (Goodfellow et al., 2015) is a well-known phenomenon where small, often imperceptible perturbations to input features can significantly alter the predictions of machine learning models, particularly deep neural networks. These perturbations can be generated through black-box queries to the model predictions (Chen et al., 2017; Ilyas et al., 2018; Guo et al., 2019). For a comprehensive review of adversarial examples, see Yuan et al. (2019); Chakraborty et al. (2021). In the proposed Outlier Attack, we employ black-box adversarial attack methods for generating adversarial examples to generate realistic outliers relative to the training dataset of the target model, without needing access to the training dataset or the model details.

# <sup>108</sup> 3 THE THREAT MODEL

# 110 3.1 THE DATA COMPENSATION SCENARIO

We consider a scenario where there is an AI Developer, and a (potentially large) set of Data Providers. The Data Providers supply training data for the AI Developer to develop an AI model. In return, the Data Providers are compensated based on their data's contribution to the model, as measured by a specific data attribution method.

Periodic Data Contribution. We assume that the Data Providers contribute data to the AI 116 Developer *periodically*, a common practice in many AI applications. For example, large lan-117 guage models need periodic updates to stay aligned with the latest factual knowledge about the 118 world (Zhang et al., 2023); recommender systems must adapt to evolving user preferences (Zhang 119 et al., 2020); quantitative trading firms rely on up-to-date information to power their predictive 120 models<sup>1</sup>; and generative models for music or art benefit from fresh, innovative works by artists to 121 diversify their creative outputs (Smith et al., 2024). However, such periodic data contribution intro-122 duces risks: a malicious Data Provider (referred to as an Adversary thereafter) could exploit 123 information from previous iterations to adversarially manipulate their future data contribution, po-124 tentially inflating their compensation unfairly.

125 To formalize this scenario, without loss of generality, we consider two consecutive iterations of data 126 contribution, denoted as time steps t = 0 and t = 1. At t = 0, there is no Adversary and 127 the training dataset consists solely of contributions from benign Data Providers. This dataset 128 is represented as  $Z_0 \in \mathcal{Z}$ , where  $\mathcal{Z}$  is the set of all possible training datasets. At t = 1, the training dataset is  $Z_1 = Z_1^b \cup Z_1^a$ , where  $Z_1^b \in \mathcal{Z}$  represents the training data provided by benign 129 Data Providers, while  $Z_1^a \in \mathcal{Z}$  is the set provided by the Adversary. We also assume that 130 131  $Z_1^b \cap Z_1^a = \emptyset$ , meaning there is no overlap between the two datasets. Finally, for a dataset  $Z \in \mathcal{Z}$ , each data point  $z \in Z$  is represented as a pair z = (x, y), where  $x \in \mathcal{X}$  is the input feature and  $y \in \mathcal{Y}$ 132 is the prediction target, with  $\mathcal{X}$  and  $\mathcal{Y}$  referring to the feature space and target space, respectively. 133

AI Training and Data Attribution. Let  $\mathcal{M}$  represent the set of AI models, and let  $\mathcal{T}: \mathcal{Z} \to \mathcal{M}$ denote a *training algorithm*, mapping a training dataset  $Z \in \mathcal{Z}$  to a model  $\mathcal{T}(Z) \in \mathcal{M}$ .

136 In data attribution, we aim to understand how individual data points from a training dataset  $Z \in \mathbb{Z}$ 137 contribute to the model output on a target (validation) data point from a validation dataset V. Given 138 Z and V, a *data attribution method* derives a *contribution function*  $\tau: Z \times V \to \mathbb{R}$  that assigns a 139 real value  $\tau(z, v)$  to each training data point  $z \in Z$  for a given validation data point  $v \in V$ . Denote 140 the set of all such function  $\tau$ 's as C, and the set of all possible validation sets as V. Therefore, a data 141 attribution method can be formalized as a function  $A: \mathbb{Z} \times M \times V \to C$ . In most cases, we will 142 consider  $\mathcal{A}(Z, \mathcal{T}(Z), V)$ , hence Z and V alone suffice to specify the resulting  $\tau$ .

**Compensation Mechanism.** In practice, the contribution function  $\tau$  derived from data attribution methods may not reliably measure the contributions of all training data points due to the inherent randomness in AI model training and the need for efficiency (Wang & Jia, 2023; Nguyen et al., 2024). Specifically, measurements of data points with smaller contributions are often less reliable than those of the most influential contributors. Following Deng & Ma (2023), we consider a compensation mechanism where only the top-k influential training data points for each validation data point  $v \in V$  receive a fixed amount of compensation.

150 3.2 THE ADVERSARY151

156 157

161

We now discuss the Adversary's objective and capability under the data compensation scenario. The Objective of the Adversary. Let  $V_1$  be the validation dataset used at step t = 1. The objective of the Adversary is to construct a dataset  $Z_1^a$  that maximizes the *compensation share* received by the Adversary, which is defined as

$$c(Z_1^a) = \frac{1}{k|V_1|} \sum_{z \in Z_1^a} \sum_{v \in V_1} \mathbf{1}[\tau_1(z, v) \in \operatorname{Top}_k\left(\{\tau_1(z', v) \mid z' \in Z_1\}\right)],\tag{1}$$

where  $\tau_1 = \mathcal{A}(Z_1, \mathcal{T}(Z_1), V_1)$  is the contribution function at t = 1; Top<sub>k</sub>(·) extracts the top-k elements from a finite set of real numbers; and  $\mathbf{1}[\cdot]$  is the indicator function.

<sup>&</sup>lt;sup>1</sup>See, for example, the Bloomberg market data feed: https://www.bloomberg.com/ professional/products/data/enterprise-catalog/market/.

The Capabilities of the Adversary. We first outline the limitations imposed on the Adversary.
 In realistic scenarios, the Adversary does not have access to any of the following:

- the exact training datasets ( $Z_0$  and  $Z_1^b$ ) and the exact validation datasets ( $V_0$  and  $V_1$ );
- white-box access to the trained models  $\mathcal{T}(Z_0)$  and  $\mathcal{T}(Z_1)$ ;
- the contribution functions  $\mathcal{A}(Z_0, \mathcal{T}(Z_0), V_0)$  and  $\mathcal{A}(Z_1, \mathcal{T}(Z_1), V_1)$ .

Then we make the following assumptions to characterize the capabilities of the Adversary.

**Assumption 1** (Persistence). Assume that  $Z_0 \subseteq Z_1^b$  and  $|Z_0|/|Z_1^b|$  is close to 1. Additionally, assume that  $V_0$  and  $V_1$  are independently sampled from the same distribution.

**Assumption 2** (Access to data distribution and training algorithm). Assume that the Adversary has access to the distributions of  $Z_0$  and  $V_0$ . Additionally, assume that the Adversary has knowledge about the training algorithm  $\mathcal{T}$  (but not the model  $\mathcal{T}(Z_0)$ ).

Assumption 3 (Black-box access to model). Assume that the Adversary has access to black-box query to the model  $\mathcal{T}(Z_0)$  to query the model predictions on any input feature  $x \in \mathcal{X}$ .

Intuitively, Assumption 1 assumes persistence between time steps t = 0, 1, so that information gained from t = 0 can inform the attack at t = 1. In Assumption 2 and Assumption 3, it is worth noting that the Adversary's access to information is limited to what is available at t = 0 only.

In practice, Assumption 1 is realistic in many real-world applications. For example, in applications that have periodic model updates, such as large language models or recommender systems, a substantial portion of the training data often remains consistent across consecutive iterations, with only incremental changes. Moreover, Assumption 2 is a common assumption in the membership inference attack literature (Shokri et al., 2017), while Assumption 3 reflects a common setup for generating adversarial examples against neural network models (Chakraborty et al., 2021).

The two proposed attack methods in Section 4 and Section 5 respectively rely on Assumption 2 and
Assumption 3, but do not depend on both simultaneously. As a result, the first method is a gray-box
attack method while the second method is a black-box attack method.

The Action Space of the Adversary. We assume that the Adversary is restricted to making only small perturbations to an existing set of real data points to construct the adversarial dataset  $Z_1^a$ . This implies that the Adversary cannot introduce entirely synthetic or arbitrary data, but rather, can modify real data points subtly. This reflects a realistic adversarial scenario as overly large or unnatural alterations would be easily detectable.

196 197

165 166

167

168

178

# 4 Shadow Attack

198 199 200

201

In this section, we introduce the proposed *Shadow Attack*, which leverages Assumption 1 and Assumption 2, and exploits the knowledge about the data distribution at t = 0 to perform attacks.

At a high level, the Adversary first performs a *shadow training* process, where models are trained on data drawn from a distribution similar to that of the training dataset  $Z_0$ , allowing the Adversary to approximate the target model  $\mathcal{T}(Z_0)$ . The Adversary then applies adversarial perturbations to the data points they plan to contribute to the AI Developer at t = 1. The adversarial perturbations are derived by maximizing the data attribution values on the models obtained through shadow training. See Figure 1 for an illustration.

208 209

4.1 SHADOW TRAINING

Given a dataset  $Z \in Z$  (that the Adversary may eventually contribute to the AI Developer as  $Z_1^a$ ), the goal of the shadow training process is to estimate the data attribution values of the elements in Z as if this dataset were contributed to the AI Developer at t = 1.

We first sample m shadow training datasets, denoted as  $Z^{(1)}, Z^{(2)}, \ldots, Z^{(m)}$ . These datasets are independent of the actual training dataset used by the AI Developer but follow the same distribution. For each shadow training dataset  $Z^{(i)}, i = 1, \ldots, m$ , we train a corresponding shadow



Figure 1: An illustration of the Shadow Attack method. Shadow training datasets  $Z^{(i)}$ 's are sampled to estimate the compensation share of a set of data points Z if it were contributed to the AI Developer, which can be leveraged to perturb data points in Z to get a higher compensation share.

model  $M^{(i)} \in \mathcal{M}$ . The shadow model  $M^{(i)}$  is trained on  $Z^{(i)} \cup Z$  using the same training algorithm  $\mathcal{T}$  as by the AI Developer.<sup>2</sup> i.e.,  $M^{(i)} = \mathcal{T}(Z^{(i)} \cup Z)$ .

In order to estimate the data attribution values, we further sample a shadow validation dataset  $V^{(0)}$ following the same distribution as  $V_0$ . We can estimate a shadow contribution function using each shadow training dataset  $Z^{(i)} \cup Z$ , the corresponding shadow model  $M^{(i)}$ , and the shadow validation dataset  $V^{(0)}$ , resulting in  $\tau^{(i)} = \mathcal{A}(Z^{(i)} \cup Z, M^{(i)}, V^{(0)})$ . Similar to Eq. (1), we consider the following shadow compensation share for the dataset of interest Z:

$$c_s(Z) = \frac{1}{mk|V^{(0)}|} \sum_{i=1}^m \sum_{z \in Z} \sum_{v \in V^{(0)}} \mathbf{1} \left[ \tau^{(i)}(z,v) \in \operatorname{Top}_k(\{\tau^{(i)}(z',v) \mid z' \in Z^{(i)} \cup Z\}) \right].$$
(2)

### 

### 4.2 ADVERSARIAL MANIPULATION BY MAXIMIZING SHADOW COMPENSATION RATE

Given the shadow compensation rate, the natural idea is to apply adversarial perturbations to Zby solving  $Z_1^a = \arg \max_{Z' \in \mathcal{N}(Z)} c_s(Z')$ , where  $\mathcal{N}(Z)$  specifies the space of datasets with unde-tectable perturbations. In practice, however, solving this optimization problem has two technical challenges. Firstly, the objective  $c_s(\cdot)$  is discrete and can be difficult to optimize. Secondly, when doing iterative-style optimization, such as gradient ascent,  $\tau^{(i)}$ 's need to be updated and re-evaluated at every step as it depends on Z and hence  $M^{(i)} = T(Z^{(i)} \cup Z)$ . This leads to two problems: On the one hand, updating  $M^{(i)}$ 's requires retraining, which is computationally heavy; on the other hand, many data attribution methods are computationally expensive, even with all the data and (re-trained) models available. Hence, maximizing  $c_s(\cdot)$  requires repeatedly retraining and running data attributions on perturbed datasets, which may be infeasible for many cases.

To address the first challenge, we replace  $c_s(\cdot)$  with the following surrogate objective

$$c'_{s}(Z) = \frac{1}{mk|V^{(0)}|} \sum_{i=1}^{m} \sum_{z \in Z} \sum_{v \in V^{(0)}} \tau^{(i)}(z, v),$$
(3)

which aims to directly maximize the contribution values of the data points in Z on the shadow validation set  $V^{(0)}$ . To address the second challenge, we first approximate the retrained models by the initial models trained with the original Z (before any gradient ascent steps) for computational efficiency, and we further adopt *Grad-Dot* (Charpiat et al., 2019), one of the most efficient data attribution methods when evaluating the contribution function.<sup>3</sup> Using this method, the contribution value  $\tau^{(i)}(z, v)$  for any training data point z and validation data point v equals the dot product between the gradients of the loss function on z and that on v, evaluated with model  $M^{(i)}$ .

<sup>&</sup>lt;sup>2</sup>In our experiment in Section 6.2, we demonstrate that the Shadow Attack remains effective when the shadow models have a slightly different architecture compared to the target model.

<sup>&</sup>lt;sup>3</sup>Empirically, this works well even when the actual data attribution method used by the AI Developer is more advanced ones, as shown in Section 6.2.

270 With such simplification, the adversarial perturbation can be derived by maximizing Eq. (3), which 271 can be solved efficiently through gradient ascent. Note that since Eq.(3) is a constrained optimization 272 problem, when  $\mathcal{N}(Z)$  is bounded, gradient ascent is guaranteed to converge to some local optimums. 273 In practice, we perform a fixed number of iterations and carefully control the overall perturbation 274 budgets when solving Eq. (3). The computation cost for each iteration is approximately the same as one forward and one backward pass on each of the m shadow models. 275

#### 5 **OUTLIER ATTACK**

276 277

278

281

283 284

287

295

296 297 298

299

279 For large-scale AI applications such as generative AI services, Assumption 2 might be overly strong as it could be difficult for the Adversary to guess the distribution of the full training data. However, in this case, Assumption 3 often holds. In this section, we further propose Outlier Attack, 282 which only relies on Assumption 1 and Assumption 3. See Figure 2 for an illustration.



Figure 2: An illustration of the Outlier Attack method. Here,  $\ell(Z)$  denotes the loss used by the model  $\mathcal{T}(Z_0)$  when evaluated on the dataset Z. The data points in Z are perturbed by maximizing the loss  $\ell(Z)$  through black-box attack methods designed to generate adversarial examples.

#### 5.1 THE OUTLIER INDUCTIVE BIAS OF DATA ATTRIBUTION

300 The core idea behind Outlier Attack leverages an inherent inductive bias present in many data attri-301 bution methods: outlier data points in the training dataset tend to be more influential. Indeed, one of the very first applications of the Influence Function developed in the statistic literature was to detect 302 outliers in training data (Cook, 1977). Consequently, if the Adversary contributes a set of outlier 303 data points, they are more likely to get higher compensation, especially given that the compensation 304 mechanism focuses on the top influential data points. 305

306 However, translating this intuition into a practical adversarial manipulation strategy poses two sig-307 nificant challenges. Firstly, the contributed data points must closely resemble real-world data; otherwise, they could be easily flagged by data quality filters employed by the AI Developer. Sec-308 ondly, we aim to develop an attack method that does not rely on direct knowledge about the training 309 dataset or its underlying distribution. This makes it challenging to determine whether certain data 310 points truly qualify as outliers relative to the training dataset gathered by the AI Developer. 311

312 313

### 5.2 GENERATING REALISTIC OUTLIERS WITH ADVERSARIAL ATTACKS

314 To tackle the first challenge, we propose starting with a set of real-world data and transforming 315 them into outliers through small adversarial perturbations. With careful design, this strategy also 316 addresses the second challenge: if after the perturbation, an AI model has low confidence in correctly 317 predicting its target, then this perturbed data point likely behaves as an outlier relative to the model's 318 training dataset. We discuss two key design aspects for achieving these goals: 319

320 1. Data component to be perturbed. Recall that a data point z = (x, y) consists of both the input feature x and the prediction target y. To obtain an outlier data point from z, 321 we perturb the feature x only, without modifying the target y. Intuitively, flipping the 322 target y could easily result in an outlier data point through mislabeling, such an outlier typically degrades model performance and is likely to be identified as negatively influential.

Moreover, mislabeled data points are easier for the AI Developer to detect. Therefore, it is more effective to perturb only the feature x while keeping the target y unchanged.

2. Objective of perturbation. For each data point z = (x, y), we aim to decrease the model confidence in predicting the annotated target y by perturbing x. In most cases, this is equivalent to increasing the loss  $\ell$  between the model prediction and the target y.

Combining these two design choices, the resulting perturbed data points coincide with what is commonly referred to as *adversarial example* in the literature of adversarial attacks against neural network models. Consequently, existing adversarial attack methods for deriving adversarial examples can be leveraged to generate realistic outliers in our problem setup. Notably, this strategy is very general and can be applied to a variety of data modalities and models by leveraging different offthe-shelf black-box adversarial attack methods.

**Choices of Adversarial Attack Methods.** We consider two concrete machine learning settings, 336 image classification, and text generation, and discuss the choices of adversarial attack methods. For 337 smaller-scale image classification settings, we leverage Zeroth Order Optimization (ZOO) (Chen 338 et al., 2017) based adversarial attacks, which approximates the gradient ascend on the loss function 339 with respect to the data features using black-box queries to the target model  $\mathcal{T}(Z_0)$ . For larger-scale 340 image classification settings, we employ a more advanced black-box adversarial attack method, 341 Simba (Guo et al., 2019), which is computationally more efficient. At a high level, Simba sequen-342 tially perturbs each scalar pixel value in the image by trying to perturb in both directions and accept 343 the perturbation once it increases the loss. For the text generation setting, we utilize TextFooler (Jin et al., 2020), a black-box adversarial attack method tailored for text data. In all the methods of 344 choice, the attack only requires black-box queries to get the predictions of the target model. 345

### 347 5.3 THEORETICAL UNDERSTANDING

The following theorem further provides theoretical insights about the effectiveness of the proposed
 Outlier Attack. The formal statement, notations, and the proof can be found in Appendix A.

**Theorem 5.1** (Informal). Consider a model trained by ERM on a dataset of size n with a smooth loss  $\ell$  with respect to model parameters  $\theta$ . Assume its corresponding influence score  $\tau$ , gradient  $\nabla_{\theta}\ell(\theta, z)$ , and Hessian  $\nabla_{\theta}^{2}\ell(\theta, z)$  are all bounded, i.e.,  $|\tau|, \|\nabla_{\theta}\ell\|_{2}, \|\nabla_{\theta}^{2}\ell\|_{op} = \Theta_{n}(1)$ . Assume the influence score  $\tau$  is based on the influence function by Koh & Liang (2017), which takes the form

 $\tau_{IF}(z_j, z_{test}; \hat{\theta}) = -\nabla_{\theta} \ell(\hat{\theta}, z_{test})^{\top} \left[ \frac{1}{n} \sum_{i=1}^n \nabla_{\theta}^2 \ell(\hat{\theta}, z_i) \right]^{-1} \nabla_{\theta} \ell(\hat{\theta}, z_j),$ (4)

where  $z_{test}$  is a test data point while  $\{z_i\}_{i=1}^n$  are the training data points and  $\hat{\theta}$  is the model parameters trained on  $\{z_i\}_{i=1}^n$ . Then, for any  $z_{test}$ , when  $z_j$  in  $\{z_i\}_{i=1}^n$  is perturbed to  $z'_j$ ,

360 361 362

364

365

359

351

352

353 354

355

356 357

324

326

327

328

$$\tau_{IF}(z_i, z_{test}; \hat{\theta}') = \tau_{IF}(z_i, z_{test}; \hat{\theta}) + O(1/n)$$
 for all  $i \neq j$ , and

• 
$$\tau_{IF}(z'_i, z_{test}; \hat{\theta}') = \tau_{IF}(z'_i, z_{test}; \hat{\theta}) + O(1/n)$$

where  $\hat{\theta}'$  refers to the model parameters trained on the perturbed dataset  $\{z_i\}_{i\neq j} \cup \{z'_j\}$ .

366 Intuitively, perturbing  $z_i$  through adversarial attack will increase the magnitude of the gradi-367 ent  $\nabla_{\theta} \ell(\hat{\theta}, z_i)$ , which tends to also increase the influence score on the original model  $\hat{\theta}$ , i.e., 368  $\tau_{\rm IF}(z'_j, z_{\rm test}; \hat{\theta}) \gg \tau_{\rm IF}(z_j, z_{\rm test}; \hat{\theta})$ . However, in order to get a higher compensation share, the per-369 turbed data point  $z'_i$  needs to have a high influence score on the model  $\hat{\theta}'$  trained on the perturbed 370 dataset, i.e.,  $\tau_{\rm IF}(z'_i, z_{\rm test}; \hat{\theta}')$ , which is not guaranteed by the adversarial attack without characteriz-371 ing the new model  $\hat{\theta}'$ . Theorem 5.1 exactly does this and asserts that  $\tau_{\rm IF}(z'_i, z_{\rm test}; \hat{\theta}')$  will be close 372 373 to  $\tau_{\rm IF}(z'_i, z_{\rm test}; \hat{\theta})$  when the datasets used to train  $\hat{\theta}$  and  $\hat{\theta}'$  are close. It further guarantees that the 374 influence scores of the rest unchanged data points will also remain similar. This explains why the effect of adversarial attacks on the original model can successfully translate to the new model. The 375 results in Theorem 5.1 can be generalized to the case where more than one data point is perturbed, 376 and a small set of clean data is also added to train the new model. The bound will change from 377 O(1/n) to O(k/n), where k is the total number of changed data points in the dataset.

# <sup>378</sup> 6 EXPERIMENTS

379 380

382

In this section, we evaluate the effectiveness of the proposed attack methods by the increase of compensation share after manipulating the data with the proposed attack methods.

# 6.1 EXPERIMENTAL SETUP

385 Tasks, Datasets, Target Models, and Data Attribution Methods. We consider two machine-386 learning tasks: image classification and text generation. For image classification, we experiment 387 on MNIST (LeCun, 1998), Digits (Jiang et al., 2023) and CIFAR-10 (Krizhevsky & Hinton, 2009) 388 datasets, with different target models including Logistic Regression (LR), Multi-layer Perceptron 389 (MLP), Convolutional Neural Networks (CNN), and ResNet-18 (He et al., 2016). We also em-390 ploy three popular data attribution methods, including Influence Function (Koh & Liang, 2017), 391 TRAK (Park et al., 2023), and Data Shapley (Ghorbani & Zou, 2019). For text generation, we conduct experiments on NanoGPT (Karpathy, 2022) trained on the Shakespeare dataset (Karpathy, 392 2015), with TRAK as the data attribution method. Finally, for the image classification settings, we evaluate both Shadow Attack and Outlier Attack, while for the text generation setting, we evaluate 394 Outlier Attack only as Assumption 2 usually does not hold for generative AI settings. For the data attribution algorithms, we adopt the implementation from the dattri library (Deng et al., 2024). The 396 experimental settings are summarized in Table 1. 397

Table 1: Summary of the experimental settings.

Setting	Task	Dataset	Target Model	Attribution Method
(a)	Image Classification	MNIST	LR	Influence Function
(b)	Image Classification	Digits	MLP	Data Shapley
(c)	Image Classification	MNIST	CNN	TRAK
(d)	Image Classification	CIFAR-10	ResNet-18	TRAK
(e)	Text Generation	Shakespeare	NanoGPT	TRAK

408 **Data Contribution Workflow.** For the image classification settings (a), (c), and (d), we set 409  $|Z_0| = 10000, |Z_1^a| = 100$ , and  $|Z_1| = 11000$ . This setup simulates the following data contribution 410 workflow: At t = 0 when there is no Adversary, and 10000 training points are used to train the 411 model  $\mathcal{T}(Z_0)$ . At t = 1, the Adversary contributes 100 perturbed training points and other Data 412 Providers contribute 900 new training data points. Together with the previous 10000 training points, a total of 11000 training points are used to train the model  $\mathcal{T}(Z_1)$ . For image classification 413 setting (b), due to the size of the dataset, we set  $|Z_0| = 1100$ ,  $|Z_1^a| = 30$  and  $|Z_1| = 1100$  for Outlier 414 Attack,  $|Z_0| = 800$ ,  $|Z_1|^a = 30$  and  $|Z_1| = 850$  for Shadow Attack. The text generation setting 415 follows a similar workflow with  $|Z_0| = 4706$ ,  $|Z_1^a| = 20$ , and  $|Z_1| = 6274$ . 416

423 To gain a more refined understanding of how the adversarial perturbations affect the top-k influential 424 data points for individual validation data points in  $V_1$ , we consider the second evaluation metric 425 named **Fraction of Change**. For each validation data point  $v \in V_1$ , it measures how many data 426 points in  $Z_1^a$  appear in the top-k influential data points for v. In comparison to when the original 427 dataset without perturbation is contributed as  $Z_1^a$ , we calculate the fraction of validation data points 428 in  $V_1$  that contain more data points from  $Z_1^a$  in the top-k influential data points when the manipulated dataset after perturbation is contributed as  $Z_1^a$ . Similarly, we calculate the fraction of validation 429 points that contain the same number of or fewer data points from  $Z_1^a$  after perturbation. We report 430 the three fractions under the categories More, Tied, and Fewer, where the higher fraction for the 431 **More** category indicates that the attack method influences the validation data points more broadly.

#### 6.2**EXPERIMENTAL RESULTS: SHADOW ATTACK**

The results of the proposed Shadow Attack method are shown in Table 2. We conduct experiments on all three image classification settings outlined in Table 1. For each of the settings, we first consider the case that the shadow models have the same architecture as the target model, as shown in the first three rows of Table 2. Additionally, for setting (d) where the target model is ResNet-18 (forth row), we further consider using ResNet-9 as the shadow model (last row), which simulates scenarios where the Adversary's knowledge about the training algorithm is limited. 

Table 2: Results of the Shadow Attack method. The target models used for evaluation in each setting are listed in Table 1 while the shadow models used in the attack are listed under **Shadow Model**. The proportion  $|Z_1^a|/|Z_1|$  of the data contributed by the Adversary relative to the full dataset at t = 1 is also reported. A higher **Ratio** indicates a more effective attack, and a higher fraction under More means that the attack influences the validation data points more broadly.

Setting	Shadow Model	$ Z_{1}^{a} / Z_{1} $	<b>Compensation Share</b>			<b>Fraction of Change</b>		
Seeing	Shudow mouth		Original	Manipulated	Ratio	More	Tied	Fewer
(a)	LR	0.0098	0.0098	0.0477	456.1%	0.955	0.038	0.007
(b)	MLP	0.0352	0.0152	0.0435	286.2%	0.533	0.333	0.134
(c)	CNN	0.0098	0.0112	0.0467	417.0%	0.781	0.195	0.024
(d)	ResNet-18	0.0098	0.0095	0.0213	217.3%	0.655	0.259	0.086
(d)	ResNet-9	0.0098	0.0095	0.0196	206.3%	0.622	0.310	0.068

The Shadow Attack method is highly effective in increasing the **Compensation Share** across all the settings. The **Ratio** of the **Manipulated** to the **Original** ranges from 206.3% to 456.1%, represent-ing a substantial increase. Notably, the last row corresponds to the setup where the shadow models have the ResNet-9 architecture while the target model is a ResNet-18, and the Shadow Attack re-mains significantly effective (206.3%), although being slightly worse than the case where shadow model and target model shares the same architecture (forth row, 217.3%.).

The Shadow Attack is also uniformly effective across a wide range of validation data points, as measured by the metrics of **Fraction of Change**. The attack is able to increase the number of top-kinfluential points from  $Z_1^a$  for more than 60% of the validation points. In contrast, only less than 9% of the validation data points have fewer top-k influential points from  $Z_1^a$  after the attack. 

**EXPERIMENTAL RESULTS: OUTLIER ATTACK** 6.3

The results of the proposed Outlier Attack method are presented in Table 3, where we include all four settings summarized in Table 1. The black-box attack method for generating the adversarial examples is chosen according to the discussion in Section 5.2.

Table 3: Results of the Outlier Attack method. The black-box adversarial attack methods for generating adversarial examples are listed under Attack Method. See Table 2 for more context.

Setting	Attack Method	$ Z_{a}^{a} / Z_{1} $	Со	mpensation Sha	Fraction of Change			
octing	Thuck Withou	$ z_1 / z_1 $	Original	Manipulated	Ratio	More	Tied	Fewer
(a)	ZOO	0.0098	0.0098	0.0631	643.9%	0.980	0.017	0.003
(b)	Simba	0.0250	0.0112	0.0218	194.6%	0.440	0.380	0.180
(c)	Simba	0.0098	0.0112	0.0668	596.4%	0.799	0.173	0.028
(d)	Simba	0.0098	0.0095	0.0176	185.2%	0.562	0.354	0.084
(e)	TextFooler	0.0013	0.0035	0.0092	262.9%	0.392	0.461	0.147

On experimental settings (a) and (c), the Outlier Attack performs even better than the Shadow Attack in terms of both Compensation Share and Fraction of Change, achieving a higher Ratio and a larger fraction under **More**. Across all four settings, the **Ratio** ranges from 185.2% to 643.9%, demonstrating the exceptional effectiveness of the Outlier Attack. Finally, the results of text gener-ation setting (e) further highlight the applicability of the proposed method to generative AI models.

# 486 6.4 BASELINE REFERENCE

495

504

505 506

507

508

529 530

531

To better understand the significance of the results in Table 2 and Table 3, we compare them with a baseline random perturbation method in Table 4 on the three image classification settings such that the baseline method applies pixel-wise random perturbation to the images, with the perturbation budget matched to that of the proposed attack methods. As shown in the results, the **Compensation Share** of the dataset with **Random Perturbation** is nearly identical to that of the **Original**. The fraction under **More** is also close to that under **Fewer**. These results highlight the effectiveness of the proposed attack methods comes from the careful design of adversarial perturbation.

Table 4: Results of the random perturbation baseline. See Table 2 for more context.

Setting	Сог	npensation Share	Fract	Fraction of Change		
Setting	Original	<b>Random Perturbation</b>	More	Tied	Fewer	
(a)	0.0098	0.0097	0.203	0.586	0.211	
(c)	0.0112	0.0117	0.385	0.326	0.289	
(d)	0.0095	0.0125	0.367	0.430	0.203	

# 6.5 VISUALIZATION OF THE PERTURBATION

Finally, Figure 3 provides visualizations of two examples of the adversarial perturbed images from MNIST and CIFAR-10. In both cases, the perturbations are barely visible to human eyes, however, they are highly effective in terms of the success of the attack.



# 7 CONCLUSION

532 This work addresses a significant gap in the current understanding of the adversarial robustness of 533 data attribution methods, which are increasingly influential in data valuation and compensation ap-534 plications. By introducing a well-defined threat model, we have proposed two novel adversarial at-535 tack strategies, Shadow Attack and Outlier Attack, which are designed to manipulate data attribution 536 and inflate compensation. The Shadow Attack utilizes knowledge of the underlying data distribu-537 tion, while the Outlier Attack operates without such knowledge, relying on black-box queries only. Empirical results from image classification and text generation tasks demonstrate the effectiveness of 538 these attacks, with compensation inflation ranging from 185% to 643%. These findings underscore the need for more robust data attribution methods to guard against adversarial exploitation.

# 540 REFERENCES

546

547

553

554

555

558

562

563

564

565

577

578

579

580

- Anirban Chakraborty, Manaar Alam, Vishal Dey, Anupam Chattopadhyay, and Debdeep Mukhopad hyay. A survey on adversarial attacks and defences. *CAAI Transactions on Intelligence Technol- ogy*, 6(1):25–45, 2021.
  - Guillaume Charpiat, Nicolas Girard, Loris Felardos, and Yuliya Tarabalka. Input similarity from the neural network perspective. *Advances in Neural Information Processing Systems*, 32, 2019.
- Pin-Yu Chen, Huan Zhang, Yash Sharma, Jinfeng Yi, and Cho-Jui Hsieh. Zoo: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models. In *Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security*, AISec '17, pp. 15–26, New York, NY, USA, 2017. Association for Computing Machinery. ISBN 9781450352024.
  - R Dennis Cook. Detection of influential observation in linear regression. *Technometrics*, 19(1): 15–18, 1977.
- Junwei Deng and Jiaqi Ma. Computational copyright: Towards a royalty model for ai music gener ation platforms. *arXiv preprint arXiv:2312.06646*, 2023.
- Junwei Deng, Ting-Wei Li, Shiyuan Zhang, Shixuan Liu, Yijun Pan, Hao Huang, Xinhe Wang,
  Pingbang Hu, Xingjian Zhang, and Jiaqi W Ma. dattri: A library for efficient data attribution. *arXiv preprint arXiv:2410.04555*, 2024.
  - Vitaly Feldman and Chiyuan Zhang. What neural networks memorize and why: Discovering the long tail via influence estimation. Advances in Neural Information Processing Systems, 33:2881– 2891, 2020.
- Amirata Ghorbani and James Zou. Data shapley: Equitable valuation of data for machine learning.
   In *International conference on machine learning*, pp. 2242–2251. PMLR, 2019.
- Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial
  examples. In Yoshua Bengio and Yann LeCun (eds.), *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceed*ings, 2015.
- 573 Chuan Guo, Jacob Gardner, Yurong You, Andrew Gordon Wilson, and Kilian Weinberger. Simple
  574 black-box adversarial attacks. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceed-*575 *ings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of*576 *Machine Learning Research*, pp. 2484–2493. PMLR, 09–15 Jun 2019.
  - Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Hongsheng Hu, Zoran Salcic, Lichao Sun, Gillian Dobbie, Philip S Yu, and Xuyun Zhang. Membership inference attacks on machine learning: A survey. *ACM Computing Surveys (CSUR)*, 54 (11s):1–37, 2022.
- Andrew Ilyas, Logan Engstrom, Anish Athalye, and Jessy Lin. Black-box adversarial attacks with
   limited queries and information. In *International conference on machine learning*, pp. 2137–2146. PMLR, 2018.
- 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 *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 1167–1176. PMLR, 2019.
- 592
  - Kevin Jiang, Weixin Liang, James Y Zou, and Yongchan Kwon. Opendataval: a unified benchmark for data valuation. *Advances in Neural Information Processing Systems*, 36, 2023.

594	Di Jin, Zhijing Jin, Joey Tianyi Zhou, and Peter Szolovits. Is BERT really robust? A strong baseline
595	for natural language attack on text classification and entailment. In The Thirty-Fourth AAAI Con-
596	ference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Arti-
597	ficial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances
598	in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020, pp. 8018–8025.
599	AAAI Press, 2020.

- 600 Hoang Anh Just, Feiyang Kang, Tianhao Wang, Yi Zeng, Myeongseob Ko, Ming Jin, and Ruoxi 601 Jia. LAVA: data valuation without pre-specified learning algorithms. In The Eleventh Interna-602 tional Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023. 603 OpenReview.net, 2023. 604
- 605 Andrej Karpathy. char-rnn. https://github.com/karpathy/char-rnn, 2015.
- Andrej Karpathy. nano-gpt, 2022. 607

606

612

619

625

637

- 608 Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In 609 Doina Precup and Yee Whye Teh (eds.), Proceedings of the 34th International Conference on 610 Machine Learning, volume 70 of Proceedings of Machine Learning Research, pp. 1885–1894. 611 PMLR, 06–11 Aug 2017.
- Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. Mas-613 ter's thesis, Department of Computer Science, University of Toronto, 2009. 614
- 615 Yongchan Kwon and James Zou. Beta shapley: a unified and noise-reduced data valuation frame-616 work for machine learning. In Gustau Camps-Valls, Francisco J. R. Ruiz, and Isabel Valera (eds.), 617 Proceedings of The 25th International Conference on Artificial Intelligence and Statistics, volume 618 151 of Proceedings of Machine Learning Research, pp. 8780-8802. PMLR, 28-30 Mar 2022.
- Yongchan Kwon and James Zou. Data-oob: Out-of-bag estimate as a simple and efficient data value. 620 In International Conference on Machine Learning, pp. 18135–18152. PMLR, 2023. 621
- 622 Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. CS 231N, 7(7):3, 2015. 623
- 624 Yann LeCun. The mnist database of handwritten digits. http://yann. lecun. com/exdb/mnist/, 1998.
- Jinkun Lin, Anqi Zhang, Mathias Lécuyer, Jinyang Li, Aurojit Panda, and Siddhartha Sen. Mea-626 suring the effect of training data on deep learning predictions via randomized experiments. In 627 International Conference on Machine Learning, pp. 13468–13504. PMLR, 2022. 628
- 629 Aleksander Mkadry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. 630 Towards deep learning models resistant to adversarial attacks. *stat*, 1050(9), 2017. 631
- Elisa Nguyen, Minjoon Seo, and Seong Joon Oh. A bayesian approach to analysing training data 632 attribution in deep learning. Advances in Neural Information Processing Systems, 36, 2024. 633
- 634 Sung Min Park, Kristian Georgiev, Andrew Ilyas, Guillaume Leclerc, and Aleksander Madry. 635 TRAK: Attributing model behavior at scale. In Andreas Krause, Emma Brunskill, Kyunghyun 636 Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), Proceedings of the 40th International Conference on Machine Learning, volume 202 of Proceedings of Machine Learning 638 *Research*, pp. 27074–27113. PMLR, 23–29 Jul 2023.
- 639 Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference at-640 tacks against machine learning models. In 2017 IEEE Symposium on Security and Privacy (SP), 641 pp. 3-18, 2017. doi: 10.1109/SP.2017.41. 642
- 643 James Seale Smith, Yen-Chang Hsu, Lingyu Zhang, Ting Hua, Zsolt Kira, Yilin Shen, and Hongxia 644 Jin. Continual diffusion: Continual customization of text-to-image diffusion with c-loRA. Trans-645 actions on Machine Learning Research, 2024. ISSN 2835-8856.
- Liwei Song and Prateek Mittal. Systematic evaluation of privacy risks of machine learning models. 647 In 30th USENIX Security Symposium (USENIX Security 21), pp. 2615–2632, 2021.

648 649	G W Stewart and Ji-Guang Sun. <i>Matrix Perturbation Theory</i> . Computer Science and Scientific Computing. Academic Press, June 1990. ISBN 9780080926131.
651 652 653	Jiachen T Wang and Ruoxi Jia. Data banzhaf: A robust data valuation framework for machine learning. In International Conference on Artificial Intelligence and Statistics, pp. 6388–6421. PMLR, 2023.
654 655	Jiachen T. Wang, Zhun Deng, Hiroaki Chiba-Okabe, Boaz Barak, and Weijie J. Su. An economic solution to copyright challenges of generative ai, 2024.
657 658 659	Xinyi Xu, Zhaoxuan Wu, Chuan Sheng Foo, and Bryan Kian Hsiang Low. Validation free and repli- cation robust volume-based data valuation. <i>Advances in Neural Information Processing Systems</i> , 34:10837–10848, 2021.
660 661	Jinsung Yoon, Sercan Arik, and Tomas Pfister. Data valuation using reinforcement learning. In <i>International Conference on Machine Learning</i> , pp. 10842–10851. PMLR, 2020.
662 663 664 665	Xiaoyong Yuan, Pan He, Qile Zhu, and Xiaolin Li. Adversarial examples: Attacks and defenses for deep learning. <i>IEEE transactions on neural networks and learning systems</i> , 30(9):2805–2824, 2019.
666 667 668 669	Yang Zhang, Fuli Feng, Chenxu Wang, Xiangnan He, Meng Wang, Yan Li, and Yongdong Zhang. How to retrain recommender system? a sequential meta-learning method. In <i>Proceedings of</i> <i>the 43rd International ACM SIGIR Conference on Research and Development in Information</i> <i>Retrieval</i> , pp. 1479–1488, 2020.
670 671 672 673 674	Zihan Zhang, Meng Fang, Ling Chen, Mohammad-Reza Namazi-Rad, and Jun Wang. How do large language models capture the ever-changing world knowledge? a review of recent advances. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pp. 8289–8311, Singapore, December 2023. Association for Computational Linguistics.
675 676 677	
678 679	
680	
681	
682	
68/	
685	
686	
687	
688	
689	
690	
691	
692	
693	
694	
695	
696	
697	
698	
699	
700	

## A FORMAL STATEMENT OF THEOREM 5.1 AND ITS PROOF

### A.1 SETUP AND THE STATEMENT

Let the original training set be  $Z = \{z_i\}_{i=1}^n$ . Given a data point z, assume the loss  $\ell(\theta, z)$  is twice differentiable. Then, the empirical loss and the Hessian is given by

$$\ell(\theta) = \frac{1}{n} \sum_{i=1}^{n} \ell(\theta, z_i), \quad H(\theta) = \frac{1}{n} \sum_{i=1}^{n} h(\theta, z_i),$$

where we define  $h(\theta, z) = \nabla_{\theta}^2 \ell(\theta, z)$ . Moreover, we consider the ERM to be

$$\hat{\theta} = \operatorname*{arg\,min}_{\theta} \frac{1}{n} \sum_{i=1}^{n} \ell(\theta, z_i)$$

<sup>716</sup> Consider the training data attribution to be the influence estimated by the *influence function* (Koh & Liang, 2017)  $\tau_{\text{IF}}$ . In particular, for a train test pair  $(z_i, z_{\text{test}})$ , Eq.(4) can be succinctly written as

$$\tau_{\mathrm{IF}}(z_i, z_{\mathrm{test}}; \hat{\theta}) = -\nabla_{\theta} \ell(\hat{\theta}, z_i)^{\top} H^{-1}(\hat{\theta}) \nabla_{\theta} \ell(\hat{\theta}, z_{\mathrm{test}}).$$

Consider perturbing the  $j^{\text{th}}$  training point such that  $z_j$  becomes  $z'_j$ . Then, the training set turns to  $\{z_1, \ldots, z'_j, \ldots, z_n\}$ , which in turn changes the empirical loss and its hessian to

$$\ell'(\theta) = \frac{1}{n} \sum_{i \neq j} \ell(\theta, z_i) + \frac{1}{n} \ell(\theta, z'_j), \quad H'(\theta) = \frac{1}{n} \sum_{i \neq j} h(\theta, z_i) + \frac{1}{n} h(\theta, z'_j).$$

It also follows that the minimizer for this perturbed training set is given by

$$\hat{\theta}' = \operatorname*{arg\,min}_{\theta} \frac{1}{n} \sum_{i \neq j} \ell(\theta, z_i) + \frac{1}{n} \ell(\theta, z'_j),$$

with the corresponding influence score being

• 
$$\tau_{\text{IF}}(z_i, z_{\text{test}}; \hat{\theta}') = -\nabla_{\theta} \ell(\hat{\theta}', z_{\text{test}})^{\top} (H'(\hat{\theta}'))^{-1} \nabla_{\theta} \ell(\hat{\theta}', z_i) \text{ for all } i \neq j; \text{ and }$$

• 
$$\tau_{\mathrm{IF}}(z'_i, z_{\mathrm{test}}; \hat{\theta}') = -\nabla_{\theta} \ell(\hat{\theta}', z_{\mathrm{test}})^{\top} (H'(\hat{\theta}'))^{-1} \nabla_{\theta} \ell(\hat{\theta}', z'_i).$$

**Remark 1.** From our definition of data attribution methods,  $\tau$  should take the form of  $Z \times V \to \mathbb{R}$ for some training dataset Z and validation dataset V. In the above notation, we are essentially considering two different training attribution functions  $\tau_{IF}$  and  $\tau'_{IF}$  where  $\tau_{IF} = \mathcal{A}(Z, \mathcal{T}(Z), V)$  and  $\tau'_{IF} = \mathcal{A}(Z', \mathcal{T}(Z'), V)$  where Z' is the perturbed dataset,  $\hat{\theta} = \mathcal{T}(Z)$  and  $\hat{\theta}' = \mathcal{T}(Z')$ . For clarity, we explicitly write  $\tau_{IF}(\cdot, \cdot) = \tau_{IF}(\cdot, \cdot; \hat{\theta})$  and  $\tau'_{IF}(\cdot, \cdot) = \tau_{IF}(\cdot, \cdot; \hat{\theta}')$  to avoid confusion.

741 Under this setup, we can now state the theorem formally.

**Theorem A.1.** Under the above setup, consider a model trained by ERM with a loss  $\ell$  that is twice-differentiable, m-strongly convex, and L-Lipschitz continuous with respect to  $\theta$ . Assume its corresponding influence score  $\tau_{IF}$ , gradient  $\nabla_{\theta}\ell(\theta, z)$ , and  $h(\theta, z)$  are all bounded, i.e.,  $|\tau_{IF}| = \Theta(1)$ ,  $||\nabla_{\theta}\ell||_2 = \Theta(1)$ ,  $||h||_{op} = \Theta(1)$ . Then, given any test data point  $z_{test}$ , we have

• 
$$\tau_{IF}(z_i, z_{test}; \hat{\theta}') = \tau_{IF}(z_i, z_{test}; \hat{\theta}) + O(1/n)$$
 for all  $i \neq j$ , and

• 
$$\tau_{IF}(z'_i, z_{test}; \hat{\theta}') = \tau_{IF}(z'_i, z_{test}; \hat{\theta}) + O(1/n).$$

A.2 TECHNICAL LEMMAS

<sup>752</sup> We first establish several technical lemmas toward proving Theorem A.1.

**Lemma 1.** Let  $\hat{\theta}, \hat{\theta}'$  be the minimizer for  $\ell(\theta), \ell'(\theta)$  respectively, then

$$\|\hat{\theta} - \hat{\theta}'\| \leq \frac{4L}{mn}$$

*Proof.* Firstly, we recall that from definition, f(x) being m-strongly convex means that for any x, y, 

$$f(y) \ge f(x) + \nabla f(x)^{\top} (y - x) + \frac{m}{2} \|y - x\|_2^2.$$

Denote  $x^* = \arg \min_x f(x)$ , then since  $\nabla f(x^*) = 0$ , for all x, we have

$$f(x) \ge f(x^*) + \frac{m}{2} ||x - x^*||_2^2$$

Hence, by strong convexity of  $\ell(\theta, z)$  and the fact that  $\hat{\theta}$  is the minimizer of  $\ell(\theta)$ , we have

$$\frac{1}{n}\sum_{i=1}^{n}\ell(\hat{\theta}', z_i) \ge \frac{1}{n}\sum_{i=1}^{n}\ell(\hat{\theta}, z_i) + \frac{m}{2}\|\hat{\theta} - \hat{\theta}'\|_2^2$$

On the other hand, we have

$$\frac{1}{n}\sum_{i=1}^{n}\ell(\hat{\theta}',z_i) = \underbrace{\frac{1}{n}\sum_{i\neq j}\ell(\hat{\theta}',z_i) + \frac{1}{n}\ell(\hat{\theta}',z'_j)}_{\ell'(\hat{\theta}')} + \frac{1}{n}(\ell(\hat{\theta}',z_j) - \ell(\hat{\theta}',z'_j))$$

 $\leq \underbrace{\frac{1}{n}\sum_{i\neq j}\ell(\hat{\theta},z_i) + \frac{1}{n}\ell(\hat{\theta},z'_j) + \frac{1}{n}(\ell(\hat{\theta}',z_j) - \ell(\hat{\theta}',z'_j))}_{\ell'(\hat{\theta})} \quad (\hat{\theta}' \text{ is a minimizer of } \ell')$ 

$$= \frac{1}{n} \sum_{i=1}^{n} \ell(\hat{\theta}, z_i) + \frac{1}{n} (\ell(\hat{\theta}', z_j) - \ell(\hat{\theta}, z_j)) + \frac{1}{n} (\ell(\hat{\theta}, z'_j) - \ell(\hat{\theta}', z'_j))$$
  
$$\le \frac{1}{n} \sum_{i=1}^{n} \ell(\hat{\theta}, z_i) + \frac{2}{n} L \|\hat{\theta} - \hat{\theta}'\|_2.$$
 ( $\ell(\theta, z)$  is *L*-Lipschitz w.r.t.  $\theta$ )

Combine the above results, we see that

$$\frac{m}{2} \|\hat{\theta} - \hat{\theta}'\|_{2}^{2} \le \frac{2L}{n} \|\hat{\theta} - \hat{\theta}'\|_{2}$$

which proves the desired result. 

**Lemma 2** (Section 2.4, Part III (Stewart & Sun, 1990)). For any  $A, E \in \mathbb{R}^{d \times d}$  with both A and A + E being invertible, if  $\sum_{k=0}^{\infty} ||E||_{op}^{k}$  converges, we have 

$$(A+E)^{-1} = A^{-1} - A^{-1}EA^{-1} + O(||E||_{op}^2).$$

**Lemma 3.** Let H and H' be the Hessian of  $\ell$  and  $\ell'$ , respectively, then

$$||H'^{-1} - H^{-1}||_{\text{op}} \le O\left(\frac{4LM}{m^3n}\right)$$

*Proof.* Since  $h(\theta, z)$  is *M*-Lipschitz w.r.t.  $\theta$ , we know that

$$||H' - H||_{\text{op}} \le M ||\hat{\theta}' - \hat{\theta}||_2 \le \frac{4LM}{mn}.$$

With Lemma 2 (let A = H and E = H' - H, it's easy to verify the conditions of Lemma 2 hold), we have  $H'^{-1} = H^{-1} - H^{-1}(H' - H)H^{-1} + O(1/n^2)$ , which gives

$$||H'^{-1} - H^{-1}||_{\rm op} \le ||H^{-1}||_{\rm op} ||H' - H||_{\rm op} ||H^{-1}||_{\rm op} + O(1/n^2).$$

Since  $\ell(\theta, z)$  is *m*-strongly convex w.r.t.  $\theta$ , it's easy to show that  $||H^{-1}||_{op} \leq 1/m$ . In all, we have

г		

A.3 PROOF OF THEOREM A.1

## 812 We can now prove Theorem A.1.

*Proof of Theorem A.1.* For all  $i \neq j$ , we want to prove that  $\tau_{\text{IF}}(z_i, z_{\text{test}}; \hat{\theta}') = \tau_{\text{IF}}(z_i, z_{\text{test}}; \hat{\theta}) + O(1/n)$ . Indeed, since

$$\begin{aligned} \tau_{\mathrm{IF}}(z_{i}, z_{\mathrm{test}}; \hat{\theta}') &= -\nabla_{\theta} \ell(\hat{\theta}', z_{\mathrm{test}})^{\top} (H'(\hat{\theta}'))^{-1} \nabla_{\theta} \ell(\hat{\theta}', z_{i}) \\ &= \left( \nabla_{\theta} \ell(\hat{\theta}, z_{\mathrm{test}}) - \nabla_{\theta} \ell(\hat{\theta}', z_{\mathrm{test}}) \right)^{\top} (H'(\hat{\theta}'))^{-1} \nabla_{\theta} \ell(\hat{\theta}, z_{i}) \\ &- \nabla_{\theta} \ell(\hat{\theta}, z_{\mathrm{test}}) \left( (H'(\hat{\theta}'))^{-1} - H^{-1}(\hat{\theta}) \right) \nabla_{\theta} \ell(\hat{\theta}, z_{i}) \\ &- \nabla_{\theta} \ell(\hat{\theta}, z_{\mathrm{test}}) H^{-1}(\hat{\theta}) \nabla_{\theta} \ell(\hat{\theta}, z_{i}) \\ &= \left( \nabla_{\theta} \ell(\hat{\theta}, z_{\mathrm{test}}) - \nabla_{\theta} \ell(\hat{\theta}', z_{\mathrm{test}}) \right)^{\top} (H'(\hat{\theta}'))^{-1} \nabla_{\theta} \ell(\hat{\theta}, z_{i}) \\ &- \nabla_{\theta} \ell(\hat{\theta}, z_{\mathrm{test}}) \left( (H'(\hat{\theta}'))^{-1} - H^{-1}(\hat{\theta}) \right) \nabla_{\theta} \ell(\hat{\theta}, z_{i}) \\ &+ \tau_{\mathrm{IF}}(z_{i}, z_{\mathrm{test}}; \hat{\theta}). \end{aligned}$$

Applying Lemmas 1 and 3, and by noting that when the Hessian is *M*-Lipschitz, so is the gradient in terms of  $\theta$ , hence we have the desired result. Specifically, we have

 $\leq \|\nabla_{\theta}\ell(\hat{\theta}, z_{\text{test}}) - \nabla_{\theta}\ell(\hat{\theta}', z_{\text{test}})\|_2 \|H'^{-1}\|_{\text{op}} \|\nabla_{\theta}\ell(\hat{\theta}, z_i)\|_2$ 

+  $\|\nabla_{\theta}\ell(\hat{\theta}, z_{\text{test}})\|_{2} \|H'^{-1} - H^{-1}\|_{\text{op}} \|\nabla_{\theta}\ell(\hat{\theta}, z_{i})\|_{2}$ 

 $+ \|\nabla_{\theta} \ell(\hat{\theta}, z_{\text{test}})\|_{2} \|H'^{-1} - H^{-1}\|_{\text{op}} \|\nabla_{\theta} \ell(\hat{\theta}, z_{i})\|_{2}$ 

 $\leq \|\nabla_{\theta} \ell(\hat{\theta}, z_i)\|_2 \left( \|H^{-1}\|_{\text{op}} + \|H'^{-1} - H^{-1}\|_{\text{op}} \right) \cdot M \|\hat{\theta}' - \hat{\theta}\|_2$ 

829 830 831

827

828

813 814

815

832 833

836 837

838 839

840 841

842 843 844

845 846

847

For the second case, the proof is the same by replacing  $z_i$  with  $z'_j$  in the above calculation.

 $= \Theta(1) \cdot \left(\frac{1}{m} + O\left(\frac{4LM}{m^3n}\right)\right) \cdot M \cdot \frac{4L}{mn} + \Theta(1) \cdot O\left(\frac{4LM}{m^3n}\right) \cdot \Theta(1) = O\left(\frac{1}{n}\right).$ 

# **B** EXPERIMENT DETAILS

In this section, we introduce the experiment details.

 $|\tau_{\text{IF}}(z_i, z_{\text{test}}; \hat{\theta}') - \tau_{\text{IF}}(z_i, z_{\text{test}}; \hat{\theta})|$ 

# B.1 DATASETS AND MODELS

Logistic Regression and Convolutional Neural Network on MNIST. For experiments on 848 MNIST, we consider two different target models, Logistic Regression (LR) and convolutional neu-849 ral network (CNN). The CNN comprises two convolutional layers: the first with 32 filters and the 850 second with 64 filters (size  $3 \times 3$ , stride 1, padding 1), both followed by ReLU and max-pooling 851 layers. The output is flattened into a vector of size  $64 \times 7 \times 7$ , which is then passed through a fully 852 connected layer with 128 units before reaching the output layer with 10 classes. Both target models 853 are trained on the first 10000 training points of the MNIST dataset. For shadow models, we train 50 854 of them, each on a randomly sampled subset of 5000 among the second 10000 training data points 855 to ensure no overlaps with the first 10000 training data points used for target model training. 856

In both the target model training and shadow model training, we train LR for 30 epochs with SGD and a learning rate of 0.01. We train CNN for 50 epochs with Adam and a learning rate of 0.001.

MLP on Digits. For experiments on Digits, we consider the target MLP model with 5 hidden layers, each has 10 hidden neurons. Due to the size limit of the Digits dataset, we train the target model for an Outlier Attack using the first 1100 training points, while for the shadow attack, we train the target model on the first 800 data points and the shadow models on the second 800 data points. In all training, we again train the MLP for 30 epochs using Adam with an initial learning rate of 0.001.

864 **ResNet18 on CIFAR-10.** For experiments on CIFAR-10, we consider the classical ResNet18 model without dropouts. We use 10000 training points for training the target model. For shadow 866 model training, we consider the classical ResNet-9 and ResNet-18 as the shadow models separately, 867 again without dropouts. We train 50 shadow models, and each on a subset of 10000 training points 868 that are disjoint with the training set of the target model. In all training, each model is trained for 100 epochs using Adam, with a learning rate of 0.001.

NanoGPT on Shakespeare. For experiments on the Shakespeare dataset, we consider the 871 NanoGPT model, which is a character-level GPT model with 4 layers, 4 heads, and 128 dimen-872 sion of the embedding. The block size is 64, and the batch size is 32. For both training at t = 0 and 873 t = 1, we train the model for 2000 epochs, both using Adam with a learning rate of  $6 \times 10^{-4}$ . Note 874 that since we only consider Outlier Attack for this setup, hence there is no shadow model training 875 involved.

876 877 878

879

880

881

882 883

884

885 886

889

891

892

893

894

899 900

901 902

903

904

905

906 907

908

909

870

#### B.2 TRAINING DATA ATTRIBUTION METHODS

In this section, we briefly introduce the attribution methods that are included in the evaluation. Given a training dataset  $\{z_i\}_{i=1}^n$ , we are interested in the data attribution of a particular training data point  $z_j$  and a test data point  $z_{\text{test}}$ .

**Influence Function based on the Conjugate Gradients.** As we have seen, the original definition of influence function (Koh & Liang, 2017) is given by Eq. (4), i.e.,

$$\tau_{\rm IF}(z_j, z_{\rm test}; \hat{\theta}) = -\nabla_{\theta} \ell(\hat{\theta}, z_j)^{\top} H(\hat{\theta})^{-1} \nabla_{\theta} \ell(\hat{\theta}, z_{\rm test})$$

887 where  $\nabla_{\theta} \ell(\hat{\theta}, z)$  is the gradient of loss of the data point w.r.t. model parameters, and  $H(\hat{\theta})^{-1}$  is 888 the inverse of the Hessian w.r.t. model parameters. We implement the conjugate gradients (CG) approach to compute the inverse-Hessian-vector-product (IHVP). 890

Tracing with the Randomly-projected After Kernel (TRAK). Introduced by Park et al. (2023), TRAK computes the attribution score for a training test pair by first linearizing the model and then applying random projection to make the computation efficient. Denote the output (e.g., raw logit) of the model  $\hat{\theta}$  to be  $f(z; \hat{\theta})$ , then the *naive TRAK influence* can be formulated as

$$\tau_{\text{TRAK}}(z_j, z_{\text{test}}; \hat{\theta}) = -(1 - p_j^{\star})\phi(z_j)^{\top} (\Phi^{\top} \Phi)^{-1} \phi(z_{\text{test}}),$$

where  $\phi(z) = P^{\top} \nabla_{\theta} f(z; \hat{\theta})$  is the random projection of  $\nabla_{\theta} f(z; \hat{\theta}) \in \mathbb{R}^p$  by some Gaussian random projection matrix  $P \sim \mathcal{N}(0, 1)^{p \times k}$ , and  $\Phi$  is the matrix formed by stacking all the  $\phi(z_i)$ , and  $p_i^{\star}$  is the predicted correct-class probability of  $z_i$  at  $\theta$ . Compare it with  $\tau_{\text{IF}}$ :

- 1. The additional factor  $1 p_j^*$  in  $\tau_{\text{TRAK}}$  is due to linearizing the model.
- 2. For linear model with feature matrix X of the training set, its Hessian is exactly  $X^{\top}X$ . For for a linearized the model, each feature  $x_i$  of training sample  $z_i$  corresponds to the gradient  $\nabla_{\theta} f(z_i; \hat{\theta})$ , i.e.,  $X = [\nabla_{\theta} f(z_i; \hat{\theta})]_{i=1}^n$ .
- 3. With random projection, the linearized features  $x_i = \nabla_{\theta} f(z_i; \hat{\theta})$  becomes  $\phi(z_i) =$  $P \nabla_{\theta} f(z_i; \hat{\theta})$ , which induces a new feature matrix  $\Phi = [\phi(z_i)]_{i=1}^n$

910 Hence, overall,  $\tau_{\text{TRAK}}$  is nothing but the influence function applied to the linearized model with ran-911 dom projection. We note that the original TRAK influence includes a step called *ensembling*, which 912 is simply averaging the above TRAK influence over multiple  $\tau_{\text{TRAK}}$  with models independently 913 trained on a subset of the training set. 914

915 Grad-Dot. Introduced by Charpiat et al. (2019), the dot product of gradient, known as Grad-Dot, is a simple, easy-to-compute training data attribution method. It is given by 916

$$\tau_{\text{Grad-Dot}}(z_j, z_{\text{test}}; \hat{\theta}) = -\nabla_{\theta} \ell(\hat{\theta}, z_j)^\top \nabla_{\theta} \ell(\hat{\theta}, z_{\text{test}})$$

**Data Shapley.** Data Shapley (Ghorbani & Zou, 2019) quantifies the influence of individual data points by considering its marginal contribution to different subsets of the training set. It is motivated by the so-called leave-one-out (LOO) influence: specifically, given any aspect  $\phi(\hat{\theta})$  we care about for the learned model  $\hat{\theta}$  (e.g., the loss  $\ell(\hat{\theta}, z_{test})$  of some test data point  $z_{test}$ ), LOO measures the influence of every individual data point  $z_j$  on  $\phi$  by the difference  $\phi(\hat{\theta}_{-z_j}) - \phi(\hat{\theta})$ , where  $\hat{\theta}_{-z_j}$  is the model learned with the dataset  $\{z_i\}_{i \neq j}$  that excludes  $z_j$ . Data Shapley builds on top of LOO by requiring additional *equitable* conditions, resulting in the following formulation

$$\tau_{\text{Data-Shapley}}(z_j, z_{\text{test}}; \hat{\theta}) = \sum_{S \subseteq \{z_i\}_{i \neq j}} \frac{\phi(\hat{\theta}_S, z_{\text{test}}) - \phi(\hat{\theta}_{S \cup \{j\}}, z_{\text{test}})}{\binom{n-1}{|S|}},$$

where we let  $\phi(\hat{\theta}, z_{\text{test}}) = \ell(\hat{\theta}, z_{\text{test}})$  and  $\hat{\theta}_S$  for some  $S \subseteq \{z_i\}_{i=1}^n \setminus \{z_j\}$  denotes the model learned on the subset S. This can be viewed as an averaged version of LOO while satisfying several equitable conditions, which we refer to Ghorbani & Zou (2019) for details.

**Remark 2.** We note that in the above formulations, we incorporate a negative sign in front of  $\tau_{TRAK}$ ,  $\tau_{Grad-Dot}$  to make them consistent to the original formulation of  $\tau_{IF}$ . Specifically, in Koh & Liang (2017), the influence change is measured as the difference between the perturbed model and the original model, while others consider the opposite.

### B.3 ATTACK METHODS

926 927 928

933

934

935

936 937

938

945

946

947 948

949

950

951

956

957

958

959

960

961

962

963

964

965

In this section, we detail the implementation of all the attack methods we have used throughout the experiments.

941 Shadow Attack. Shadow attack is a classical adversarial attack method that originates from 942 (Shokri et al., 2017). The details of shadow model training can be found in Appendix B.1. As 943 introduced in Section 4, after 50 shadow models are trained, we optimize Eq. (3) using gradient 944 ascent with respect to input feature x with 10 iterations and a step size of  $\epsilon = 0.01$ .

**Outlier Attack.** For Outlier Attack, we utilize the following black-box attack methods to produce adversarial examples:

Zeroth Order Optimization (ZOO): ZOO is a black-box attack method that approximates gradient through finite numerical methods. In our experiment, given an input z = (x, y), we perturb it by x' ← x + ε ⋅ sgn(g(θ̂, z)), where g(θ̂, z) is an estimation of the loss of the gradient with respect to x, given by the symmetric difference quotient

$$\left(g(\hat{\theta},z)\right)_i = \frac{\ell(\hat{\theta},(x+h\mathbf{e}_i,y)) - \ell(\hat{\theta},(x-h\mathbf{e}_i,y))}{2h},$$

along the  $i^{\text{th}}$  standard basis direction  $\mathbf{e}_i$  with  $h \in \mathbb{R}$ . We note that  $\ell(\hat{\theta}, (x, y))$  can be obtained by black-box queries of the target model. In experiments, we set  $\epsilon = 0.03$ .

- Simba: Introduced by Guo et al. (2019), Simba proposes a black-box attack method to perturb each input pixel sequentially after permutation. More precisely, each input pixel is attempted to be perturbed in both directions, i.e.,  $(x')_i \leftarrow x_i \pm \epsilon$ , respectively. After a perturbation is attempted, the target model is queried again, and the perturbation is accepted as long as this perturbation increases the loss. In practice, we set  $\epsilon = 0.1$ .
- **TextFooler**: TextFooler is originally an attack method for text *classification* task (Jin et al., 2020). It perturbs texts following a two-step approach: for a piece of text that it wants to perturb, it first sorts the words by their importance, and then replaces the influential words to increase the loss of prediction while preserving text similarity.
- 966 We modify the method as follows. Firstly, for a character-level GPT model (in our case, 967 NanoGPT), we work on characters rather than words. Secondly, as we work on text *gen-*968 *eration* task, for a sequence of characters that we want to perturb, we take the sum of the 969 negative log-likelihood of predicting its m future characters as a loss. We then sort the 970 characters by their importance, where importance is measured by the loss increase after the 971 character is masked. Finally, the k-most important characters are replaced by characters 972 that maximize the increment of the loss. In experiments we set m = 20 and k = 15.

#### С ADDITIONAL EXPERIMENTS

In this section, we conduct additional experiments under various settings, including ablation studies and more.

## C.1 Ablation study: changing value of $|Z_0|$ and $|Z_1|$

We test the performance of our methods under different sizes of  $|Z_0|$  and  $|Z_1|$  to understand the effect of data set sizes. For settings (a), (c), and (d) in Table 1, we consider decreasing  $(|Z_0|, |Z_1|)$  from (10000, 11000) to (5000, 6000) and increasing  $(|Z_0|, |Z_1|)$  from (10000, 11000) to (15000, 16000), experimenting with both the Shadow Attack and the Outlier Attack. For the setting (e), i.e., the text generation setup, the original  $|Z_0|$  was already small (originally,  $(|Z_0|, |Z_1|) = (4706, 6274))$ , which takes only 30% of the data. Hence, we consider increasing the data size to two different extents. Note that since the Shadow Attack is infeasible in this setting, only the Outlier Attack's results are shown for setting (2). The results for the two attack methods are respectively shown in Tables 5 and 6. We see that the results demonstrate the effectiveness of our methods across different sizes of the dataset. Overall, with a few exceptions, the **Ratio** of the **Compensation Share** appears to become even higher when the dataset size is larger.

Table 5: Results of Shadow Attack for various  $(|Z_0|, |Z_1|)$  settings.

992	Setting	Setting $( Z_0 ,  Z_1 )$		<b>Compensation Share</b>			Fraction of Change		
993 994	Setting	( 20 ,  21 )	Original	Manipulated	Ratio	More	Tied	Fewer	
995		(5000, 6000)	0.0153	0.0764	499.3%	0.996	0.003	0.001	
996	(a)	(10000, 11000)	0.0098	0.0477	456.1%	0.955	0.038	0.007	
997		(15000, 16000)	0.0050	0.0373	746.0%	0.963	0.036	0.001	
998		(5000, 6000)	0.0168	0.0539	320.7%	0.857	0.129	0.014	
999	(c)	(10000, 11000)	0.0112	0.0467	417.0%	0.781	0.195	0.024	
1000		(15000, 16000)	0.0002	0.0062	3100.0%	0.431	0.568	0.001	
1000		(5000, 6000)	0.0206	0.0413	200.7%	0.696	0.174	0.130	
1001	(d)	(10000, 11000)	0.0095	0.0213	217.3%	0.655	0.259	0.086	
1002		(15000, 16000)	0.0057	0.0092	161.5%	0.264	0.616	0.120	

Table 6: Results of Outlier Attack for various  $(|Z_0|, |Z_1|)$  settings.

Setting	$( Z_0   Z_1 )$	Co	ompensation Sha	are	Fract	Fraction of Chang		
Setting	( 20 ,  21 )	Original	Manipulated	Ratio	More	Tied	Fewer	
	(5000, 6000)	0.0153	0.0747	488.9%	0.998	0.002	0.000	
(a)	(10000, 11000)	0.0098	0.0631	643.9%	0.980	0.017	0.003	
	(15000, 16000)	0.0050	0.0400	800.0%	0.964	0.036	0.000	
	(5000, 6000)	0.0168	0.0334	198.8%	0.720	0.241	0.039	
(c)	(10000, 11000)	0.0112	0.0668	596.4%	0.799	0.173	0.028	
	(15000, 16000)	0.0002	0.0051	2550.0%	0.397	0.603	0.000	
	(5000, 6000)	0.0206	0.0411	199.5%	0.761	0.133	0.106	
(d)	(10000, 11000)	0.0095	0.0176	185.2%	0.562	0.354	0.084	
	(15000, 16000)	0.0057	0.0219	384.2%	0.731	0.192	0.077	
	(4706, 6274)	0.0013	0.0035	262.9%	0.392	0.461	0.147	
(e)	(7843, 9411)	0.0016	0.0064	400.0%	0.420	0.507	0.073	
	(12549, 14116)	0.0029	0.0214	737.9%	0.400	0.543	0.057	

### C.2 ABLATION STUDY: DATA ATTRIBUTION METHOD

Next, we conduct an ablation study by varying the data attribution methods used in the evaluation under the experiment setting (b) in Table 1. Specifically, apart from Data Shapley, we consider

1026 two additional data attribution methods, Influence Function, and TRAK, for evaluating the com-1027 pensation share, experimenting with both the Shadow Attack and the Outlier Attack. Note that 1028 the Adversary has no knowledge about what data attribution method will be used by the AI 1029 Developer. The results are shown in Tables 7 and 8. The results below show that the proposed 1030 attack methods are highly effective for all three data attribution methods, as reflected by the Ratio of the Compensation Share. 1031

Table 7: Results of Shadow Attack with different data attribution methods of setting (b).

Attribution Method	$ Z_{1}^{a} / Z_{1} $	Co	mpensation Sha	re	Fract	tion of C	hange
		Original	Manipulated	Ratio	More	Tied	Fewer
Data Shapley	0.0352	0.0152	0.0435	286.2%	0.533	0.333	0.134
Influence Function	0.0352	0.0496	0.1004	202.4%	0.980	0.000	0.020
TRAK	0.0352	0.0392	0.0936	238.8%	0.820	0.100	0.080

Table 8: Results of Outlier Attack with different data attribution methods of setting (b).

Attribution Method	$ Z_{1}^{a} / Z_{1} $	Co	Fraction of Change				
		Original	Manipulated	Ratio	More	Tied	Fewer
Data Shapley	0.0250	0.0112	0.0218	194.6%	0.440	0.380	0.180
Influence Function	0.0250	0.0186	0.0442	237.6%	0.920	0.060	0.020
TRAK	0.0250	0.0160	0.0412	257.5%	0.780	0.180	0.040

1050 1051

1052

1032

1033 1034 1035

1039 1040 1041

1042

#### $C_3$ WHITE-BOX ATTACK

In this section, we further consider *white-box* attacks as an oracle reference, where the Adversary 1053 has full knowledge of the target model and has access to the model parameters. For the choice of 1054 attack method under the white-box threat model, we experiment with the Fast Gradient Sign Method 1055 (FGSM) (Goodfellow et al., 2015) for setting (a) and the Projected Gradient Descent(PGD) (Mkadry 1056 et al., 2017) for setting (d). The results are shown in Table 9. Overall, compared to this oracle white-1057 box attack, the proposed black-box attacks are only slightly worse in terms of the **Ratio** of the 1058 **Compensation Share**. This further confirms that the proposed methods are highly effective.

Table 9: Results of white-box attacks under settings (a) and (d).

Setting	Attack Method	Co	mpensation Sha	Fraction of Change			
Seeing		Original	Manipulated	Ratio	More	Tied	Fewer
(a)	FGSM	0.0098	0.0649	662.2%	0.990	0.007	0.003
(d)	PGD	0.0095	0.0222	233.7%	0.689	0.204	0.107

1067 1068

#### C.4 ADVERSARIAL ATTACK UNDER DATA AUGMENTATION 1069

1070 In this section, we test the performance of our proposed attack methods when the AI Developer 1071 utilize data augmentation techniques when training the model. Specifically, after the AI 1072 Developer gathers data, the training dataset is then formed by a combination of the gathered data and its augmented version. For simplicity, we consider settings (a), (c), and (d) in Table 1 and 1074 consider standard image augmentation methods such as random cropping and random affine trans-1075 formation to images. With data augmentation, the compensation share will be attributed back to the **original** data point if its augmented version is identified as influential. The results are shown in Tables 10 and 11. Compared to the original setting without data augmentation, we cannot draw a 1077 definite conclusion on whether data augmentation helps defend the proposed attacks since the trend 1078 is unclear. However, overall, we can conclude that the proposed attacks are still highly effective 1079 even when data augmentation is utilized by the AI Developer.

Setting	$ Z_{1}^{a} / Z_{1} $	Co	mpensation Sha	re	Fract	tion of C	hange
Setting		Original	Manipulated	Ratio	More	Tied	Fewer
(a)	0.0098	0.0099	0.0209	211.1%	0.632	0.244	0.124
(c)	0.0098	0.0084	0.0516	614.3%	0.894	0.095	0.011
(d)	0.0098	0.0091	0.0747	820.9%	0.959	0.035	0.006

Table 10: Results of Shadow Attack under data augmentation

Table 11: Results of Outlier Attack under data augmentation.

Setting	$ Z_{1}^{a} / Z_{1} $	Co	mpensation Sha	re	Fract	tion of C	hange
Seeing		Original	Manipulated	Ratio	More	Tied	Fewer
(a)	0.0098	0.0099	0.0287	289.9%	0.776	0.158	0.066
(c)	0.0098	0.0084	0.0691	822.6%	0.853	0.113	0.034
(d)	0.0098	0.0091	0.0481	528.6%	0.916	0.065	0.019

1099 C.5 COMPENSATING ONLY CORRECT PREDICTION

1101 In this section we consider the setup when the AI Developer only rewards training data points 1102 that are highly influential to validation samples with *correct* predictions. In detail, after training the 1103 model and before calculating of compensation share, the AI Developer will test the model on a private validation set, select those validation points that are correctly predicted, and only attribute 1104 compensation to influential training points for these validation points. We test this on settings (a), 1105 (c), and (d). Note that the **Fraction of Change** metric is no longer feasible in this setup since the 1106 model trained on the original and manipulated training set may not share the same validation set. 1107 The results are shown in Table 12. Overall, our proposed attacks are still highly effective based on 1108 the Ratio of Compensation Share.

1109

1080

1081

1089

1090 1091

1093 1094 1095

1110 1111 1112

Table 12: Results when compensating only correctly predicted data points.

Attack Type	Setting	$ Z_{1}^{a} / Z_{1} $	Co	mpensation Sha	re
Tituek Type	Setting		Original	Manipulated	Ratio
	(a)	0.0098	0.0050	0.0382	764.0%
Shadow Attack	(c)	0.0098	0.0090	0.0752	835.6%
	(d)	0.0098	0.0110	0.0353	320.9%
	(a)	0.0098	0.0050	0.0273	546.0%
Outlier Attack	(c)	0.0098	0.0090	0.0552	613.3%
	(d)	0.0098	0.0110	0.0227	206.4%

1120 1121 1122

1124

### 1123 C.6 ADVERSARIAL ATTACK ON LARGE-SCALE SETUP

In this section, we scale up our experiment and see how the proposed attacks generalize in this large-scale experiment. Specifically, we consider the ResNet-18 model (He et al., 2016) on the Tiny ImageNet dataset (Le & Yang, 2015) with  $|Z_0|, |Z_1|$ ) = (50000, 60000), and we perturb  $|Z_1^a| = 100$ training points. Both the Shadow Attack and Outlier Attack are tested, and the results are shown in Table 13. It is evident that from the **Ratio** of **Compensation Share**, two proposed attacks are still highly effective in a large-scale experiment.

1130

1131

Table 13	3: Results of S	Shadow Atta	ick and Outlier A	.ttack under	large-sca	lle setup	
Table 13	3: Results of $\frac{ Z_a }{ Z_1 }$	Shadow Atta	ick and Outlier A	ttack under	large-sca Frac	ale setup	Chan
Table 13 Attack Type	3: Results of S $\frac{ Z_1^a / Z_1 }{ Z_1^a }$	Shadow Atta Co Original	ck and Outlier A ompensation Sha Manipulated	ttack under are Ratio	large-sca Fract More	ale setup tion of C Tied	Chan Fe
Table 13 Attack Type Outlier Attack	3: Results of $S$ $ Z_1^a / Z_1 $ 0.0017	Shadow Atta Co Original 0.0009	ck and Outlier A mpensation Sha Manipulated 0.0056	attack under are Ratio 622.2%	large-sca Fract More	tion of C Tied	<b>Chan</b> <b>Fe</b> 0.