

Data Attribution for Diffusion Models: Timestep-induced Bias in Influence Estimation

Tong Xie^{*1}

Haoyu Li^{*1}

Andrew Bai²

Cho-Jui Hsieh²

TONGXIE@UCLA.EDU

HAOYULI02@UCLA.EDU

ANDREWBAI@CS.UCLA.EDU

CHOHSIEH@CS.UCLA.EDU

¹*Department of Mathematics;* ²*Department of Computer Science*
University of California, Los Angeles

Abstract

Data attribution methods trace model behavior back to its training dataset, offering an effective approach to better understand “black-box” neural networks. While prior research has established quantifiable links between model output and training data in diverse settings, interpreting diffusion model outputs in relation to training samples remains under-explored. In particular, diffusion models operate over a sequence of timesteps instead of instantaneous input-output relationships in previous contexts, posing a significant challenge to extend existing frameworks to diffusion models directly. Notably, we present Diffusion-TracIn that incorporates this temporal dynamics and observe that certain timesteps induce a dominating gradient norm, causing a prominent bias in influence estimation that leads specific training samples to emerge as generally influential. To mitigate this effect, we introduce Diffusion-ReTrac as an adaptation that enables the retrieval of training samples more targeted to the test sample of interest, facilitating a localized measurement of influence and considerably more intuitive visualization. We demonstrate the efficacy of our approach through various evaluation metrics and auxiliary tasks.

Keywords: Diffusion Models, Interpretability, Data Attribution, Influence

1. Introduction

Deep neural networks have emerged to be powerful tools for the modeling of complex data distributions and intricate representation learning. However, their astounding performance often comes at the cost of interpretability, leading to an increasing research interest to better explain these “black-box” methods. Instance-based interpretation is one approach to explain why a given machine learning model makes certain predictions by tracing the output back to training samples. While these methods have been widely studied in supervised tasks and demonstrated to perform well [Koh and Liang, 2017; Yeh et al., 2018; Pruthi et al., 2020], there is limited exploration of their application in unsupervised settings, especially for generative models [Kingma and Welling, 2013; Goodfellow et al., 2020; Ho et al., 2020]. In particular, diffusion models represent a state-of-the-art advancement in generative models and demonstrate remarkable performance in a variety of applications such as image generation, audio synthesis, and video generation [Kong et al., 2020; Dhariwal and Nichol, 2021; Ho and Salimans, 2022; Saharia et al., 2022; Hertz et al., 2022; Li

*. Equal contribution

et al., 2022; Ho et al., 2022]. The prevailing generative agents in creative arts such as Stable Diffusion [Rombach et al., 2022] also call for fair attribution methods to acknowledge the training data contributors. Nonetheless, the interpretability and attribution of diffusion models remain an under-explored area [Georgiev et al., 2023; Dai and Gifford, 2023].

Compared to traditional supervised settings, the direct extension of instance-based interpretation to diffusion models is challenging due to the following factors. First, the diffusion objective involves an expectation over the injected noise $\epsilon \sim \mathcal{N}(0, I)$, hence a precise computation is impractical. Second, diffusion models operate over a sequence of timesteps instead of instantaneous input-output relationships. Although each timestep is weighted equally during the training process, we observe that certain timesteps can exhibit the *dominating gradient norm effect*. This means the gradient of the diffusion loss function with respect to model parameters is dominantly large relative to all other timesteps (Figure 2). As most instance-based explanation models utilize this first-order gradient information, such biased gradient norms can propagate its domination into the influence estimation for diffusion models. In practice, specifically, timesteps are often uniformly sampled during training. Therefore a training sample that happens to be trained on certain timesteps may exhibit higher-than-usual gradient norms, and thus be characterized as “generally influential” to various completely different test samples.

We present Diffusion-TracIn and Diffusion-ReTrac to demonstrate and address the existing difficulties. Diffusion-TracIn is a designed extension of TracIn [Pruthi et al., 2020] to diffusion models that incorporates the denoising timestep trajectory. This approach showcases instances where influence estimation is biased. Subsequently, we introduce Diffusion-ReTrac as a re-normalization of Diffusion-TracIn to alleviate the dominating-norm effect.

Our contributions are summarized as follows:

1. Propose Diffusion-TracIn as a designed extension to diffusion models that incorporates the timestep dynamics.
2. Identify and investigate the timestep-induced gradient norm bias in diffusion models, providing preliminary insights into its impact on influence estimation.
3. Introduce Diffusion-ReTrac to mitigate the timestep-induced bias, offering fairer and targeted data attribution.
4. Illustrate and compare the effectiveness of the proposed approach on auxiliary tasks.

2. Related Work

Data attribution methods trace model interpretability back to the training dataset, aiming to answer the following counterfactual question: *which training samples are most responsible for shaping model behavior?*

2.1 Influence Estimations

Influence functions quantify the importance of a training sample by estimating the effect induced when the sample of interest is removed from training [Koh and Liang, 2017]. This method involves inverting the Hessian of loss, which is computationally intensive and can be fragile in highly non-convex deep neural networks [Basu et al., 2020]. Representer Point is

another technique that computes influence using the representer theorem, yet also relies on the assumption that attribution can be approximated by the final layer of neural networks, which may not hold in practice [Yeh et al., 2018]. For diffusion models, the application of influence functions is significantly hindered by its computational expense while extending the representer point method is ambiguous due to the lack of a natural “final layer” in diffusion models. Pruthi et al. introduced TracIn to measure influence based on first-order gradient approximation that does not rely on optimality conditions [Pruthi et al., 2020]. Recently, TRAK (Tracing with the Randomly-projected After Kernel) is introduced as an attribution method for large-scale models, which requires a designed ensemble of models and hence is less suitable for naturally trained models [Park et al., 2023].

In this paper, we extend the TracIn framework to propose an instance-based interpretation method specific to the diffusion model architecture. For a fairer attribution, we present Diffusion-ReTrac that re-normalizes the gradient information to mitigate bias. Previous works have utilized this re-normalization technique to enhance influence estimator performance in supervised settings. Barshan et al. reweight influence function estimations using optimization objectives that place constraints on global influence, enabling the retrieval of explanatory examples more localized to model predictions [Barshan et al., 2020]. Gradient aggregated similarity (GAS) leverages re-normalization to better identify adversarial instances [Hammoudeh and Lowd, 2022]. These works align well with our studies in understanding the localized impact of training instances on model behavior.

2.2 Influence in Unsupervised Settings

The aforementioned methods address the counterfactual question in supervised settings, where model behavior may be characterized in terms of model prediction and accuracy. However, extending this framework to unsupervised settings is non-trivial due to the lack of labels or ground truth. Prior works explore this topic and approach to compute influence for Generative Adversarial Networks (GAN) [Terashita et al., 2021] and Variational Autoencoders (VAE) [Kong and Chaudhuri, 2021]. Previous work in diffusion models quantifies influence through the use of ensembles, which requires training multiple models with subsets of the training dataset, making it unsuitable for naturally trained diffusion models [Dai and Gifford, 2023]. Additionally, the attribution method Journey TRAK applies TRAK [Park et al., 2023] to diffusion models and attributes each denoising timestep individually [Georgiev et al., 2023], which is less interpretable since the diffusion trajectory spans multiple timesteps and a single-shot attribution is more holistic. These works are complementary to our studies and contribute to a more comprehensive understanding of instance-based interpretation methods in unsupervised settings.

2.3 Areas of Application

Data attribution methods prove valuable across a wide range of domains, such as outlier detection, data cleaning, data curating, and memorization analysis [Khanna et al., 2019; Liu et al., 2021; Kong et al., 2021; Lin et al., 2022; Feldman, 2020; van den Burg and Williams, 2021]. The adoption of diffusion models in artistic pursuits, such as Stable Diffusion and its variants, has also gained substantial influence [Rombach et al., 2022; Zhang et al., 2023]. This then calls for fair attribution methods to acknowledge and credit artists whose works

have shaped these models’ training. Such methods also become crucial for conducting further analyses related to legal and privacy concerns [Carlini et al., 2023; Somepalli et al., 2023].

3. Preliminaries

3.1 Diffusion Models

Denoising Diffusion Probabilistic Models (DDPMs) [Ho et al., 2020] are a special type of generative models that parameterized the data as $p_\theta(x_0) = \int p_\theta(x_{0:T}) dx_{1:T}$, where x_1, \dots, x_T are latent variables of the same dimension as the input. The inner term $p_\theta(x_{0:T})$ is the reverse process starting at the standard normal $p(x_T) = \mathcal{N}(x_T; 0, I)$, which is defined by a Markov chain:

$$p_\theta(x_{0:T}) = p(x_T) \prod_{t=1}^T p_\theta(x_{t-1}|x_t), \quad (1)$$

where $p(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t))$. The reverse process is learned to approximate the forward process $q(x_{1:T}|x_0)$, which is fixed to a Markov Chain based on a variance scheduler β_1, \dots, β_T :

$$q_\theta(x_{1:T}|x_0) = \prod_{t=1}^T q_\theta(x_t|x_{t-1}), \quad (2)$$

where $q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t}, \beta_t I)$. Being conditioned on the clean image, one notable property of the forward process is that each x_t at timestep t can be sampled directly from the knowledge of x_0 , independently from the previous timesteps. The distribution of x_t is given as follows

$$q(x_t|x_0) = \mathcal{N}(x_t; \sqrt{\alpha_t}x_0, (1 - \alpha_t)I), \quad (3)$$

where $\alpha_t = \prod_{i=1}^t (1 - \beta_i)$. Therefore, efficient training can be achieved by stochastically selecting timesteps for each sample. DDPM further simplifies the loss by re-weighting each timestep, leading to the training objective used in practice,

$$L_{\text{simple}}(\theta) = \mathbb{E}_{x_0, t, \epsilon} [d(\epsilon, \epsilon_\theta(\sqrt{\alpha_t}x_0 + \sqrt{(1 - \alpha_t)}\epsilon, t))], \quad (4)$$

where $\epsilon \sim \mathcal{N}(0, I)$ and d is the loss function such as l_1 or l_2 distance.

3.2 TracIn

TracIn is proposed as an efficient first-order approximation of a training sample’s influence. It defines the idealized version of influence of a training sample z to a test sample z' as the total reduction of loss on z' when the model is trained on z . Formally, let k denote each occurrence of model update using the training sample z , and w_k denote the model parameters after the update. Then the idealized influence is calculated as

$$\text{Ideal-Influence}(z, z') = \sum_{k: z_k = z} \ell(w_k, z') - \ell(w_{k+1}, z'). \quad (5)$$

For tractable computation, the change in loss on the test sample is approximated by the following Taylor expansion centered at w_k

$$\ell(w_{k+1}, z') - \ell(w_k, z') = \nabla \ell(w_k, z') \cdot (w_{k+1} - w_k) + \mathcal{O}(\|w_{k+1} - w_k\|^2). \quad (6)$$

If Stochastic Gradient Descent (SGD) is utilized in model training, then the model parameter update is measured by $w_{k+1} - w_k = -\eta_t \nabla \ell(w_k, z_k)$. Therefore, the first-order approximation of Ideal-Influence 5 is derived to be

$$\text{TracIn}(z, z') = \sum_{k: z_k = z} \eta_k \nabla \ell(w_k, z') \cdot \nabla \ell(w_k, z). \quad (7)$$

To reduce computational costs, TracIn approximates the influence with saved checkpoints from training to replay the entire training process. It is also expected that the practical form of TracIn remains the same across variations in training, such as optimizers, learning rate schedules, and handling of minibatches [Pruthi et al., 2020]. A training sample that exerts positive influence over the test sample is called the proponent, and opponent otherwise.

4. Method

4.1 Diffusion-TracIn

In this section, we present Diffusion-TracIn to provide an efficient extension of TracIn designed for diffusion models. Specifically, two adjustments keen to diffusion models are critical to enable this extension. First, the diffusion objective is an expectation of denoising losses over different timesteps $1 \leq t \leq T$. Second, the objective involves an expectation over the added noise $\epsilon \sim \mathcal{N}(0, I)$. To address these challenges, we first apply TracIn conditioned on each timestep t , then we compute a Monte Carlo average over m randomly sampled noises ϵ .

From equation 4, it is possible to treat the diffusion model learning objective as a combination of T loss functions. If we denote $L_t(\theta, \epsilon, x_0) = L_{\text{simple}}(\theta, \epsilon, x_0, t) := d(\epsilon, \epsilon_\theta(\sqrt{\alpha_t}x_0 + \sqrt{(1 - \alpha_t)}\epsilon), t)$ to be a distinct loss function on each timestep t and noise, we can treat L_{simple} as an expectation over all the L_t . Subsequently, we compute a TracIn influence score over each of the timestep t

$$\begin{aligned} \text{TracIn}(z, z', t) &:= \mathbb{E}_\epsilon \left(\sum_{k: z_k = z} \eta_k \nabla_\theta L_t(\theta_k, \epsilon, z') \cdot L \right) \\ &\approx \frac{1}{m} \sum_{i=1}^m \sum_{k: z_k = z} \eta_k \nabla_\theta L_t(\theta_k, \epsilon_i, z') \cdot L, \end{aligned} \quad (8)$$

where $L = \nabla_\theta L_{t_{\text{train}}}(\theta_k, \epsilon_{\text{train}}, z)$, t_{train} is the training timesteps, and ϵ_{train} is the noise utilized when training on the sample z . This enables a true replay of the training dynamics. Then we define Diffusion-TracIn to be the expectation over T timesteps to provide an

one-shot attribution score that covers the full diffusion process,

$$\begin{aligned}
\text{Diffusion-TracIn}(z, z') &:= \mathbb{E}_t(\text{TracIn}(z, z', t)) \\
&= \frac{1}{T} \sum_{t=1}^T \text{TracIn}(z, z', t) \\
&= \frac{1}{T} \sum_{t=1}^T \frac{1}{m} \sum_{i=1}^m \sum_{k: z_k=z} \eta_k \nabla_{\theta} L_t(\theta_k, \epsilon_i, z') \cdot L \\
&= \sum_{k: z_k=z} \eta_k \left(\frac{1}{mT} \sum_{t=1}^T \sum_{i=1}^m \nabla_{\theta} L_t(\theta_k, \epsilon_i, z') \right) \cdot L.
\end{aligned} \tag{9}$$

The practical form of Diffusion-TracIn also employs training checkpoints, as suggested by [Pruthi et al., 2020] to enhance computational efficiency,

$$\text{Diffusion-TracIn}(z, z') := \sum_{k=1}^s \eta_k \left(\frac{1}{mT} \sum_{t=1}^T \sum_{i=1}^m \nabla_{\theta} L_t(\theta_k, \epsilon_i, z') \right) \cdot L, \tag{10}$$

where s is the number of checkpoints. Since the number of timesteps T can be large (e.g. 1000) in practice, we leverage evenly spaced timesteps ranging from 1 to T for the test sample z' . It is empirically observed that this yields an effective approximation to the full diffusion process, and provides similar attribution results.

4.2 Diffusion-ReTrac

Although Diffusion-TracIn is a direct extension of TracIn for diffusion models derived from the mathematical definition, we identify the *dominating loss gradient norm effect* which can lead to bias in influence estimation.

Intuition. By Cauchy-Schwarz inequality, it can be noticed from Equation 7 that

$$|\text{TracIn}(z, z')| \leq \sum \eta_t \|\nabla \ell(w_t, z')\| \|\nabla \ell(w_t, z)\|.$$

Hence, training samples with disproportionately large gradient norms tend to have significantly higher influence score $|\text{TracIn}(z, z')|$. This suggests that these samples are more likely to be characterized as either a strong proponent or opponent to the given test sample z' , depending on the direction alignment of $\nabla l(w_t, z')$ and $\nabla l(w_t, z)$.

Motivation. In most machine learning models, a dominating gradient norm can be largely attributed to the training sample itself. For example, outliers and samples near the decision boundary may exhibit higher gradient norms than usual. However, while sample-induced variance in gradient norms is informative for influence estimation, we demonstrate that the variance in gradient norms for diffusion models can also be an artifact of the diffusion training dynamics. In particular, empirical results show that the loss function component L_t from certain timesteps is more likely to have a larger gradient norm as shown in Figure 2.

Definition 1 We define the *highest norm inducing timestep* for sample z to be

$$t_{\max}(z) = \arg \max_i \|L_{t_i}(\epsilon_{\text{train}}, z)\|.$$

In other words, if during the training process, the stochastically chosen timestep for a training sample z happens to fall closer to $t_{\max}(z)$, then z will exhibit a biased larger norm that propagates into the influence calculation. Since the natural training of diffusion models sample timesteps randomly for each z , different degrees of “timestep-induced” norm biases are introduced, leading to unfair influence estimates across training samples.

Approach. Since the gradient norm is not solely a property attributed to the sample but rather also caused by the norm bias inherent to diffusion models, an ideal instance-based interpretation should not overestimate the influence of samples with large norms and penalize those with small norms.

In fact, this dominating norm effect can be introduced by the choice of timesteps for both test sample z' and each training sample z , whose loss gradient norms are computed using timestep t and t_{train} respectively. For test sample z' , the gradient information $\sum_{t=1}^T \sum_{i=1}^m \nabla_{\theta} L_t(\theta_k, \epsilon_i, z')$ derives from an expectation over all timesteps $t \in [1, T]$. Therefore, influence estimation inherently upweights timesteps with larger norms and downweights those with smaller norms. For each training sample z , the timestep t_{train} was stochastically sampled during the training process, hence incorporating varying degrees of timestep-induced norm bias. To this end, we propose Diffusion-ReTrac which introduces normalization that reweights the training samples to address the dominating-norm effect. We normalize these two terms and define

$$\text{Diffusion-ReTrac}(z, z') = \sum_{k: z_k = z} \eta_k \left(\frac{1}{mT} \sum_{t=1}^T \sum_{i=1}^m \frac{\nabla_{\theta} L_t(\theta_k, z', \epsilon_i)}{\|\nabla_{\theta} L_t(\theta_k, z', \epsilon_i)\|} \right) \cdot \frac{L}{\|L\|} \quad (11)$$

The bias introduced to influence estimation due to timestep-induced norms is thus mitigated. In this way, we minimize the vulnerability that the calculated influences are dominated by training samples with a disproportionately large gradient norm arising from stochastic training.

5. Experiments

To illustrate our observation, we provide evidence showcasing the dominating-norm bias in influence estimation. We further present instances where this effect may be unnoticed in common benchmarks, and evaluate the performance of Diffusion-TracIn and ReTrac. Our discussion addresses the following questions:

1. **Timestep-induced Bias:** How does timestep affect the influence estimation?
2. **Outlier Detection:** Why might the timestep-induced bias be unnoticed in detecting outliers or atypical samples by calculating *self-influence*?
3. **Image Tracing:** How effective is each method at attributing the learning source of an image to the training data through *test-influence*?
4. **Targeted Attribution:** How does Diffusion-ReTrac outperform Diffusion-TracIn by addressing the norm bias?

5.1 Timestep-induced gradient norm bias

The *dominating gradient norm effect* refers to when influence estimation is biased by the sample’s loss gradient norm arising from diffusion timesteps. Since the training timestep for each instance is stochastically sampled, every instance receives a varying degree of such timestep-induced bias which propagates into the influence calculation. Such bias is particularly evident in samples whose training timestep falls close to t_{\max} (Definition 1), leading them to be characterized as “generally influential” (Figure 6). This suggests that samples’ norms may be a suboptimal source of information, because it can be induced by diffusion timesteps rather than fully attributed to the sample itself.

Using a diffusion model trained on Artbench-2 with 10,000 samples as illustration, we demonstrate the presence of such timestep-induced norm by showing:

1. Notable trends and statically significant correlation between training samples’ gradient norm and their training timestep.
2. For each individual sample, its gradient norm is highly dependent on the timestep.
3. Susceptibility of influence estimation: significant alteration in a training sample’s influence score via timestep manipulation.

Norm vs. Timestep. We examine the distribution of training samples’ loss gradient norm and the training timestep (Figure 2). This is conducted for multiple checkpoints, each yielding a distribution displaying norms at its specific stage of the training process. The distributions demonstrate a notable upward trend that peaks at, in this case, the later timestep region (i.e. timesteps closer to noise). This suggests that samples whose training timestep falls within the later range tend to exhibit higher norms.

We further examine the impact and quantify the relationship between loss gradient norms and training timesteps. Using randomly selected samples $\{x_i\}_{i=1}^{50}$, we measure the correlation between the proximity of sample x_i ’s training timestep to $t_{\max}(x_i)$ and ranking of x_i ’s norm among the entire training dataset. The detailed procedure is included in Algorithm 1. This is conducted for every checkpoint used to compute influence. As an illustration, Figure 1 shows a checkpoint with a correlation of 0.7 and p -value 1.38×10^{-7} . A linear regressor is also fitted to the 50 data points, giving a slope of 6.038 and p -value 2.55×10^{-8} . This suggests a statistically significant positive correlation between the norms and training timesteps. Consequently, it indicates a notable training timestep-induced norm bias that could well dominate over sample-induced norms, which will then propagate into influence estimation.

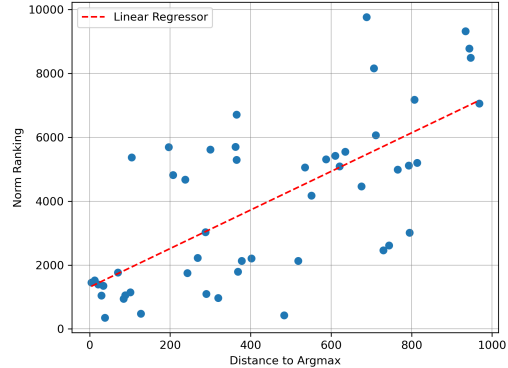


Figure 1: **Correlation of Training Timestep and Norm Ranking.** We plot the norm ranking and distance between training timestep to t_{\max} for 50 randomly selected samples. The resulting correlation is 0.7 and the linear regressor (red) has a slope of 6.038.

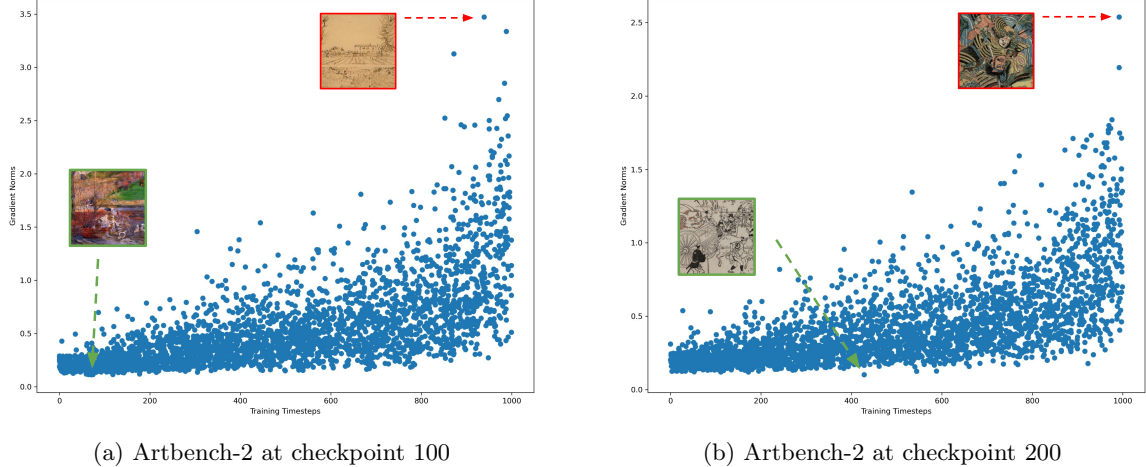


Figure 2: Samples’ Norm vs. Training Timestep. We plot the norm and timestep of 2,000 randomly selected training samples. We observe that loss gradient norms tend to increase when the training timestep falls in the later range (towards noise). This upward trend is consistent at other checkpoints tested. The sample with the largest norm (red) and smallest norm (green) are shown; no exceptional visual patterns are noticed.

Varying Timestep for a Single Sample. We further analyze the norm distribution for an individual training sample. At a given model checkpoint, we compute the loss gradient norm for a fixed sample x at every timestep. We plot the norm distribution for randomly selected samples and observe similar trends. Figure 3 shows the distribution for an example training data at three different model stages. The highest norm-inducing region at these three checkpoints all falls within the later timestep range, regardless of where the training timestep is sampled at. The further implication is that for each sample, its loss gradient norm is highly dependent on the chosen timestep. It is also observed that within the same epoch, various samples share similar trends in norm distribution. This suggests a systematic pattern (e.g. artifact of diffusion learning dynamics) beyond individual instances, supporting the intuition that over-reliance on gradient norms may not be ideal.

Timestep Manipulation. We further illustrate the timestep-induced bias by exploring the susceptibility of influence estimation to the manipulation of timesteps. We conduct the experiment on 500 training samples that are characterized by Diffusion-TracIn as “un-influential” to a random test sample (i.e. influence score is close to 0, neither proponent nor opponent). For each un-influential sample x_i , we compute its influence using $t_{\max}(x_i)$ instead of the original training timestep. The result shows that after deliberately modifying timestep, the ranking of the magnitude of influence for these samples increases by 4,287 positions on average. Given that there are only 10,000 images in the dataset, this notable fluctuation indicates that the timestep-induced bias is significant enough to flip a training sample from un-influential to proponents or opponents. It further indicates that Diffusion-TracIn’s attribution results could well arise from timestep-induced norms in general. Such findings highlight the potential vulnerability within Diffusion-TracIn as a data attribution method, emphasizing the need for more robust influence estimation techniques.

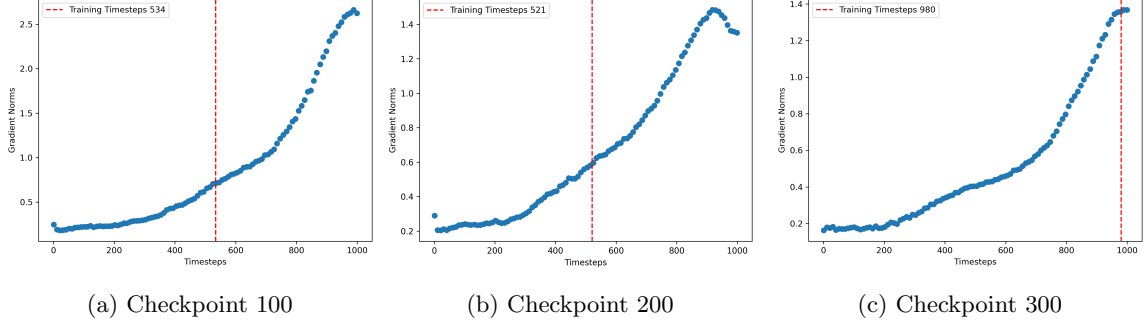


Figure 3: **One Sample’s Norms Varying Timestep.** Example of norm distributions for one randomly selected sample is shown. On each checkpoint, we observe that the norm distribution is skewed to later timesteps for each individual sample.

5.2 Outlier Detection

Influence estimation is often used to identify outliers that deviate notably from the rest of the training data. Intuitively, outliers independently support the model’s learning at those sparser regions of the input space to which they belong, whereas learning of the typical samples is supported by a wide range of data. Hence in an ideal influence method, outliers tend to exhibit high self-influence, indicating that they exert a high contribution in reducing their own loss. Because of such an outlier-induced norm, we observe that biased estimations may easily go unnoticed in this common metric of outlier detection.

Setup. We begin by training a diffusion model on a combination of the entire 5,000 samples from CIFAR-10 airplane subclass and 200 samples from MNIST zero subclass. Subsequently, we compute the self-influence of each of the training instances, and sort them by descending order. Since the 200 MNIST samples are outliers that independently support a region, we evaluate whether our methods assign high self-influence to the 200 samples of MNIST zero.

Result. The results show that both Diffusion-TracIn and Diffusion-ReTrac successfully rank outlier samples with high self-influence (Table 1). However, the bias introduced by diffusion timesteps is unnoticed in this experiment. Since outliers naturally exhibit larger norms compared to the typical inliers, the timestep-induced norm becomes a more obscure confounding factor and hence is less subtle in the computation of self-influence.

	Top 100	Top 200	Top 300
Diffusion-TracIn	0.880	0.880	1.000
Diffusion-ReTrac	0.860	0.845	1.000

Table 1: **Outlier Detection.** This table measures the proportion of MNIST samples among top- k identified samples with the highest self-influence. We observe that the performance of these two methods is on par with each other. Both methods assign high self-influence to the 200 MNIST outliers out of the 5,000 CIFAR planes.

Visualization. Examining high-ranking samples further shows that the non-outlier airplane samples with high self-influence (among the top 200) are images with large contrast

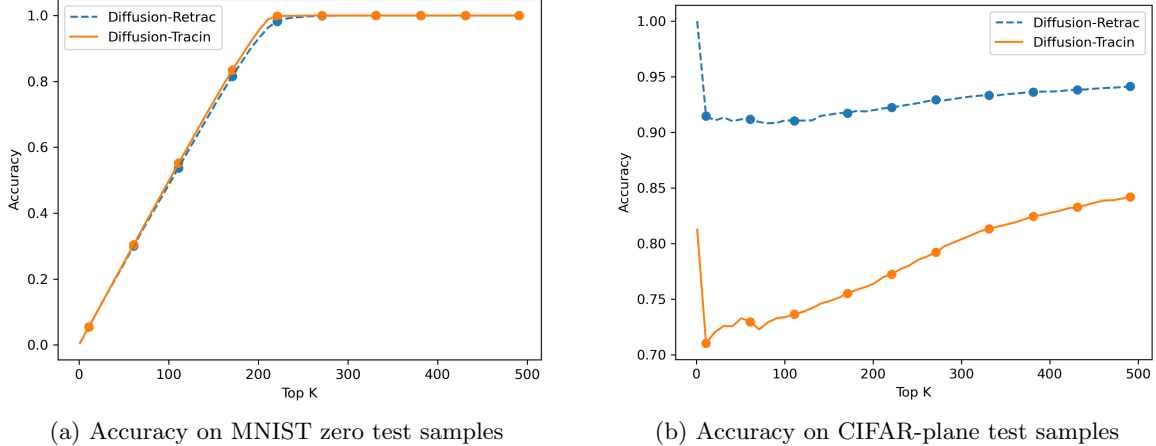


Figure 4: **Image Tracing Accuracy.** We evaluate the proportion of correctly attributed training samples among the top- k influential samples identified by the two methods. While (a) shows both methods successfully attribute MNIST zero test samples, (b) shows that Diffusion-TracIn fails to attribute CIFAR-plane test samples since the outlier MNIST samples with large norms received biased estimations.

and atypical backgrounds compared to airplane samples with low self-influence. Overall, samples with high self-influence tend to exhibit high visual contrast or are difficult to recognize. This observation is consistent with patterns revealed in previous work on influence estimation for VAE [Kong and Chaudhuri, 2021]. Visualization for plane samples with high self-influence is included in Appendix C.1.

5.3 Image Tracing

One fundamental role of data attribution methods is to trace the model’s outputs back to their origins in the training samples. This idea is also utilized for analyzing *memorization* [Feldman, 2020], a behavior where the generated sample is attributed to a few nearly identical training samples. In essence, *Image source tracing* helps pinpoint specific training samples that are responsible for a generation. Thus we evaluate our methods on the question: Given a test sample, which instances in the training dataset is the model’s knowledge of the test sample derived from?

Setup. We extend our experiment using the aforementioned model trained on CIFAR-10 airplane and MNIST zero subclass [Krizhevsky et al., 2009; LeCun and Cortes, 2010]. Given a test sample of MNIST zero, it is expected that the 200 MNIST samples in the training dataset serve as ground truth for the image source. Similarly, a test sample of CIFAR-plane should be attributed to the 5,000 CIFAR training samples. We thus obtain an accuracy score by measuring the correctly attributed proportion among the top- k influential sample.

Results. The top- k accuracy scores for Diffusion-TracIn and ReTrac are reported in Figure 4. While both methods successfully attribute the MNIST test samples to the 200 MNIST training samples, we note that Diffusion-TracIn is also more likely to attribute MNIST training samples to a CIFAR-plane test sample (Figure 5). This aligns with the expectation that Diffusion-TracIn tends to assign higher influence to training samples with large norms,

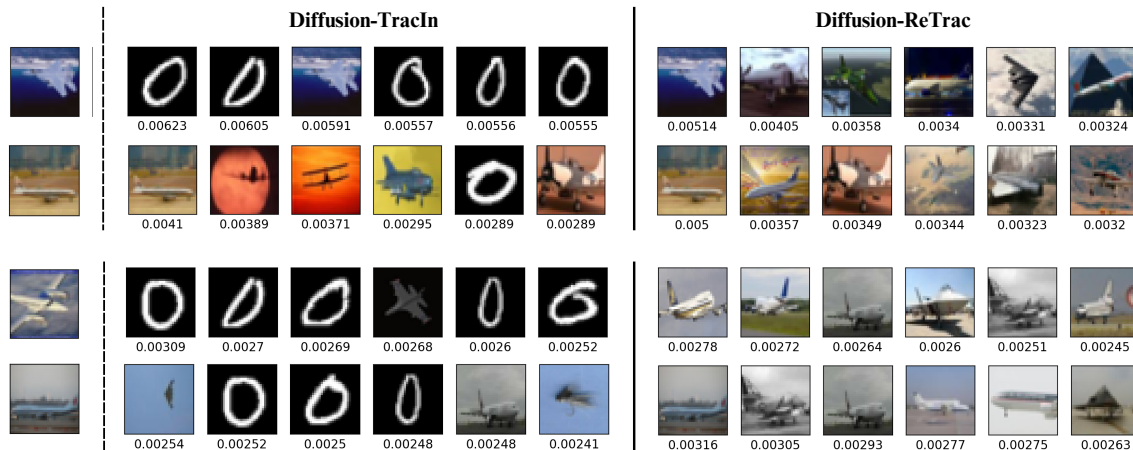


Figure 5: **Image Tracing on Outlier Model.** We evaluate the two attribution methods on both the training samples (top 2 rows) and generated samples (bottom 2 rows). While both methods successfully attribute MNIST zero test samples, Diffusion-TracIn also incorrectly characterizes MNIST training samples as influential for a CIFAR-plane test sample. This is notably mitigated with Diffusion-ReTrac.

in this case, the outlier MNIST zeros. Further analysis of these influential zeros indicates that their training timestep is sampled exactly at or close to the region that induces a large norm. This then becomes the confounding factor that further amplifies these outliers’ norms, exacerbating the bias introduced when attributing the test samples. On the other hand, Diffusion-ReTrac successfully attributes both MNIST-zero and CIFAR-plane test samples.

Visualization. Upon closer examining Diffusion-TracIn attribution results, the set of MNIST zero samples exerting influence on plane is relatively consistent across various CIFAR plane test samples. For instance, the MNIST sample with highest influence to the two generated planes are identical in Figure 5. For the CIFAR-plane samples that Diffusion-TracIn successfully attributes (without influential MNIST zeros), there still appear to be generally influential planes. This phenomenon is alleviated for samples retrieved using ReTrac, with sets of influential samples being more distinct and visually intuitive. Additionally, the CIFAR-plane instances with high influence (e.g. among the top 200) to MNIST test samples tend to be planes with black backgrounds, which to an extent also resemble the MNIST zero. Visualization for these proponent planes is included in Appendix C.2. It is also worth highlighting that Diffusion-ReTrac identifies potentially memorized samples for the generated image, such as the last row in Figure 5.

5.4 Targeted Attribution

We then provide a comprehensive analysis of the influential samples retrieved by Diffusion-TracIn and ReTrac. In this experiment, we compute the influence over Artbench-2 and CIFAR-10 datasets. Compared to previous settings, this experiment minimizes the effects of unusually large “sample-induced” gradient norms due to the deliberately introduced outliers. This experiment further compares the capability of Diffusion-TracIn and ReTrac in tasks with different emphases or objectives.

		Top 10	Top 50	Top 100
ArtBench-2	D-TracIn	0.293	0.261	0.248
	D-ReTrac	0.812	0.646	0.605
CIFAR-10	D-TracIn	0.725	0.663	0.636
	D-ReTrac	0.856	0.800	0.768

Table 2: **Targeted Attribution.** This table shows the average proportion of unique samples retrieved over multiple test samples. It can be noticed that Diffusion-TracIn overall extracts far less unique samples compared to Diffusion-ReTrac, especially on the two-class ArtBench.

Setup. We compute test-influence on two diffusion models trained with datasets 1). Artbench-2 consisting of “Post-impressionism” and “ukiyo-e” subclasses from Artbench [Liao et al., 2022], each containing 5,000 training samples with resolution 64×64 , and 2). CIFAR-10 [Krizhevsky et al., 2009] consisting of 50,000 training samples with resolution 32×32 .

Results. To quantify the targeted-ness of the data attribution method, we assess the prevalence of generally influential samples. For a given k , we measure the proportion of distinct samples among the top- k influential samples identified by the two methods. The results are shown in Table 2. We note that Diffusion-TracIn yields extremely homogenous influential samples. This trend is particularly evident in ArtBench-2, where the two-class setting is less diverse and more prone to bias induced. In this case, Diffusion-TracIn incurs an $\frac{2}{3}$ of overlaps within the top 10 influential samples. This observation aligns with our argument that the stochastically chosen timesteps have amplified the number of samples that exhibit larger gradient norms, therefore causing more generally influential samples.

Visualization. From Figure 6, it is visually evident that Diffusion-TracIn retrieves numerous generally influential training samples. The same or similar sample can be attributed to test samples that are completely different (e.g. in terms of subclass or visual similarities such as color and structure).

Further analysis of these generally influential samples suggests that their associated timestep tends to be close to t_{\max} . As an illustration, we show distribution of norm vs. timestep for an example generally influential sample in Figure 7. This specific sample emerges as influential at checkpoint 80, which coincides with its training timestep falling within the T_{\max} region. It then becomes generally influential as shown in Figure 6a (first proponent on the last row). This phenomenon is notably mitigated after normalization in Diffusion-ReTrac. The revised approach retrieves influential training samples that bear greater visual resemblance to the test samples, highlighting ReTrac’s targeted attribute and reinforcing that dissimilar test samples are more likely to be influenced by a distinct set of training samples.

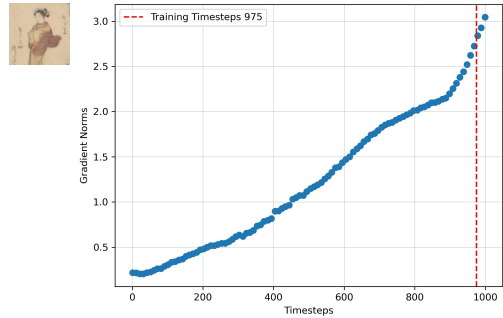
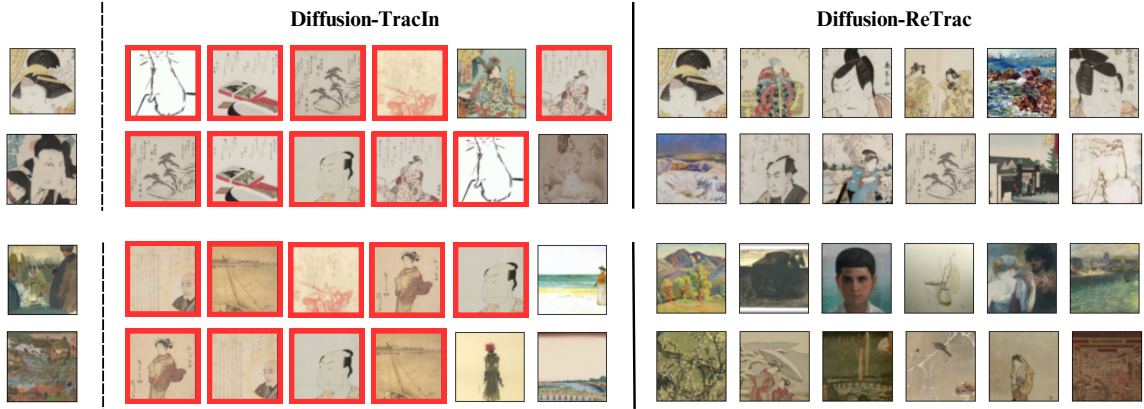
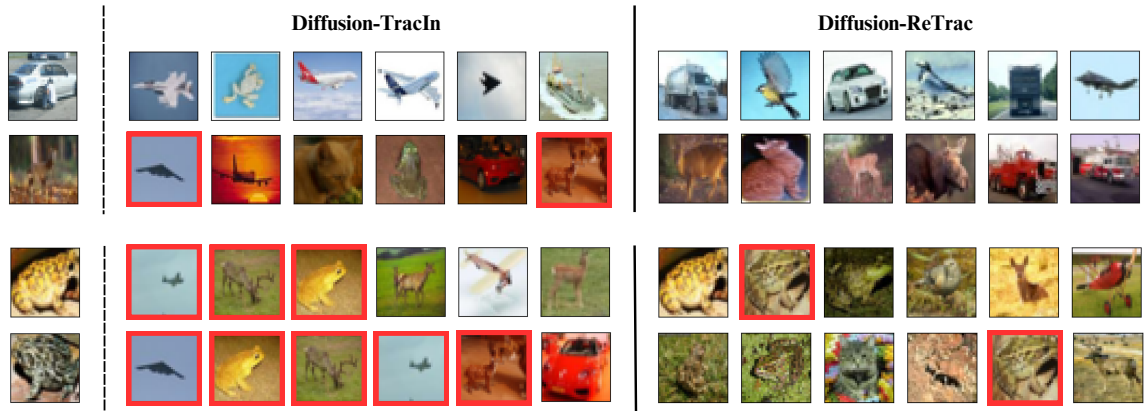


Figure 7: **Timestep and Generally Influential.** Example of a generally influential sample (left) whose training timestep falls exactly within the $t_{\max}(x)$ region.



(a) **Artbench-2**. Top 2 test samples are of the class “Ukiyo-e” and the bottom 2 are of “Post Impressionism”.



(b) **CIFAR-10**. Top 2 test samples are of the subclass “automobile,” “deer” and bottom 2 are of “frog”.

Figure 6: **Targeted Attribution**. We attribute the 8 test samples (leftmost) using both Diffusion-TracIn and Diffusion-ReTrac. The top 6 proponents are shown. The generally influential images that appear multiple times for different test samples are indicated in red. It is visually evident that Diffusion-ReTrac provides more distinct attribution results.

6. Conclusion

In this work, we extend data attribution framework to diffusion models and identify a prominent bias in influence estimation originating from loss gradient norms. Our detailed analysis elucidates how this bias propagates into the attribution process, revealing that gradient information harbors undesired bias caused by diffusion model dynamics. Subsequent experiments validate Diffusion-ReTrac as an effective attempt to mitigate this effect, offering fairer and targeted attribution results.

Limitations and future work. A theoretical explanation for the large-norm-inducing timesteps better pinpoints the causes and provides ad hoc solutions for the problem. While renormalization mitigates the dominating norm effect and “generally influential” samples, further examination of the gradient alignments may also be beneficial. Analysis of other potential confounding factors gives further insights into a fair attribution method.

References

- Elnaz Barshan, Marc-Etienne Brunet, and Gintare Karolina Dziugaite. Relatif: Identifying explanatory training samples via relative influence. In *International Conference on Artificial Intelligence and Statistics*, pages 1899–1909. PMLR, 2020.
- Samyadeep Basu, Philip Pope, and Soheil Feizi. Influence functions in deep learning are fragile. *arXiv preprint arXiv:2006.14651*, 2020.
- Nicolas Carlini, Jamie Hayes, Milad Nasr, Matthew Jagielski, Vikash Sehwal, Florian Tramèr, Borja Balle, Daphne Ippolito, and Eric Wallace. Extracting training data from diffusion models. In *32nd USENIX Security Symposium (USENIX Security 23)*, pages 5253–5270, 2023.
- Zheng Dai and David K Gifford. Training data attribution for diffusion models. *arXiv preprint arXiv:2306.02174*, 2023.
- Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021.
- Vitaly Feldman. Does learning require memorization? a short tale about a long tail. In *Proceedings of the 52nd Annual ACM SIGACT Symposium on Theory of Computing*, pages 954–959, 2020.
- Kristian Georgiev, Joshua Vendrow, Hadi Salman, Sung Min Park, and Aleksander Madry. The journey, not the destination: How data guides diffusion models. 2023.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020.
- Zayd Hammoudeh and Daniel Lowd. Identifying a training-set attack’s target using renormalized influence estimation. In *Proceedings of the 2022 ACM SIGSAC Conference on Computer and Communications Security*, pages 1367–1381, 2022.
- Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626*, 2022.
- Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022.

- Rajiv Khanna, Been Kim, Joydeep Ghosh, and Sanmi Koyejo. Interpreting black box predictions using fisher kernels. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 3382–3390. PMLR, 2019.
- Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In *International conference on machine learning*, pages 1885–1894. PMLR, 2017.
- Shuming Kong, Yanyan Shen, and Linpeng Huang. Resolving training biases via influence-based data relabeling. In *International Conference on Learning Representations*, 2021.
- Zhifeng Kong and Kamalika Chaudhuri. Understanding instance-based interpretability of variational auto-encoders. *Advances in Neural Information Processing Systems*, 34:2400–2412, 2021.
- Zhifeng Kong, Wei Ping, Jiaji Huang, Kexin Zhao, and Bryan Catanzaro. Diffwave: A versatile diffusion model for audio synthesis. *arXiv preprint arXiv:2009.09761*, 2020.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- Yann LeCun and Corinna Cortes. MNIST handwritten digit database. 2010. URL <http://yann.lecun.com/exdb/mnist/>.
- Xiang Li, John Thickstun, Ishaan Gulrajani, Percy S Liang, and Tatsunori B Hashimoto. Diffusion-lm improves controllable text generation. *Advances in Neural Information Processing Systems*, 35:4328–4343, 2022.
- Peiyuan Liao, Xiuyu Li, Xihui Liu, and Kurt Keutzer. The artbench dataset: Benchmarking generative models with artworks. *arXiv preprint arXiv:2206.11404*, 2022.
- Jinkun Lin, Anqi Zhang, Mathias Lécuyer, Jinyang Li, Aurojit Panda, and Siddhartha Sen. Measuring the effect of training data on deep learning predictions via randomized experiments. In *International Conference on Machine Learning*, pages 13468–13504. PMLR, 2022.
- Zhuoming Liu, Hao Ding, Huaping Zhong, Weijia Li, Jifeng Dai, and Conghui He. Influence selection for active learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9274–9283, 2021.
- Sung Min Park, Kristian Georgiev, Andrew Ilyas, Guillaume Leclerc, and Aleksander Madry. Trak: Attributing model behavior at scale. *arXiv preprint arXiv:2303.14186*, 2023.
- Garima Pruthi, Frederick Liu, Satyen Kale, and Mukund Sundararajan. Estimating training data influence by tracing gradient descent. *Advances in Neural Information Processing Systems*, 33:19920–19930, 2020.

- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022.
- Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan Ho, Tim Salimans, David Fleet, and Mohammad Norouzi. Palette: Image-to-image diffusion models. In *ACM SIGGRAPH 2022 Conference Proceedings*, pages 1–10, 2022.
- Gowthami Somepalli, Vasu Singla, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Diffusion art or digital forgery? investigating data replication in diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6048–6058, 2023.
- Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020.
- Naoyuki Terashita, Hiroki Ohashi, Yuichi Nonaka, and Takashi Kanemaru. Influence estimation for generative adversarial networks. *arXiv preprint arXiv:2101.08367*, 2021.
- Gerrit van den Burg and Chris Williams. On memorization in probabilistic deep generative models. *Advances in Neural Information Processing Systems*, 34:27916–27928, 2021.
- Chih-Kuan Yeh, Joon Kim, Ian En-Hsu Yen, and Pradeep K Ravikumar. Representer point selection for explaining deep neural networks. *Advances in neural information processing systems*, 31, 2018.
- Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3836–3847, 2023.

Appendix

A Timestep-Induced Bias

A.1 NORM VS. TIMESTEP

To demonstrate that diffusion timesteps have a significant impact on loss gradient norms, we plot the distribution of 2,000 randomly selected training samples' norms and their training timesteps. Visualization for the distribution is shown in Figure 8. There is a notable upward trend that peaks at the later range of the timesteps (i.e. timesteps closer to noise), suggesting that samples trained during these later timesteps tend to exhibit larger norms. Additionally, it is also observed that the trend in norm distribution gradually diminishes at the model convergence. This further supports that such variance due to timestep is an artifact of the training dynamic, rather than a property of the training sample. However, Diffusion-TracIn utilizes gradient information throughout the entire learning process instead of focusing solely on those near convergence. This approach is due to the tendency of the latter to contain minimal information, resulting in an inevitable trend in norms affecting influence estimation. This also motivates the renormalization technique in Diffusion-ReTrac.

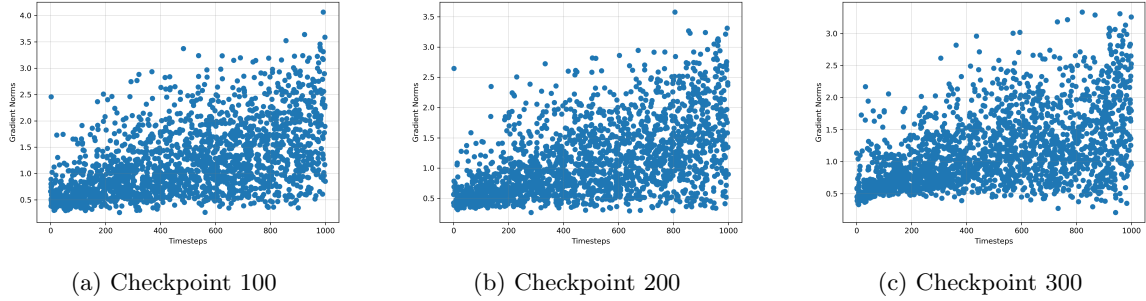


Figure 8: **Norm Distribution.** We plot the loss gradient norm and training timestep of 2,000 samples. The distributions at checkpoints 100, 200, and 300 all demonstrate an upward trend. This suggests that samples whose training timesteps fall within the later timestep region tend to have a larger norm.

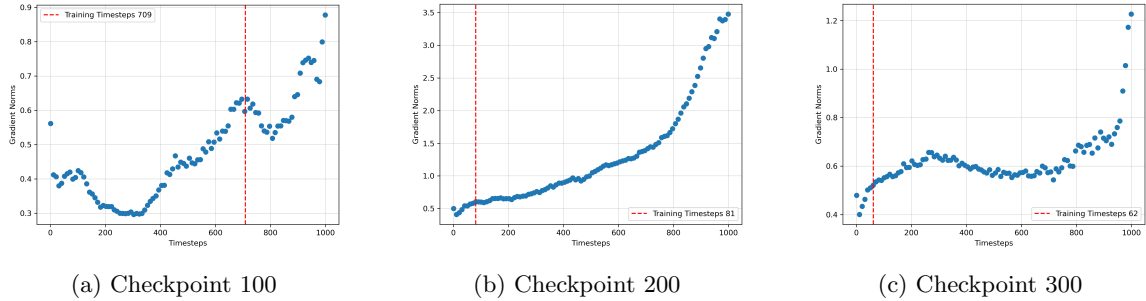


Figure 9: **Varying Timestep for a Single Sample.** The norms of sample #0 are computed at different timesteps. The distribution obtained at checkpoints 100, 200, and 300 all demonstrates a similar trend that peaks at the later timestep region.

Algorithm 1 Training Timesteps and Norm Ranking

```
 $y_1 \leftarrow \{\}$   
 $y_2 \leftarrow \{\}$   
for  $k \leftarrow 0$  to  $N$  do:  
   $x \leftarrow k$ -th training sample  
   $y_1 \leftarrow y_1 \cup |t_{\max}(x) - t_{\text{train}}(x)|$   
   $y_2 \leftarrow y_2 \cup \text{Rank}(x)$   
end for  
return spearman-rank( $y_1, y_2$ )
```

A.2 VARYING TIMESTEP FOR A SINGLE SAMPLE

We further show that for a fixed training sample, the gradient norm with respect to its loss computed at different timesteps varies significantly. This reinforces the effect of training timestep in the estimation of influence, indicating that each sample receives a varying degree of bias since the training timestep is stochastically sampled. An example norm distribution for a fixed sample at different checkpoints is shown in Figure 9.

A.3 CORRELATION

We provided quantitative analysis addressing the question: If the training timestep of a sample x falls closer to $t_{\max}(x)$, does x also have a relatively larger norm compared to the rest of the training dataset? To analyze the relationship between the stochastically chosen training timestep and the sample’s overall norm ranking among the rest, we obtain a correlation score by i). compute the distance between a sample’s training timesteps t_{train} and the timestep that yields the maximum norm $t_{\max}(x)$, ii). the ranking of this sample’s gradient norm among all the training samples, and iii). calculate a Spearman Rank correlation score between distance and ranking (Algorithm 1). Figure1 in the main text shows a visualization of the measured correlation.

B Generally Influential Samples

As additional motivation for renormalization, we observe that Diffusion-TracIn assigns dominantly high influence to samples with a large norm, even at a single checkpoint. Such a large norm is often associated with training timesteps close to the t_{\max} region, signifying a strong timestep-induced bias in the loss gradient norm. Furthermore, these particular samples only emerge as influential when such checkpoints are utilized, and are likely to persist as strong proponents or opponents throughout the attribution process. This suggests that a substantial norm in one checkpoint can significantly overshadow and dominate attribution results, which is suboptimal if the domination arises from systematic timestep patterns rather than sample-induced variance. However, this phenomenon is notably alleviated after renormalization in Diffusion-ReTrac, providing more consistent and distinct influence estimations.

C Supplemental Visualizations

C.1 CIFAR-PLANES WITH HIGH SELF-INFLUENCE

Self-influence is used to identify outliers in the training dataset. While Diffusion-TracIn and ReTrac assign high self-influence to most of the 200 MNIST samples, certain CIFAR-plane samples also received high self-influence scores and are ranked among the top 200. These plane samples tend to have dark backgrounds or high contrast, which are also visually distinct from typical samples in the CIFAR-plane subclass (Figure 10).



Figure 10: **CIFAR-planes with High Self-Influence.** These four samples are assigned high self-influence scores by both Diffusion-TracIn and ReTrac. They are visually distinct from the typical plane samples in the training dataset.

C.2 CIFAR-PLANES INFLUENTIAL TO MNIST SAMPLES

The auxiliary task of Image Source Tracing pinpoints specific training samples that are responsible for the generation of a test sample. For an MNIST zero, while most of the retrieved proponents are MNIST zeros, some planes are also assigned high influences (Figure 11). We noticed that these planes are visually distinct from the other, and visually resemble the MNIST samples. They tend to exhibit a black background and the planes are centered in the middle, which highly resembles the layout of MNIST zeros. This further proves the effectiveness of Diffusion-ReTrac in identifying highly influential samples.

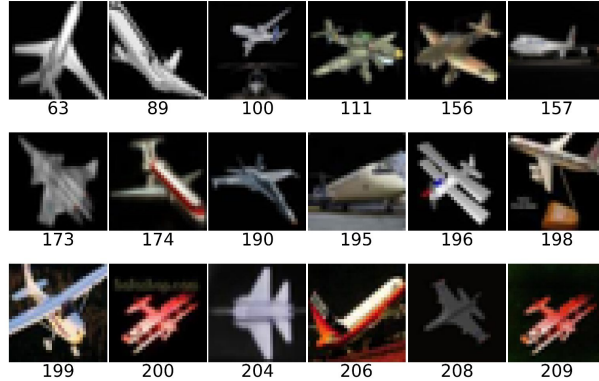


Figure 11: **CIFAR-planes Influential to MNIST Samples.** These CIFAR-Planes are assigned high influence scores to an MNIST zero test sample. The attribution results are in descending order and the corresponding ranking for each sample is labeled.

D Implementation Details

D.1 MODEL DETAILS

We trained a Diffusion Denoising Implicit Model (DDIM) [Song et al., 2020] with 1,000 denoising timesteps and 50 inference steps using an Adam optimizer. However, it is noted that our approach should remain consistent across variations of diffusion models, since the methods are designed based on the training process which is largely unaffected by differences in inference procedures. It may also be modified to accommodate the variations in training. Nonetheless, the practical form of TracInCP is expected to remain the same across these variations [Pruthi et al., 2020].

D.2 CHECKPOINT SELECTION

When estimating influences, it is ideal to select checkpoints with consistent learning and a steady decline in loss. Checkpoints that are early in the model’s learning stage often yield fluctuating gradient information, while those near model convergence offer limited insights into the attribution. Influence estimation at these early/late epochs of the learning process can introduce noise and compromise the accuracy of attribution results.

Attribution methods that rely on loss gradient norm information are also particularly sensitive to checkpoint selection. We observe that certain samples may exhibit an unusually large norm at specific checkpoints. When this checkpoint is used in Diffusion-TracIn, such samples emerge as generally influential with notably high influence on various test samples, overshadowing attribution results from previous checkpoints. This effect is mitigated in Diffusion-ReTrac due to renormalization, reducing the method’s susceptibility to dominant norms.

D.3 TIMESTEP SELECTION

To approximate the expectation over timesteps in the attribution efficiently, 50 linearly spaced timesteps over the denoising trajectory are used. This provides similar results to estimating influences across the entire trajectory using T timesteps. It is also observed that the loss induced is relatively stable at neighboring timesteps, while significant variation persists among distant timesteps. This provides justification for reducing computational costs by employing an adequate number of evenly spaced timesteps to approximate the loss over the entire trajectory.