TRUSTWORTHY DATASET PROOF: CERTIFYING THE AUTHENTIC USE OF DATASET IN TRAINING MODELS FOR ENHANCED TRUST

Anonymous authors

006

008 009 010

011

013

014

015

016

017

018

019

021

024

025

026

027

028 029

031

Paper under double-blind review

ABSTRACT

In the realm of deep learning, the veracity and integrity of the training data are pivotal for constructing reliable and transparent models. This study introduces the concept of Trustworthy Dataset Proof (TDP), which tackles the significant challenge of verifying the authenticity of training data as declared by trainers. Existing dataset provenance methods, which primarily aim at ownership verification rather than trust enhancement, often face challenges with usability and integrity. For instance, excessive operational demands and the inability to effectively verify dataset authenticity hinder their practical application. To address these shortcomings, we propose a novel technique termed Data Probe, which diverges from traditional watermarking by utilizing subtle variations in model output distributions to confirm the presence of a specific and small subset of training data. This modelagnostic approach improves usability by minimizing the intervention during the training process and ensures dataset integrity via a mechanism that only permits probe detection when the entire claimed dataset is utilized in training. Our study conducts extensive evaluations to demonstrate the effectiveness of the proposed data-drobe-based TDP framework, marking a significant step toward achieving transparency and trustworthiness in the use of training data in deep learning.

1 INTRODUCTION

Against the backdrop of the rapid development of deep learning technologies, the reliability and
 transparency of models are increasingly being scrutinized, and the authentic use of training data
 is the cornerstone of building effective and trustworthy models (Shayne, 2024; Anomalo, 2024;
 Aldoseri et al., 2023). However, the actual usage of training data is often based on the trainer's
 self-report, which is extremely difficult to verify in practice.

037 The authenticity and integrity of training data are significant security concerns in the field of deep 038 learning, as even minor tampering with the training data can lead to significant changes in model behavior. For instance, numerous studies (Chen & Babar, 2024; Khaddaj et al., 2023; Mengara et al., 2024; Saha et al., 2020) have shown that attackers can embed a small number of backdoor samples 040 in training data, which could be activated during the model's deployment, posing security risks. On 041 the other hand, trainers may introduce illegal or non-compliant data to enhance the performance 042 of the model without informing the users, especially in competitions or commercial applications, 043 raising concerns regarding the transparency and fairness. For example, high-profile lawsuits have 044 been initiated against major tech companies like Google and StabilityAI, where plaintiffs argue that 045 their copyrighted or personal data were used without permission to train AI systems (Sabine, 2024). 046

In this paper, we define the problem of Trustworthy Dataset Proof (TDP) to formally describe this
challenging task of verifying the integrity of training datasets, and further explore possible solutions.
As illustrated in Fig. 1, unlike the widely studied concept of data provenance, which focuses on the
ownership verification of a personal dataset, the TDP problem seeks to verify the authentic use of a
trustworthy dataset by the trainer, thereby enhancing trust.

Challenges: Although existing dataset provenance techniques hold promise for addressing the TDP problem, they still face several formidable challenges. O Usability: these techniques often fail to meet practical requirements, such as demanding that the trainer provide all details of the training



Figure 1: **The Overview of the Trustworthy Dataset Proof (TDP) Problem.** Compared to the widely studied data provenance techniques, the TDP presents significant differences in terms of dataset utilization and verification scope, making it a more challenging and novel problem.

process (Choi et al., 2024; Jia et al., 2021), requiring the training of a dedicated classifier for each verification request (Maini et al., 2021; Dziedzic et al., 2022), or utilizing a backdoor for watermarking (Adi et al., 2018; Tang et al., 2023), which could potentially be exploited maliciously (Mengara et al., 2024; Khaddaj et al., 2023). **9** Integrity: existing technologies are inadequate to certify the integral use of the training dataset claimed by the trainers. In contrast, they can only verify whether the actual dataset used approximates the distribution of the claimed dataset (Maini et al., 2021; Dziedzic et al., 2022; Choi et al., 2024; Jia et al., 2021), or whether it contains specific samples from the claimed dataset (Adi et al., 2018; Tang et al., 2023).

Motivations: To address the challenges of utility, we propose a novel concept named *Data Probe*.
This technique, distinct from traditional watermarking, does not require a model to produce a predetermined output. Instead, it leverages subtle distinctions in model outputs distribution to validate
its existence. It is designed to be model-agnostic, thereby reducing constraints on the training process. For further tackling the integrity challenges, we bind the integrity of the training dataset to
the *data probe selection strategy*. This ensures that successful probe implantation and detection are
contingent upon the use of the complete dataset for training, serving as a possible solution to TDP.

Contributions: The contributions of our study are concluded as follows:

- To the best of our knowledge, this is *the first exploration* of the challenging security problem associated with verifying the comprehensive and authentic use of the training dataset.
- We formalize this issue as a Trustworthy Dataset Proof (TDP) problem, analyze the challenges posed by existing technologies, and derive technical insights for potential solutions.
 - We innovatively design a watermarking-like technique called Data Probe, which underpins a TDP framework that is highly available and capable of verifying the integrity of training data.
- Extensive evaluations demonstrate the effectiveness of our proposed methods, validating our approach in various experimental settings.
- 089 090

079

080

081

082

084

085

087

2 RELATED WORK

091

Existing research on the protection and traceability of training datasets primarily focuses on verifying the ownership of datasets Although this objective differs from the goal discussed in this paper,
which is to verify the integrity of training datasets to enhance credibility, we review these studies to
better understand the existing challenges and potential technical motivations, with a more detailed
discussion to follow in the Sec. 4.

Dataset watermarking. It achieves ownership authentication by embedding backdoors into training 098 datasets (Adi et al., 2018). Specifically, the owner of a dataset can select and modify a small subset of 099 samples to serve as backdoor triggers, which can induce atypical outputs from the model. Once the model is trained on this dataset, the backdoor is automatically implanted. Subsequently, the copy-100 right holder of the dataset can verify whether a suspected model exhibits the expected anomalous 101 behavior by activating the trigger, thereby confirming the occurrence of the training. Building on 102 this, many studies focus on enhancing the stealth and harmlessness of these backdoors. For instance, 103 Tang et al. (2023) propose to construct a clean-label backdoor by applying adversarial perturbations 104 to samples, which embeds a specific backdoor pattern without altering the labels of the samples. 105

- 106 Membership and Dataset Inference. Membership Inference (MI) (Shokri et al., 2017; Salem et al.,
- 107 2019; Yeom et al., 2018; Song & Mittal, 2021) is an attack targeting the privacy of training datasets, but its technical approach can also be reversed for tracing training sets. MI exploits the phenomenon

108 of overfitting in deep learning (Yeom et al., 2018), whereby a model better memorizes samples from 109 its training set, distinguishing them from non-training samples. For instance, Shokri et al. (Shokri 110 et al., 2017) proposed training numerous shadow classifiers to reveal such distinctions. Other studies 111 analyze various score metrics produced by the model when different samples are inputted, such 112 as confidence (Salem et al., 2019), loss (Yeom et al., 2018), and entropy (Song & Mittal, 2021). To further enhance the stability of verification, Maini et al. (2021) proposed the Dataset-Inference 113 (DI) technique, which actually conducts MI at the distribution level and utilizes the characteristic 114 differences in the model's responses to training and test samples. 115

116 Proof of Training Data. Proof-of-Training-Data (PoTD) (Choi et al., 2024) is a protocol enabling 117 a model trainer to assure a verifier of the specific training data responsible for generating a set of 118 model weights, which can confirm both the quantity and type of data. The protocol mandates that the model trainer meticulously document and provide comprehensive details throughout the training 119 process, referred to as a training transcript, including the dataset, training codes, hyperparameters, 120 and intermediate checkpoints. The goal of PoTD is similar to our work. However, it focuses on 121 ensuring a precise correspondence between the training processes and model parameters, but fails to 122 verify the presence of subtle manipulations in the training dataset, as stated in their work limitations. 123

124 125

3 FORMALIZE TDP

126 3.1 PROBLEM DEFINITION

128 In the Trustworthy Data Proof (TDP) problem, two key roles are identified: model trainer and 129 verifier. The trainer is capable of utilizing a publicly credible dataset \mathcal{D} to train a model \mathcal{M} . On 130 the other hand, the verifier is tasked with providing a mechanism to verify whether \mathcal{M} was trained 131 using \mathcal{D} in a reliable and trustworthy manner. We formally define two functions to implement TDP:

Definition 1 (Trusted Training) T-Train(\mathcal{D}) $\rightarrow \mathcal{M}, \mathcal{C}$

Trusted Training is a training mechanism specified by verifier, denoted as T-Train(). It could be further represented as T-Train = $\mathbb{T} \cup \mathbb{O}$, where \mathbb{T} refers to general model training procedures, and \mathbb{O} represents additional operations required by the verifier. The trainer need to use T-Train() to train models to meet the requirements for subsequent trustworthy verification. Specifically, it performs T-Train(\mathcal{D}) to obtain a trained model \mathcal{M} , and at the same time, get a certificate \mathcal{C} for the subsequent verification.

140 **Definition 2** (Verification) Verify $(\mathcal{D}, \mathcal{M}, \mathcal{C}) \rightarrow \{0, 1\}$

Verification is the process by which the verifier provide the proof, denoted as Verify(). The trainer want to assert that \mathcal{M} is trained on \mathcal{D} . He then provides the verifier with \mathcal{M} and \mathcal{C} obtained through T-Train, as well as the \mathcal{D} to be validated. Verify() outputs 1 to indicate that it judges the model is indeed trained on \mathcal{D} , while 0 indicates it is not.

Model

Trainer

We justify the rationality of the above definitions: **①** Generality: the additional operations \mathbb{O} and \mathcal{C} can both be empty \emptyset . At this point, T-Train and Verify describe the most basic process of solving the TDP problem without any assumptions, making them universal. **②** Necessity: accord-

Figure 2: The conceptual framework of TDP.

ing to the consensus in the relevant research field (Shayne, 2024; Aldoseri et al., 2023; Choi et al., 2024), due to the high-dimensional complexity of DL models and the non-convex optimization training process, it is difficult to directly verify the fact that \mathcal{M} is trained on the complete \mathcal{D} through \mathcal{M} alone. Therefore, the verifier should be allowed to intervene and restrict the training process to enhance their ability during the verification process.

- 155
- 156 3.2 THREAT MODEL

In TDP, the verifier is generally served by trusted authorities to ensure the fairness and effectiveness of the verification process, and to defend against potential acts that may subvert the verification. Therefore, the verifier is defined as the defender. On the other hand, the trainer may maliciously exploit the TDP mechanism for unjust benefits, thus the dishonest trainer is defined as the attacker. In this section, we will analyze the goals and capabilities of each party separately.

Defender's Goals. The fundamental goal of the defender is to ensure the correctness of the verification (G1), accompanied by several performance-related goals (G2, G3, G4):

- **G1. Fidelity.** If the trainer indeed carries out T-Train(\mathcal{D}) $\rightarrow \mathcal{M}, \mathcal{C}$, then the defender, executing Verify($\mathcal{D}, \mathcal{M}, \mathcal{C}$), must output 1 with high probability. In all other cases, output 0.
- **G2.** Low-Invasiveness. In T-Train = $\mathbb{T} \cup \mathbb{O}$, it is advisable to minimize \mathbb{O} 's restrictions on \mathbb{T} , as excessive limitations can affect the trainer's flexibility in defining their own models and training processes, hindering the generalizability of the TDP mechanism.
- **G3. Harmlessness.** Implementing TDP should not compromise the model's performance significantly. Specifically, the performance of the trained model $\mathcal{M} = \text{T-Train}(\mathcal{D})$ should be close to the performance of $\mathcal{M}_{\mathbb{T}} = \mathbb{T}(D)$, which is obtained using only the standard training process.
- G4. Efficiency. The verification process should be computationally efficient.

Defender's Capabilities. Based on practical scenarios, we make the following assumptions and restrictions on the defender's capabilities. In T-Train = $\mathbb{T} \cup \mathbb{O}$, we assume that the additional operations required by the verifier, the \mathbb{O} , is *model-agnostic*. \mathbb{O} will not modify the model architecture, nor obtain specific procedures or hyperparameters of \mathbb{T} . In Verify, we assume that the defender only has *black-box access* to the model \mathcal{M} to be verified, and can only query the model's prediction probabilities, i.e., $y = \mathcal{M}(x)$, where x is the input sample, $y \in \mathbb{R}^d$, and d is the number of classes.

Attacker's Goals. The attacker is a dishonest model trainer, whose goal is to attempt to subvert the defender's fidelity goal (G1), formally defined as follows:

When the attacker use the modified dataset $\mathcal{D}^* \neq \mathcal{D}$ in trusted training, they would obtain T-Train $(\mathcal{D}^*) \rightarrow \mathcal{M}^*, \mathcal{C}^*$. But they purposely claim to the defender (verifier) that *the model* \mathcal{M}^* *is trained on* \mathcal{D} , attempting to make Verify $(\mathcal{D}, \mathcal{M}^*, \mathcal{C}^*) \rightarrow 1$.

It is important to emphasize that the attacker's primary objective is to successfully complete the TDP. By obtaining the proof that *the model's training set is a credible dataset* \mathcal{D} , the attacker can gain the trust of users. As a result, the model \mathcal{M}^* it publicly releases is more likely to be downloaded, deployed, and used. On this basis, the attacker attempts to 'hide' the fact that it actually tampered with the training dataset for training.

Attacker's Capabilities. We assume that the attacker can only modify the dataset \mathcal{D} but cannot manipulate the verifier's additional operations \mathbb{O} in T-Train. The reasonableness of this assumption is consistent with the discussion in problem definition (Sec. 3.1). The defender (verifier) needs to intervene in the model training with certain operations \mathbb{O} to facilitate verification. Once the attacker can manipulate \mathbb{O} , they can easily neutralize various verification mechanisms of the defender. Therefore, we reasonably limit the capabilities of the attacker, requiring them to at least fully execute T-Train. In other words, any modifications made by the attacker to \mathcal{M} must go through T-Train.

198 199 200

201

202

165

166

167

168

169

170

171

172

173

174

181

182

183

184

185

187

188

189

190

191

4 CHALLENGES AND MOTIVATIONS

4.1 CHALLENGES FOR CURRENT TECHNIQUES

Based on the formal definition of TDP in Sec. 3, we utilize representative existing dataset provenance
technologies, including Watermarking, MI/DI, and PoTD, to propose some hypothetical solutions
for TDP and analyze the challenges involved. Detailed definitions and analyses of each approach
are elaborated in the Appendix A, from which we conclude that the application of existing schemes
to TDP primarily confronts challenges in two aspects: *integrity* and *usability*.

Challenge 1: Intergrity. All existing schemes fail to ensure the integrity of the training dataset D. Specifically, if an attacker claims to have used D for training but actually employs D^* , current schemes can only verify whether D^* approximates the distribution of D (MI/DI, PoTD), or whether D^* contains several specific samples from D (Watermarking). However, none of these approaches can confirm that $D^* \neq D$.

This disadvantage fundamentally stems from differences in the definition of the threat model, which
is also detailed in Fig. 1. In most of the aforementioned works, the party holding the data primarily
aims at *ownership verification*. Therefore, they are typically regarded as trustworthy or defenders,
with no incentive to actively alter their dataset. However, in TDP, the main goal of the data holder

and user is to *enhance trust*, which might lead them to actively modify the dataset while still claiming
 to use an authentic dataset to gain undue benefits. This disparity is the root of the vulnerabilities in
 existing methods, rendering them unable to meet the fundamental goal of fidelity (G1) in TDP.

Challenge 2: Usability. Some existing techniques may exhibit deficiencies in usability within practical scenarios. For instance, the PoTD scheme does not meet the defender's goal: low-invasiveness (G2), as it requires accessing all information from the model trainer during training. DI, to some extent, fails to satisfy the efficiency goal (G4), since it may necessitate training a dedicated classifier for each verification request. The watermarking approach could potentially compromise the defender's goal of harmlessness (G3), as watermarks that utilize backdoors might be maliciously exploited, thereby introducing inherent security risks.

226 227

228

4.2 MOTIVATIONS AND INSIGHTS

Tackling usability challenge. Among various approaches, we consider the concept of watermarking
 to exhibit the highest usability. It imposes minimal constraints during the training phase and offers
 the highest efficiency during verification, with only requiring black-box access to the model under
 test. However, its shortcomings originate from the use of backdoors (Adi et al., 2018), typically
 requiring the model to produce anomalous outputs when presented with specific trigger samples,
 which often leads to degradation in model performance or security risks (Mengara et al., 2024)

We aim to relax assumptions regarding watermark capabilities to enhance usability while ensuring
sufficient capability to determine the occurrence of watermark. We also draw inspiration from MI,
noting that the membership of samples could be determined based on discrepancies in model output.
Consequently, we introduce a new concept: Data Probe, for conducting watermarking operations.
Its formal definition is as follows:

240 Definition 3 (Data Probe). Data probe is defined as a small subset of samples, \mathbf{x}_p , within the train- **241** ing dataset \mathcal{D} . Based on this, the training dataset can be divided into two parts: $\mathcal{D} = \{\mathbf{x}_p \cup \mathbf{x}_{np}\}$, **242** where \mathbf{x}_{np} refers to 'non-probe' samples. After training the model \mathcal{M} on \mathcal{D} , we expect a noticeable **243** difference in the output distribution of the model when inputting data from these two subsets, i.e.,

$$\Pr(\mathbf{y}|\mathbf{x}_p; \mathcal{M}) \neq \Pr(\mathbf{y}|\mathbf{x}_{np}; \mathcal{M})$$

Data probe can be considered as a weakened version of a backdoor, it only necessitates a discernible difference in the output compared to that from non-probe inputs \mathbf{x}_{np} , such as a slight increase in prediction confidence, but not a 'directed' output. From the perspective of MI, although \mathbf{x}_p and \mathbf{x}_{np} both belong to the training set, we expect \mathbf{x}_p to behave more like a 'special member' of the dataset.

Tackling intergrity challenge. Simple data probe selection and implanting strategies may still fall into the dilemmas encountered by watermarking schemes. Hence, our insight is to *bind the integrity of the training dataset with the data probe selection strategy*, ensuring that successful probe implantation and detection can only occur when the dataset is used in its entirety for training.

- 5 IMPLEMENT TDP
- 256 257 258

259

244

245

250

251

252

253

254 255

Conceptual Overview. Our proposed data-probe-based TDP framework is illustrated in Fig. 3 and formally described in Algorithm 1. Relevant functions are denoted with the subscript DP.

260 In T-Train_{DP}, a pseudo-random mechanism ProbeSelect is introduced to select the data probe 261 \mathbf{x}_p . It needs to perform a keyed-hash on training dataset \mathcal{D} based on a user-specific key k. Next, 262 operation \mathbb{O}_{DP} is exerted on the selected probe to facilitate its 'implantation' into the model \mathcal{M} 263 during training. Subsequently, the trainer can carry out the normal training process \mathbb{T} and submit the 264 key k as a certificate \mathcal{C} for verification.

In Verify_{DP}, the verifier first reproduces the pseudo-random process ProbeSelect based on C to select data probe \mathbf{x}_p and non-probe \mathbf{x}_{np} from \mathcal{D} claimed by the trainer. He then calculates the scores \mathbf{s}_p , \mathbf{s}_{np} through a function ProbeScore to measure the difference in the output distribution of these two groups of data in model \mathcal{M} . If there is a noticeable difference between \mathbf{s}_p and \mathbf{s}_{np} , it is considered that the data probe \mathbf{x}_p has been detected, authenticating the genuine use of the training dataset \mathcal{D} .

281

283

284

287



Figure 3: The Conceptual Overview of the Proposed TDP Framework based on Data Probe.

Verify_{DP} satisfies the main objective of the defender: fidelity (G1). The principle lies in that, if the trainer indeed use \mathcal{D} for training \mathcal{M} , the data probe \mathbf{x}_p selected and implanted in T-Train_{DP} will be completely consistent with the probe calculated in Verify_{DP}. Assuming that the data probe can correctly perform the function defined in Definition 3, the verifier can reliably certify its claim. 285 Conversely, if the trainer makes any minor modifications to D, as shown in the lower part of Fig. 3, the implanted data probe \mathbf{x}_{p}^{*} will not correspond with \mathbf{x}_{p} , thus failing the validation.

288 **Probe Selection.** The probe selection rule is based on pseudo-randomness, which utilizes the 289 *uniqueness* of hash functions (Rivest, 1992). Initially, it calculates a hash of the complete dataset \mathcal{D} to obtain a unique hash value. Subsequently, we use this hash value as a random seed to randomly 290 select a small subset of samples from the training dataset as data probe x_p , as described in Lines 3-4 291 of Algorithm 1. Clearly, \mathcal{D} and \mathbf{x}_p are bound together, and any minor modification of \mathcal{D} will result 292 in changes to the hash value, which in turn leads to changes in the selection of data probe x_p . 293

294 However, simple hash calculations pose security risks. The same dataset will yield the same data probe for all users, which could be easily exploited by attackers. Therefore, we introduce the user-295 specific key k and replace the hash with keyed-hash, such that different users using the same trusted 296 dataset (such as CIFAR10 (Krizhevsky & Hinton, 2009)) will generate different data probes. 297

298 Probe Implantation. Reflecting on the 299 functionality of data probe described in 300 Definition 3, our goal is to elicit a special 301 response from the trained model \mathcal{M} to the data probe \mathbf{x}_p . At the same time, we also 302 need to fully consider the usability of the 303 scheme. \mathbb{O}_{DP} does not need to change the ar-304 chitecture of model \mathcal{M} , nor does it need to 305 obtain the hyperparameters in training pro-306 cess \mathbb{T} . Therefore, we have determined that 307 \mathbb{O}_{DP} should only perform data-level opera-308 tions on data probe \mathbf{x}_p , without needing to 309 change the normal training process \mathbb{T} .

Table 1: Data Probe Types and Principles. PS represents the generic scoring process, calculating the probe score s_p and non-probe score s_{np} . l_p and l_{np} denote the ground-truth labels. l_p^t and l_{np}^t are targeted labels defined in TP.

Туре	\mathbf{s}_{p}	\mathbf{s}_{np}	Expectation
PP	$PS(\mathcal{M},\mathbf{x}_p,\mathbf{l}_p)$	$PS(\mathcal{M},\mathbf{x}_{np},\mathbf{l}_{np})$	$\mathbf{s}_p^{PP} > \mathbf{s}_{np}^{PP}$
AP	$PS(\mathcal{M},\mathbf{x}_p,\mathbf{l}_p)$	$PS(\mathcal{M}, \mathbf{x}_{np}, \mathbf{l}_{np})$	$\mathbf{s}_p^{AP} < \mathbf{s}_{np}^{AP}$
UP	$PS(\mathcal{M},\mathbf{x}_p,\mathbf{l}_p)$	$PS(\mathcal{M},\mathbf{x}_{np},\mathbf{l}_{np})$	$\mathbf{s}_p^{UP} < \mathbf{s}_{np}^{UP}$
TP	$PS(\mathcal{M},\mathbf{x}_p,\mathbf{l}_p^t)$	$PS(\mathcal{M}, \mathbf{x}_{np}, \mathbf{l}_{np}^t)$	$\mathbf{s}_p^{\mathrm{TP}} > \mathbf{s}_{np}^{\mathrm{TP}}$

310 Inspired by existing dataset provenance research, we have developed four different types of data 311 probe: Prominent Probe (PP), Absence Probe (AP), Untargeted Probe (UP), and Targeted Probe 312 (**TP**). We summarize the principles of various probes in Tab. 1 and detail their concepts and imple-313 mentations in this section. The unique operations of different types of probes are marked by super-314 scripts. Besides, the performance and characteristics of the four probes are compared in Sec. 6.2. 315

O Prominent Probe (PP). 'Prominent' means the model's response to the data probe x_p is more 316 significant compared to non-probe \mathbf{x}_{np} . In principle, we aim to make the model more overfit on the 317 \mathbf{x}_p , providing *more confident* scores. \mathbb{O}_{DP}^{PP} can be implemented through the built-in data sampling 318 mechanism of the deep learning computation library, by assigning a higher weight to the x_p . The 319 \mathbf{x}_p is therefore more likely to be selected during the training , resulting in better fitting to the \mathbf{x}_p . 320

2 Absence Probe (AP). The principle of AP is exactly opposite to that of PP. To maximize the 321 insignificance, we consider an extreme case where the probe weight is 0, meaning the model has 322 never encountered \mathbf{x}_p during training. Then the trained model \mathcal{M} should provide *less confident* 323 scores for \mathbf{x}_p . $\mathbb{O}_{\mathsf{DP}}^{\mathsf{AP}}$ can be implemented in a similar manner to $\mathbb{O}_{\mathsf{DP}}^{\mathsf{PP}}$ by setting the probe weight to 0.

324 **③** Untargeted Probe (UP). Untargeted Probe utilizes the concept of data poisoning (Adi et al., 325 2018; Mengara et al., 2024), where \mathbb{O}_{DP}^{UP} assigns random, incorrect labels to the probe \mathbf{x}_p . In this 326 case, we expect the model to be *less confident* in its predictions for the probe x_p . For instance, a 327 decrease is witnessed in the predicted probabilities for their ground truth classes.

328 **④** Targeted Probe (TP). In contrast to UP, Targeted Probe is where \mathbb{O}_{DP}^{TP} uniformly assigns a random label, marked as l^t , to the probe. In this case, we expect the trained model to be *more confident* in 330 predicting the \mathbf{x}_p as class l^t . Assuming the number of probe and non-probe are N_p and N_{np} , we 331 construct two sets of labels: $\mathbf{l}_p^t = \{l^t\}^{N_p}$ and $\mathbf{l}_{np}^t = \{l^t\}^{N_{np}}$ and obtain scores as stated in Tab. 1. 332

Probe Score Calculation. The detection of data probe relies on comparing the saliency scores of 333 probe s_p and non-probe s_{np} . In this section, we will demonstrate possible scoring methods. Inspired 334 by the techniques for assessing sample saliency adopted in the MI work, we propose the following 335 four possible methods: confidence-based, loss-based, entropy-based, and modified-entropy-based. 336

337 According to the aformentioned ProbeScore function protocol, the inputs include a model \mathcal{M} , a set of data x, and a set of labels l. For the sake of brevity in expression, we describe the score calculation 338 method for individual samples $x \sim x$, without distinguishing between probe and non-probe, as they 339 actually use the same calculation method. Besides, we convert the corresponding label $l \sim l$ into the 340 index of the model's output for convenience in expression, i.e., when l indicates that x belongs to 341 class y, we use $\mathcal{M}(x)_y$ to represent the \mathcal{M} 's prediction probability for x being in that class. 342

O Confidence-based score (Conf) For a more *significant* sample during training, the model should 343 make predictions with higher confidence in it (Salem et al., 2019). we define the confidence-based 344 score as: $|Conf(\mathcal{M}, x) = \max(\mathcal{M}(x))|$ which will be directly return as the probe score. 345

346 **2** Loss-based score (Loss) For a more *significant* sample during training, the model should has a 347 lower prediction loss on it (Yeom et al., 2018). We mark the loss function, such as the cross-entropy, 348 as \mathcal{L} , and we define the loss-based score as: $|Loss(\mathcal{M}, x, l) = \mathcal{L}(\mathcal{M}(x), l)|$. The more significant 349 the sample x, the smaller the Loss. In order to make the return values of ProbeScore continuous, 350 where larger values represent more significance, we return -Loss as the probe score. 351

352 **③** Entropy-based score (Entr) For a more *significant* samples during training, the model's predic-353 tion on it should be close to the one-hot encoded label, i.e., its entropy will be close to 0 (Salem et al., 2019). We define the entropy-based score as: $|Entr(\mathcal{M}, x) = -\sum_i \mathcal{M}(x)_i \log (\mathcal{M}(x)_i)|$. Similar 354 355 to Loss, we return **-Entr** as the probe score.

356 **9** Modified-entropy-based score (Mentr) It is an enhanced version of Entr by considering the 357 ground-truth label l (Song & Mittal, 2021). We define the modified-entropy-based score as: 358 Mentr $(\mathcal{M}, x, l) = -(1 - \mathcal{M}(x)_y) \log (\mathcal{M}(x)_y) - \sum_{i \neq y} \mathcal{M}(x)_i \log (\mathcal{M}(x)_i)$. Similar to Loss, we 359 return -Mentr as the probe score. 360

361 **Probe Detection Metric.** According to the definition of Data Probe, probe scores are expected to differ in distribution from non-probe scores. Therefore, we aggregate the aforementioned sample-362 level scores into distribution-level metrics for probe detection. The following two metrics are employed: **O** Probe Saliency AUC (PSA). It is a new metric proposed in this work. It utilizes probe 364 scores to plot the ROC curve, further calculating the Area Under ROC curve (AUC) as the metric. The greater the PSA exceeds 0.5, the more it indicates that probe scores are separable from non-366 probe scores. Statistical test p-value (pV). It has been widely adopted in previous works for mea-367 suring distribution differences. A pV less than a certain level of significance, such as 0.1, indicates 368 a substantial difference in score distributions, signifying that the probe has been detected. We detail 369 the principles of these two metrics and their calculation methods for each probe in Appendix B. 370

6 EVALUATIONS 372

371

- 373 **Overview.** In the evaluations, we aim to investigate the following four research questions (RQs):
 - **RQ1.** Whether various types of data probes can effectively verify the integrity of the dataset?
- 375 **RQ2.** How many probes need to be implanted during training to achieve the verification ? 376
- **RQ3.** How do different probe score calculation strategies impact the effectiveness of detection ? 377
 - **RQ4.** How robust is the verification mechanism when attackers launch adaptive attacks ?

Table 2: **Comprehensive Probe Perfomance Evaluations.** Best performance of PSA and pV under each setting are hilighted as **BLUE** and **RED**. And the metric under probe-mismatch cases are marked as GRAY. Scores are reported as % except for pV.

							-				•										
	•	Ori.	200	DS A ↑	PP PSA*	nV∣	nV*	200	DS V \	AP PSA*	nV∣	nV*	200	DS V \	UP PS A *	nV∣	nV*	200	PS∆↑	TP	υĮ
		acc	acc	r3A	rsA	h∙↑	P۷	acc	I SA	ISA	CIFAR-1	0	acc	r SA	rsA	р∙↓	P۲	acc	r SA	rsa	Р
ResNet	8 8	5 22	85 19	56.15	50.58	10^{-21}	0.35	85.08	52.14	49 97	0.01	0.42	84 68	57.09	49.60	0.02	0.55	84 92	59.22	50.18	10
Mobilel	Vet 8	4.55	84.12	55.95	49.60	10^{-16}	0.54	84.64	52.59	50.09	0.02	0.39	84.25	57.02	49.99	10^{-4}	0.37	84.43	57.17	50.16	10
Shufflel	Vet 8	4.69	83.06	57.20	50.41	10^{-43}	0.50	82.57	52.66	49.67	0.03	0.41	86.77	59.96	49.53	10^{-5}	0.39	82.83	57.87	49.78	10
DenseN	et 8	6.32	85.27	55.58	50.11	10^{-27}	0.51	86.18	52.60	49.42	0.08	0.55	83.04	57.63	48.95	10^{-3}	0.35	86.49	59.51	49.81	1(
											SVHN										
ResNet	8 9	2.03	91.86	51.28	49.25	10^{-11}	0.65	91.88	50.68	50.84	0.29	0.28	92.10	53.38	51.04	0.27	0.24	92.15	53.85	50.59	10
Mobile	Vet 9	1.98	91.43	51.17	49.33	10^{-6}	0.64	91.67	50.74	50.41	0.28	0.56	91.81	52.79	50.48	0.26	0.29	91.91	52.83	50.40	0
Shufflel	Vet 9	2.21	91.50	51.32	49.44	10^{-4}	0.78	91.66	51.01	50.44	0.30	0.35	92.41	52.84	50.72	0.35	0.21	91.75	53.56	50.23	10
DenseN	et 9	2.45	91.89	51.40	49.32	10^{-3}	0.67	92.37	51.21	50.72	0.26	0.29	92.00	53.29	50.58	0.19	0.25	92.03	53.44	50.46	1
										(CIFAR-1	00									
ResNet	8 6	64.62	63.57	69.20	50.06	10^{-99}	0.66	63.83	58.07	49.29	10^{-10}	0.48	63.95	71.28	49.65	10^{-27}	0.62	64.39	79.83	50.67	10
Mobile	Vet 6	2.45	61.17	69.36	50.72	10^{-99}	0.40	61.68	56.75	49.03	10^{-7}	0.64	62.30	63.42	48.89	10^{-11}	0.63	62.35	69.04	50.09	10
Shufflel	Vet 6	0.22	59.40	70.13	50.33	10^{-99}	0.51	59.86	57.39	50.02	10^{-8}	0.53	60.30	63.27	49.49	10^{-12}	0.69	60.48	71.14	50.66	10
DenseN	et 6	4.42	63.44	69.21	50.40	10 ⁻⁹⁹	0.57	64.39	57.67	49.81	10^{-10}	0.65	64.33	68.12	49.27	10^{-12}	0.66	64.61	77.56	50.64	10
										Tiny	ImageN	et-200)								
ResNet	8 5	3.00	51.03	76.26	50.17	10^{-54}	0.47	51.83	62.89	50.30	10^{-29}	0.41	53.35	49.82	50.61	0.44	0.38	53.18	49.94	49.21	0
Mobile	Vet 5	1.06	49.41	75.38	49.47	10-75	0.62	50.12	57.35	50.44	10^{-6}	0.46	50.83	50.20	50.56	0.41	0.42	51.00	49.87	49.68	0
Shufflel	Vet 5	0.11	47.85	76.02	49.85	10-78	0.44	48.47	57.99	50.31	10-9	0.51	50.04	49.58	50.20	0.45	0.55	49.81	49.38	50.15	0
DenseN	et 5	5.26	53.16	74.63	50.02	10-55	0.45	53.86	60.51	50.00	10^{-20}	0.55	55.72	49.58	50.02	0.40	0.51	55.59	50.05	49.63	- 0

6.1 Setup

Datasets: Four datasets are adopted: CIFAR-10 (Krizhevsky & Hinton, 2009), SVHN (Netzer et al., 2011), CIFAR-100 (Krizhevsky & Hinton, 2009), and Tiny-ImageNet-200 (Le & Yang, 2015).

403 Models: We employed various architectures: ResNet18 (He et al., 2016), MobileNet (Howard
404 et al., 2017), ShuffleNet (Zhang et al., 2018), and DenseNet (Huang et al., 2017). These models
405 were chosen to ensure a broad evaluation of performance across different architectural dynamics.

406
 407
 408
 408
 406
 407
 408
 408
 408
 408
 408
 409
 408
 409
 408
 409
 408
 409
 408
 409
 408
 409
 409
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400

409 For further specifics, such as hyperparameters, please refer to the Appendix C.

410 411

412

399

400

378

379

380

6.2 EVALUATION RESULTS

Comprehensive performance comparison (RQ1). We comprehensively evaluated the perfor-413 mance of various architectures when trained on different datasets using the four types of data probes 414 proposed for TDP. The results are displayed in the Tab. 2. All experiments were repeated five times, 415 and the mean values were recorded. We tested the accuracy of models trained directly without any 416 probes as a baseline for comparison, denoted as "Ori.". For each type of probe, in addition to train-417 ing accuracy, we recorded two metrics when the declared dataset matched the actual training dataset: 418 **PSA** and p-value (**pV**). Additionally, we tested metrics for mismatched declared and training models 419 as a contrast (**PSA**^{*} and \mathbf{pV}^*). In each experimental group, we selected 1% of the training set to 420 serve as probes. All probe scores were calculated using Mentr.

421 Through extensive experimental comparisons, we summarize the following results: **1** The implan-422 tation of data probes has almost no impact on the performance of the original model. As shown 423 by the accuracy metrics in Tab. 2, the performance of models with data probes deviates minimally 424 from the original performance, with the majority of variations within an acceptable range ($\leq \pm 1\%$). 425 This meets the defender's goal of harmlessness (G3). ⁽²⁾ Data probes generally achieve the task 426 of verifying dataset completeness. Specifically, for the PSA metric, we expect it to approach 0.5 427 when probes do not match, and significantly exceed 0.5 when they do match. For the pV, values 428 below 0.1 indicate that probes have been detected with high confidence, while in cases of mismatch, we expect the p-value to be as high as possible. It can be seen that most probes meet the afore-429 mentioned requirements across various training sets and models, successfully implementing TDP. 430 Additionally, we observed that among the four types of probes, PP exhibits high significance in the 431 p-value metric, while TP demonstrates a clear advantage in the PSA metric.



436

437

438 439

440

441

442

451



Figure 4: Evaluations of the Impact of Probe Implant Quantity We use a gray dashed line to mark two critical reference values: PSA = 0.5 and pV = 0.1. In PSA, lighter-colored dashed lines are non-probe scores. Here we present the results of PP, with the complete results shown in Fig. 7.

Impact of probe quantity (RQ2). We assessed how the performance metrics of TDP vary with an 443 increase in the number of probes, ranging from 0.1% to 2% of the training set, as shown in Fig. 4. 444 We trained various models on the CIFAR-10 dataset. All experiments were repeated five times, and 445 the mean values as well as the standard deviation were recorded and displayed. The probe scores 446 were also calculated using Mentr. From the evaluation results, we can conclude that selecting 447 approximately 1% of the training set samples is sufficient to reliably implement TDP. As the 448 number of implanted samples increases, the performance of the model may experience a negligible 449 decline, while the metrics of PSA and p-value tend to stabilize. We believe that using about 1% of 450 the training data as probes effectively achieves TDP, representing an appropriate trade-off between model verifiability and performance. 452

Impact of probe score calculation methods (RQ3). 453

We analyzed the viability of various possible probe 454 score calculation methods introduced in Sec. 5, as dis-455 played in Fig. 5. We trained ShuffleNet on CIFAR-10 456 using various default settings for different probes and 457 calculated the probe scores using four different meth-458 ods. All experiments were repeated five times, and the 459 mean values were reported. From the experimental re-460 sults, we found that Mentr exhibits the best gener-461 alizability and detection effectiveness. It effectively



Figure 5: Evaluations of Various Probe Score Calculation Methods.

detected data probes when applied across all four probe schemes. Additionally, we discovered that 462 Conf and Entr are completely inapplicable to TP, as TP relies on inducing samples to point towards 463 a specific label, whereas the calculations for Conf and Entr are independent of the label. 464

465 **Robustness to adaptive attackers (RQ4).** We evaluate a worst-case security scenario in which the 466 attacker knows the key k used in the probe selection process and the specific type of probe employed. Specifically, the attacker can replicate the probe \mathbf{x}_p corresponding to \mathcal{D} after training the model \mathcal{M}^* 467 using \mathcal{D}^* , and attempt to illegitimately pass the verification for \mathcal{D} by embedding the probe into the 468 model. We therefore designed four probe forging attacks for each probe, with the underlying princi-469 ples detailed in Appendix D, denoted with the prefix F for "Forged". Evaluations adopt ShuffleNet 470 on CIFAR-10. We trained 10 different models randomly and designated 10 different target probes 471 for each model, resulting in a total of 100 experimental groups per attack. We recorded the mean 472 and standard deviation of various metrics before and after the attack. Additionally, we set the PSA 473 threshold at 0.51 and the p-value at 0.1, and documented the Attack Success Rate (ASR). 474

The complete training results are displayed in 475 Tab. 3. We found that all probe schemes, ex-476 cept for PP, exhibit a certain degree of ro-477 bustness. This indicates that in most cases, it 478 is challenging for attackers to easily deceive 479 the verification mechanism through forged 480 probes without compromising model perfor-481 Additionally, we observed certain mance. changes in metrics after the attacks, revealing 482 483 potential security risks. Therefore, keeping the user's key k hidden from the users, such as by 484 implementing it through a server API, might 485 be a solution to enhance robustness.

Table 3: Ev	aluations	s of the R (obustnes	ss to Ad	lap-
tive Probe	Forging	Attacks.	Metrics	before	and
after attack	are shown	n in GRA	Y and \mathbf{B}	OX.	

		PSA		pV		
Attacks	Acc	score AS		score	ASR	
EDD	84.0±0.4	50.1±1.7	600	0.28±0.25	6701	
FPP	84.0±0.5	54.5±1.6	02%	$10^{-6} \pm 10^{-5}$	01%	
EAD	83.6±0.3	50.7±1.7	1107-	0.45±0.27	270%	
I'AF	82.3±0.5	51.4±1.7	14%	0.20 ± 0.21	5270	
EUD	83.6±0.3	50.7±1.7	507	0.49±0.27	007	
FUP	83.0±0.4	51.2±1.7	5%	0.49 ± 0.28	0%	
ETD	83.6±0.3	50.8±1.5	0.07	0.37±0.26	100%	
FTP	82.7±0.5	51.4±1.5	9%	0.27±0.22	10%	

		Original	Water	mark	Dataset I	nference		Data Pro	be (Ours)	
Description	Sample	Trigger	Verify	Trigger	Verify	Trigger	Verify	Trigger	HASH	HASH
		Score	Result	Score	Result	Score	Result	Score	(train)	(test)
Automobile ↓ Deer		99.82 Deer	PASSED 1.000	84.24 Deer	PASSED 2.3×10 ⁻⁷	99.82 Deer	BLOCKED 0.85	98.06 Deer	90db8061ae4 Odcbd47f31bf 17e173651	9504b703430 662aad7aa7c 4a9f0333b2

Figure 6: **Case Studies for Comparison with Existing Techniques.** The scores for each validation are presented in GRAY, with watermarking representing the predicted probabilities, and DI and Data Probe indicating the p-values. HASH(train) shows the hash value computed on the tampered dataset during probe implantation, which differs from the HASH(test) obtained on the declared untampered dataset (CIFAR-10) during testing. More results are exhibited in Fig. 8 in the Appendix.

496 497 498

499

492

493

494

495

7 CASE STUDY

500 We conducted several case studies to validate the efficacy of the proposed Data-Probe-based TDP 501 (DP) when genuine modifications occur within datasets. **PP** is adopted for evaluations. Meanwhile, 502 we conducted comparative analyses with existing technologies, selecting two representative ap-503 proaches: Watermarking (Tang et al., 2023) (WM) and Dataset Inference (Maini et al., 2021) (DI). Although they are not originally designed for the TDP mission, we adapted them using the hypo-504 thetical schemes proposed in Appendix A, utilizing their official open-source implementations. We 505 employed CIFAR-10 and ResNet18 to test the verification results when an attacker claims to have 506 trained on CIFAR-10, while subtly modifying it in training, which is expected to fail verification. 507

508 Simulated modification. We simulated 509 two typical scenarios of dataset tampering: introducing additional data and embedding 510 backdoors. For the former, we randomly se-511 lected a small proportion of samples, from 512 0.01% to 1%, and duplicated them. For 513 the latter, we chose a small subset of sam-514 ples and applied minimal noise (bounded by 515 $l_{\infty} = 8/255$), adhering to the common con-516 figurations used in backdoor attacks. For

Table 4: Case Study of Simulated Modifications. Success rate (%) is shown in \boxed{BOX} with socres in GRAY, consistent with Fig. 6.

1.1.1	o · ·	Extra Data			Ba	ickdoor	Usability Eva.		
Method	Origin	0.01%	0.10%	1%	0.01%	0.10%	1%	Sec. Risk	Time
WM	100	0	0	0	0	0	5	35.60%	8.85
	1.0	1.0	1.0	0.99	0.99	0.99	0.97	55.00 %	0.05
DI	100	0	0	0	0	0	0	-	58min
	10	10	10	10	10	10	10		
DP(Ours)	100	100	100	100	95	100	95	97.20%	6.5s
	10 ~	0.53	0.58	0.46	0.54	0.50	U.46		

517 each setup, we conducted 20 repeated experiments and recorded the success rates (success is de-518 fined as *block from verification*, unless the trainer did not modify the dataset). Additionally, we 519 evaluated the usability metrics for each approach, namely the runtime and security risks. Regarding 520 security risks, DI was not assessed because it does not alter the training process. For WM we tested the prediction accuracy of the backdoor samples it used. For DP, we evaluated the prediction accu-521 racy of the data probe. The comparative results from Tab. 4 demonstrate that only the Data-Probe 522 approach successfully denied verification requests from attackers while exhibiting the lowest 523 time expenditure and minimal security risks. 524

Practical modification. We employed a representative and effective backdoor attack named
Witches' Brew (Geiping et al., 2021). This attack alters 1% of the samples so that the trained model
incorrectly classifies a targeted image, known as the trigger, as the wrong category. For instance, an
automobile would be recognized as a deer. The results are displayed in Fig. 6, and the conclusions
are consistent with those from the simulated experiments.

530 531

532

8 CONCLUSION

In this study, we highlight the importance of the Trustworthy Dataset Proof (TDP) in enhancing the veracity and integrity of training data for deep learning models. By introducing the novel Data Probe technique, this research successfully addresses the limitations of existing dataset provenance methods, which often falter in usability and integrity. The Data Probe, by leveraging subtle variations in model output distributions to verify the inclusion of specific training subsets, offers a model-agnostic and minimally invasive approach to dataset verification. Our extensive evaluations validate the effectiveness of our Data-Probe-based TDP framework, significantly advancing the pursuit of transparency and trustworthiness in training data usage.

540 REFERENCES

- Yossi Adi, Carsten Baum, Moustapha Cisse, Benny Pinkas, and Joseph Keshet. Turning your weakness into a strength: Watermarking deep neural networks by backdooring. In 27th USENIX secu-*rity symposium (USENIX Security 18)*, pp. 1615–1631, 2018.
- Abdulaziz Aldoseri, Khalifa N Al-Khalifa, and Abdel Magid Hamouda. Re-thinking data strategy and integration for artificial intelligence: concepts, opportunities, and challenges. *Applied Sciences*, 13(12):7082, 2023.
- 548 Anomalo. Data reliability in machine learning: Ensuring trustwor-549 hhttps://www.anomalo.com/blog/ 2024. thy model outputs, Apr. 550 data-reliability-in-machine-learning-ensuring-trustworthy-model-outputs/. 551
- Huaming Chen and M Ali Babar. Security for machine learning-based software systems: A survey of threats, practices, and challenges. *ACM Computing Surveys*, 56(6):1–38, 2024.
- Dami Choi, Yonadav Shavit, and David K Duvenaud. Tools for verifying neural models' training data. Advances in Neural Information Processing Systems (NIPS), 36, 2024.
- Adam Dziedzic, Haonan Duan, Muhammad Ahmad Kaleem, Nikita Dhawan, Jonas Guan, Yannis
 Cattan, Franziska Boenisch, and Nicolas Papernot. Dataset inference for self-supervised models.
 Advances in Neural Information Processing Systems (NIPS), 35:12058–12070, 2022.
- Jonas Geiping, Liam Fowl, W Ronny Huang, Wojciech Czaja, Gavin Taylor, Michael Moeller, and
 Tom Goldstein. Witches' brew: Industrial scale data poisoning via gradient matching. In *International Conference on Learning Representations*, 2021.
- 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.
- Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. IEEE, 2017.
- Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4700–4708, 2017.
- Hengrui Jia, Mohammad Yaghini, Christopher A Choquette-Choo, Natalie Dullerud, Anvith Thudi,
 Varun Chandrasekaran, and Nicolas Papernot. Proof-of-learning: Definitions and practice. In
 2021 IEEE Symposium on Security and Privacy (SP), pp. 1039–1056. IEEE, 2021.
- Alaa Khaddaj, Guillaume Leclerc, Aleksandar Makelov, Kristian Georgiev, Hadi Salman, Andrew Ilyas, and Aleksander Madry. Rethinking backdoor attacks. In *International Conference on Machine Learning*, pp. 16216–16236. PMLR, 2023.
- Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. *Master's thesis, Department of Computer Science, University of Toronto*, 2009.
- Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. *CS 231N*, 7(7):3, 2015.
- Pratyush Maini, Mohammad Yaghini, and Nicolas Papernot. Dataset inference: Ownership resolution in machine learning. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021.
- Orson Mengara, Anderson Avila, and Tiago H Falk. Backdoor attacks to deep neural networks: A survey of the literature, challenges, and future research directions. *IEEE Access*, 2024.
- Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Baolin Wu, Andrew Y Ng, et al.
 Reading digits in natural images with unsupervised feature learning. In *NIPS workshop on deep learning and unsupervised feature learning*, volume 2011, pp. 4. Granada, 2011.

594 595	Ronald Rivest. Rfc1321: The md5 message-digest algorithm, 1992.
596	Neschke Sabine. Legal challenges against generative ai: Key
507	takeaways Ian 2024 https://bipartisanpolicy.org/blog/
508	legal-challenges-against-generative-ai-kev-takeawavs/.
500	
599	Aniruddha Saha, Akshayvarun Subramanya, and Hamed Pirsiavash. Hidden trigger backdoor at-
601	tacks. In <i>Proceedings of the AAAI conference on artificial intelligence (AAAI)</i> , volume 34, pp.
602	11957–11965, 2020.
602	Ahmed Salem, Yang Zhang, Mathias Humbert, Pascal Berrang, Mario Fritz, and Michael Backes,
604	Ml-leaks: Model and data independent membership inference attacks and defenses on machine
605	learning models. In Proceedings of the Network and Distributed System Security Symposium
606	(NDSS), 2019.
607	
602	Longpre Shayne. Data authenticity, consent, and provenance for ai are all broken: What will it take
600	to fix them?, Apr. 2024. https://mit-genal.pubpub.org/pub/uk/op825/release/2.
610	Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference at-
611	tacks against machine learning models. In 2017 IEEE symposium on security and privacy (SP),
610	pp. 3–18. IEEE, 2017.
612	Lives Song and Dratack Mittal Systematic avaluation of microsy risks of machine learning models
61/	In 30th USENIX Security Symposium (USENIX Security 21) pp. 2615–2632, 2021
615	In som Oselvix security symposium (Oselvix security 21), pp. 2015–2052, 2021.
616	Ruixiang Tang, Qizhang Feng, Ninghao Liu, Fan Yang, and Xia Hu. Did you train on my dataset?
617	towards public dataset protection with cleanlabel backdoor watermarking. ACM SIGKDD Explo-
610	rations Newsletter, 25(1):43–53, 2023.
610	Samuel Veom Irane Giacomelli Matt Fredrikson and Somesh Iba. Drivacy risk in machine learn
620	ing. Analyzing the connection to overfitting In 2018 IEEE 31st computer security foundations
621	symposium (CSF), pp. 268–282, IEEE, 2018.
622	
623	Xiangyu Zhang, Xinyu Zhou, Mengxiao Lin, and Jian Sun. Shufflenet: An extremely efficient
624	convolutional neural network for mobile devices. In <i>Proceedings of the IEEE conference on</i>
625	computer vision and pattern recognition, pp. 0848–0850, 2018.
626	
627	
628	
629	
630	
631	
632	
633	
634	
635	
636	
637	
638	
639	
640	
641	
642	
643	
644	
645	
646	
647	

648 CHALLENGES FOR CURRENT TECHNIQUES А 649

650 In this section, based on the formal definition of TDP in Sec. 3, we combine existing dataset traceability or copyright authentication technologies to propose some hypothetical solutions for TDP, and 652 analyze the challenges involved.

653 **Case 1: Watermarking** 654

If we assume using watermarking technology (Adi et al., 2018; Tang et al., 2023) to achieve TDP, 655 656 then the verifier can validate the occurrence of training behaviors based on two steps: watermark embedding and watermark detection. We use the subscript WM to represent the relevant functions of 657 watermarking-based TDP. 658

659 In T-Train_{WM}, the trainer is required to embed watermarks into the model through training, and this 660 operation can be denoted as \mathbb{O}_{WM} . \mathbb{O}_{WM} allows the trainer to select some data samples from \mathcal{D} as 661 triggers \mathbf{x}_t according to certain rules, which can make the trained model produce a specific pattern of output \mathbb{P} . We represent this characteristic as: $\mathbb{P} = \mathcal{M}(\mathbf{x}_t)$. Then the trainer sets $\mathcal{C} \leftarrow \mathbf{x}_t$ and 662 submits it along with the trained model \mathcal{M} to the verifier. 663

664 In the verification, the verifier could make the judgement:

 $\mathsf{Verify}(\mathcal{D}, \mathcal{M}, \mathcal{C})_{\mathsf{WM}} \leftarrow \mathbb{1}\left(\mathbb{P} = M(C) \land C \in D\right).$

However, Verify_{WM} cannot meet the defender's primary goal: fidelity (G1), and we provide the 667 simplest counterexample. The attacker can slightly manipulate the remaining data $\mathcal{D}_{-} = \mathcal{D} \setminus \mathbf{x}_t$ 668 except for the triggers. It is noted that the current dataset $\mathcal{D}^* = \{\mathcal{D}_- \cup \mathbf{x}_t\} \neq \mathcal{D}$. Then the attacker 669 execute T-Train(\mathcal{D}^*)_{MM} $\rightarrow \mathcal{M}^*, \mathcal{C}^*$. Here the $\mathcal{C}^* = \mathcal{C} = \mathbf{x}_t$ because the \mathbf{x}_t remains unchanged. 670 Since the watermark embedding mainly relies on x_t , which means that the attacker can achieve 671 $\mathbb{P} = \mathcal{M}^*(\mathbf{x}_t)$. Therefore, it is obvious that attackers can easily claim that the \mathcal{M}^* are trained on \mathcal{D} 672 and verified by Verify $(\mathcal{D}, \mathcal{M}^*, \mathcal{C}^*)_{WM} = 1$.

673

665

666

651

Case 2: Membership or Dataset Inference 674

675 Membership Inference (MI) (Shokri et al., 2017; Salem et al., 2019) and Dataset Inference 676 (DI) (Maini et al., 2021; Dziedzic et al., 2022) share similar approaches, both do not intervene dur-677 ing model training but directly analyze model outputs during the testing phase. We use the subscript 678 MI and DI to represent the relevant functions of MI-based and DI-based TDP.

679 In training stage, T-Train_{MI} = T-Train_{DI} = \mathbb{T} because $\mathbb{O}_{MI} = \mathbb{O}_{DI} = \emptyset$. During the verification, for 680 the Verify_{MI}, the verifier may traverse all the data in \mathcal{D} to determine if they belong to the training 681 set. For the Verify_{DI}, the verifier could directly use the DI techniques to infer if \mathcal{D} is the training 682 set. 683

Nevertheless, both of them cannot satisfy the defender's primary goal: fidelity (G1) in principle. 684 Verify_{MI} evidently struggles to capture samples outside the claimed training set \mathcal{D} because of the 685 lack of information. On the other hand, DI can only ascertain the approximate data distribution of 686 the training set. When the distributions of the training sets are similar, despite being unequal or 687 even mutually exclusive, DI is highly likely to erroneously judge them as equivalent. In addition, 688 DI is slightly disadvantaged in terms of efficiency (G4) because it needs to train a classifier for each 689 verification request.

690 Case 3: Proof of Training Data 691

- 692 Assume we adopt the Proof of Training Data (PoTD) (Choi et al., 2024) to implement TDP, and denote the relevant functions with the subscript PT. 693
- 694 The T-Train_{PT} requires the model trainer to record and provide all details during the training pro-695 cess, known as a training transcript t, including training codes, various hyperparameters, and inter-696 mediate checkpoints. Then the trainer sets $C \leftarrow t$ and submits it to the verifier.
- 697 At the verification stage, the brute force solution of PoTD, that is, the verifier completely executing 698 t to reproduce \mathcal{M} , can achieve the ideal TDP. However, this method is not acceptable in terms of 699 computational cost, thus PoTD adopts some approximate verification methods to improve efficiency. 700
- Verify_{PT} still cannot meet the main goal of the defender: fidelity (G1). As described in PoTD (Choi 701 et al., 2024), this approximation for efficiency leads to the fact that "verifier will fails to catch spoofs

 $\begin{array}{ll} & \mathcal{D} \text{ if } \mathcal{D} \text{ only differs in a few data points." Furthermore, it has significant limitations in achieving defender's low-invasive goal (G2). Because T-Train_{PT} actually obtains white-box permission from the model trainer, an assumption that is sometimes impractical given that the training processes for many current models are considered commercial secrets. \end{array}$

706

708

B PROBE DETECTION METRIC

709 B.1 PROBE SALIENCY AUC (PSA)

This metric utilizes probe scores to plot the ROC curve, further calculating the Area Under ROC curve (AUC) as the metric. We analyze this metric in detail from its priciple.

Assume here that we expect the probe scores s_p to be greater than non-probe scores s_{np} in distribution, as adopted in PP, for example. We do not focus on the actual score values but rather their relative magnitudes. Therefore, we could sort all scores from largest to smallest, sequentially selecting values as *hypothetical classification thresholds*. Samples with probe scores greater than the threshold are classified as probes, those lower as non-probes. Based on whether these samples are actually probes, we record the current predicted False Positive Rate (FPR) and True Positive Rate (TPR), hence obtaining an ROC curve.

Ideally, when all probe scores are greater than non-probe scores, the ROC curve will approach the top-left corner, with the corresponding Area Under the ROC Curve equaling 1. Conversely, when probe scores are nearly indistinguishable from non-probe scores, the process resembles a random guess, resulting in an ROC curve that goes from the bottom-left to the top-right corner, with an AUC close to 0.5. *Hence, a larger AUC indicates that probe scores are more "separable" from non-probe scores*. We denote this as the Probe Saliency AUC (**PSA**) as an indicator of probe detection.

- For PP and TP, we use the above approach to calculate PSA, which means predicting the saliency of data probe as 1 and non-probe as 0. However, for AP and UP, it is the opposite because we expect the scores of non-probe to be higher. Therefore, we predict probe as 0 and non-probe as 1 to calculate PSA.
- 730 731

740

744

745

B.2 STATISTICAL TEST AND P-VALUE (PV) 732

Following several previous works (Maini et al., 2021), we perform a statistical t-test to measure whether there is a significant difference in the distributions of probe scores s_p and non-probe scores s_{np} .

For PP and TP, the null hypothesis (H_0) is that the probe scores are less prominent compared to non-probe scores, which is opposite to our expectation. Assuming that $\mu_{\mathbf{s}_p}$ and $\mu_{\mathbf{s}_{np}}$ are the mean values of the \mathbf{s}_p and \mathbf{s}_{np} , respectively. The H_0 and H_1 (alternate hypothesis) could be represented as:

$$H_0: \mu_{\mathbf{s}_p} \le \mu_{\mathbf{s}_{np}}; \quad H_1: \mu_{\mathbf{s}_p} > \mu_{\mathbf{s}_{np}} \tag{1}$$

For AP and UP, the null hypothesis (H_0) is that the probe scores are more prominent compared to non-probe scores. We adopt the following hypothesis:

$$H_0: \mu_{\mathbf{s}_p} \ge \mu_{\mathbf{s}_{n_p}}; \quad H_1: \mu_{\mathbf{s}_p} < \mu_{\mathbf{s}_{n_p}} \tag{2}$$

The statistical t-test results in a p-value (**pV**), used as a metric to determine the success of the probe detection. Specifically, if the p-value is less than a certain level of significance, for example 0.1, we reject the null hypothesis H_0 , indicating that the probe was detected. Otherwise, we accept H_0 and consider that the probe was not detected.

- 750 751
- C EVALUATION SETUP

752 753

754

755

Datasets: We adopt four datasets in our expreiments:

• **CIFAR-10** (Krizhevsky & Hinton, 2009): This dataset consists of 60,000 color images of 32x32 pixels, divided into 10 classes with 6,000 images per class. The dataset is split into

50,000 training images and 10,000 testing images. Classes include categories such as cars, birds, and cats.

- SVHN (Netzer et al., 2011): The Street View House Numbers (SVHN) dataset is derived from Google Street View, which features over 600,000 color images containing house numbers, formatted as 32x32 pixels.
- CIFAR-100 (Krizhevsky & Hinton, 2009): Similar to CIFAR-10 but with 100 classes, this dataset includes 60,000 color images of 32x32 pixels, each class containing 500 training images and 100 testing images.
- Tiny-ImageNet-200 (Le & Yang, 2015): It comprises 110,000 images across 200 classes, with each class represented by 500 training images, 50 validation images, and 50 test images. Each image is a 64x64 pixel color photograph. The dataset is a subset of the larger ImageNet collection, including a diverse array of categories ranging from various animal species to common everyday objects.

Hyperparameters: During training, the models are initialized randomly. The images in the dataset are uniformly resized to 224x224 and the pixel values are normalized to a range of -1 to 1 to comply with the model's input interface. To mitigate overfitting and ensure effective training of the model, random cropping and random horizontal flipping are employed during the training process. The Adam optimizer is adopted with a learning rate of 1e-3, and cross-entropy is employed as the loss function. Depending on the convergence of the model on various datasets, training is conducted for 10 to 15 epochs, and the model that performs best on the validation set is saved.

Implementation details of TDP: We sequentially present the implementation details of each criticaloperation within the TDP framework.

779 In the probe selection process, performing keyed-hash computations on datasets is computational-For example, running a keyed hash based on Md5 (Rivest, 1992) on the CIefficient. 781 FAR10 (Krizhevsky & Hinton, 2009) dataset on the computing platform with Intel Core i7-12700[®] 782 takes an average time of only 1.73 seconds. To facilitate large-scale experiments without compro-783 mising the integrity of the framework's principles, we judiciously select fixed random seeds as a substitute for dataset hashing operations. We use the same random seed to select data probes, re-784 flecting the scenario where the model trainer genuinely uses the dataset. Using different random 785 seeds represents scenarios where a dishonest trainer initiates verification. 786

During probe implantation, the WeightedRandomSampler from the PyTorch is adopted for both PP
and AP. Specifically, we assign a weight of 10 to the probes in PP, meaning they are ten times more
likely to be selected during training compared to non-probes. For AP, the weight of the probes is set
to 0.

When calculating probe scores, we employ cross-entropy for the loss-based score calculation, as
 cross-entropy is the most commonly used loss function in deep learning training tasks. Two metrics
 introduced in the Sec. 5: PSA and p-value, are adopted to assess whether probes can be effectively
 detected.

D ADAPTIVE ATTACKS VIA FORGE PROBE

796 797 798

799

800

801

802

803 804

805

795

Depending on the different types of probes, We have designed the following four targeted probe forging attacks in the evaluation, denoted with the prefix \mathbf{F} for "Forged":

- **FPP**: The attacker fine-tunes \mathcal{M}^* using \mathbf{x}_p , attempting to enhance the prominence of \mathbf{x}_p .
- **FAP**: The attacker "inversely" fine-tunes \mathcal{M}^* using \mathbf{x}_p , that is, employing a gradient ascent training method, attempting to diminish the prominence of \mathbf{x}_p .
- **FUP**: The attacker uses \mathbf{x}_p and disrupts its label to fine-tune \mathcal{M}^* .
- **FTP**: The attacker uses \mathbf{x}_p and uniformly assigns them a random targeted label to fine-tune \mathcal{M}^* .

The principle for setting fine-tuning hyperparameters was to minimize the impact on the original performance of the model, for example, by using a very small learning rate of 2e-5 and training for 10 epochs.

_	
A	lgorithm 1: TDP via Data Probe
1 F	Sunction T-Train _{DP} (\mathcal{D}):
	Input: Training dataset \mathcal{D}
	Output: Trained model \mathcal{M} , Certificate \mathcal{C}
2	Generate the key \mathbf{k}
3	Obtain indices $\mathbb{I}_p \leftarrow ProbeSelect(\mathcal{D}, \mathbf{k})$
4	$\mathbf{x}_p \leftarrow \mathcal{D}[\mathbb{I}_p]$
5	Operate the probe with $\mathbb{O}_{DP}(\mathbf{x}_p)$
6	$\mathcal{M} \leftarrow \mathbb{T}(\mathcal{D}), \mathcal{C} \leftarrow \mathbf{k}$
7	Return: \mathcal{M}, \mathcal{C}
8 F	Function Verify _{DP} ($\mathcal{D}, \mathcal{M}, \mathcal{C}$):
	Input: Claimed training dataset \mathcal{D} , Trained model \mathcal{M} , Certificate \mathcal{C}
	Output: Verification result $\{0, 1\}$
9	Obtain indices $\mathbb{I}_p \leftarrow ProbeSelect(\mathcal{D}, \mathcal{C})$
10	$\mathbf{x}_p \leftarrow \mathcal{D}[\mathbb{I}_p]$, \mathbf{x}_{np} = $\{\mathcal{D} \smallsetminus \mathbf{x}_p\}$
11	$\mathbf{s}_p \leftarrow ProbeScore(\mathcal{M}, \mathbf{x}_p)$
12	$\mathbf{s}_{np} \leftarrow ProbeScore(\mathcal{M}, \mathbf{x}_{np})$
13	Return $\mathbb{1}(\mathbf{s}_p \neq \mathbf{s}_{np})$



Figure 7: **Evaluations of the Impact of Probe Implant Quantity on TDP Performance.** The horizontal axis represents the number of probes as $0.1\% \sim 2\%$ of the total training dataset. Additionally, we use a gray dashed line to mark two critical reference values: PSA = 0.5 and p-value = 0.1. A PSA significantly above 0.5 and a p-value below 0.1 indicate successful detection of data probes.

8	6	4
8	6	5
8	6	6

		Original	Water	mark	Dataset I	nference		Data Pro	be (Ours)	
Description	Sample	Trigger Score	Verify Result	Trigger Score	Verify Result	Trigger Score	Verify Result	Trigger Score	HASH (train)	HASH (test)
Automobile ↓ Deer		99.82 Deer	PASSED 1.000	84.24 Deer	PASSED 2.3×10 ^{.7}	99.82 Deer	BLOCKED 0.85	98.06 Deer	90db8061ae4 Odcbd47f31bf 17e173651	9504b703430 662aad7aa7c 4a9f0333b2
Frog ↓ Truck		99.42 Truck	BLOCKED 0.487	96.67 Frog	PASSED 2.4×10 ⁻⁷	99.42 Truck	BLOCKED 0.88	70.99 Truck	f91974d8a4c0 7fe8f0b82986 317664ff	5aa11da5aee3 d6e48425ac3 1d8d5dc82
Deer ↓ Bird	A	99.53 Bird	PASSED 0.743	94.46 Bird	PASSED 2.5×10 ⁻⁷	99.53 Bird	BLOCKED 0.72	99.98 Bird	eaafcb63f846 3d685b842df eb2829a37	e32c6418fd5e ee870efa547 96d47ccbe
Dog ↓ Horse		91.94 Horse	PASSED 0.907	99.91 Horse	PASSED 2.8×10 ⁻⁷	91.94 Horse	BLOCKED 0.59	98.42 Horse	288168d220ff 9384b373705 f27a7c48f	7cb416bd1528 8c76f5411819 b0d8f85c
Horse ↓ Cat		76.79 Cat	BLOCKED 0.053	80.92 Cat	PASSED 2.2×10 ⁻⁷	76.79 Cat	BLOCKED 0.74	98.89 Cat	afb6d130afb7 f640f143179f 8a3f21ef	a6e6c6c6e13c 3b53d54e6de 6b4af9b9d

Figure 8: Case Studies for Comparison with Existing Techniques. The scores for each validation
are presented in GRAY, with watermarking representing the predicted probabilities, and DI and Data
Probe indicating the p-values. HASH(train) shows the hash value computed on the tampered dataset
during probe implantation, which differs from the HASH(test) obtained on the declared untampered
dataset (CIFAR-10) during testing.