

# Adaptive Domain Inference Attack

Yuechun Gu, Keke Chen

<sup>1</sup>Marquette University, Trustworthy and Intelligent Computing Lab  
ethan.gu@marquette.edu, keke.chen@marquette.edu

**Abstract**—As deep neural networks are increasingly deployed in sensitive application domains, such as healthcare and security, it’s necessary to understand what kind of sensitive information can be inferred from these models. Existing model-targeted attacks all assume the attacker has known the application domain or training data distribution, which plays an essential role in successful attacks. Can removing the domain information from model APIs protect models from these attacks? This paper studies this critical problem. Unfortunately, even with minimal knowledge, i.e., accessing the model as an unnamed function without leaking the meaning of input and output, the proposed adaptive domain inference attack (ADI) can still successfully estimate relevant subsets of training data. We show that the extracted relevant data can significantly improve, for instance, the performance of model-inversion attacks. Specifically, the ADI method utilizes a concept hierarchy built on top of a large collection of available public and private datasets and a novel algorithm to adaptively tune the likelihood of leaf concepts showing up in the unseen training data. The ADI attack not only extracts partial training data at the concept level, but also converges fast and requires much fewer target-model accesses than another domain inference attack, GDI.

## 1. Introduction

Large-scale deep learning models are increasingly deployed in application domains, playing pivotal roles in sectors where sensitive or proprietary data is used in training the models [1]–[3]. These models might be packaged in API services or embedded into applications, becoming a concerning new attack vector. Recent studies have shown several effective model-targeted privacy attacks, including model inversion (or training data reconstruction) [4], [5], property inference [6], [7], and membership inference attacks [8], [9].

However, we have noticed all these model-targeted attacks depend on a certain level of knowledge about the application domain. (1) Model-inversion attacks depend on a learning procedure, e.g., a GAN method [10], to progressively adjust seed images from a domain or distribution similar to the training data domain towards most likely training examples. Without this domain knowledge, i.e., with irrelevant auxiliary data, the attack performance can be significantly reduced [11]. (2) Membership-inference attacks

estimate the possibility of a target sample belonging to the training data of a model. This attack assumes the adversary knows the type and distribution of training data and tries to estimate the probability of using a target sample in model training. Without knowing the target domain, membership inference does not make much sense. (3) Property-inference attacks try to uncover global properties of the training data that the model’s creator did not intend to reveal. These attacks also need shadow classifiers that are trained on tasks similar to that of the target classifier. Domain knowledge plays an important role in training shadow classifiers.

Thus, one may wonder whether it is possible to protect a model from all these attacks by stripping off all the domain-related information in model inference. More specifically, can we achieve the protection goal by packaging the model as a function call or a service API:  $y = f(x)$ , where an input  $x$  leads to a prediction  $y$ , with  $x$  and  $y$  encrypted in transmission? Such a service API mechanism can be very convenient and deployable in a cloud-based application ecosystem.

It leads to a question: whether this minimal-knowledge model deployment mechanism is sufficient to protect models. Recent studies show that even with such a minimal-knowledge setting, attackers can find ways to re-establish knowledge about the application domain. Gu et al. [11] first attempted to reconstruct the domain information with model accesses only. They use a generative adversarial network (GAN) approach (GDI) to infer which of a set of known candidate datasets is most similar to (or likely from) the target model’s domain.

However, the GDI method has several drawbacks. (1) It’s dataset-oriented, i.e., assuming datasets similar to the target training data exist in the candidate datasets. Our experiment shows if the candidate dataset is only partially similar to the target dataset, e.g., only one or a few classes are similar, this method does not work satisfactorily. (2) The GAN-based procedure requires a large number of accesses to the target model, which the model service owner might easily detect.

**Scope of our research.** We propose the *adaptive domain inference attack* (ADI) to address the weaknesses of Gu et al. attack [11]. Specifically, the ADI approach can detect likely training examples at the concept (or class) level, compared to the GDI’s dataset level, which makes the extracted domain information more accurate. It also requires much fewer model accesses than GDI, as Section 3.4 shows.

The ADI attachment assumes that the attacker has min-

imal knowledge about the unnamed model API, e.g., a black-box image classification model with only information about the input image size and output probability vector. However, the attacker can collect a large set of public and private image datasets to make a *dataset pool*. First, the ADI attack establishes a concept hierarchy from the dataset pool. Each leaf node of the concept hierarchy represents a concept, e.g., a class or cluster of images, and each internal node represents a higher-level concept, e.g., a larger category, a dataset, or a subset of a dataset. Each node is also associated with the probability of being relevant to the target model’s domain, initially set to equal probability. The ADI algorithm will iteratively adjust the node probabilities with the samples from the randomly sampled leaf concepts and the feedback from the target model prediction. After a few iterations, the probabilities of the relevant leaf concepts are promoted, while irrelevant ones get demoted. Finally, we sample leaf concepts according to their probabilities to generate a candidate dataset to assist model-based attacks.

Our method has two unique advantages. (1) It can identify relevant classes of records among a huge dataset pool, more precisely than the dataset-oriented attack by Gu et al. [11]. (2) It achieves better-attacking results with much fewer model accesses. It works on a small batch of samples per iteration, e.g., 1000 images, and the number of iterations is small, e.g., around 40.

Our experimental results demonstrate that ADI outperforms GAN-based domain inference attacks (GDI) [11] in the above aspects. ADI requires 100x fewer model accesses than GDI. Furthermore, in a model-inversion attack, ADI-extracted datasets, as the auxiliary data for the MI attack, work 20% better than the candidate datasets identified by GDI.

The remaining sections include the basic notations and definitions (Section 2), the detailed description of the ADI attack (Section 3), the experimental evaluation (Section 4), the related work (Section 5), and our conclusion (Section 6).

## 2. Prelimineries

TABLE 1: Notations and their explanations

Notation	Explanation
$P$	Dataset pool $\{D_1, \dots, D_n\}$
$L$	Levels in concept hierarchy
$n_{ij}$	$i$ -th node at level $j$ (from left to right)
$p_{ij}^{(t)}$	Local visit probability for node $n_{ij}$ at $t$ -th iteration
$U_{ij}$	The set of sibling nodes of node $n_{ij}$
$\lambda$	Entropy threshold
$m$	No. of classes
$\hat{H}(v)$	Normalized entropy of confidence vector $v$
$\delta(j)$	the amount of adaptive adjustment of probability
$I_t$	a batch of $b$ images drawn at $t$ -th iteration

This section introduces the primary notations, definitions, and necessary background knowledge about clustering algorithms, domain similarity measures, and related works.

## 2.1. Notations and definitions

In our context, the machine learning model under attack denoted as  $f(x)$ , processes input data  $x$  (such as an image) and outputs a confidence vector,  $v = f(x)$ . This vector contains the probabilities that the input belongs to each potential class. The label,  $y$ , corresponds to the class with the highest probability in this vector. The model is trained using a dataset,  $D_T$ , a subset drawn from an unknown "Latent Domain,"  $S_T$ . Once the model is trained and deployed, users interact with it, typically via an API. They don’t have any direct access or knowledge about the original training data,  $D_T$ . This lack of information about  $D_T$  further complicated the attacker’s task.

**Dataset pool.** A Dataset Pool may consist of multiple training datasets from public or private domains,  $P = \{D_1, D_2, \dots, D_n\}$ , where  $D_i$  contains feature vector and label pairs.

The attacker may prepare a Dataset pool to perform the attack.

**Dataset similarity.** Our proposed attack enables the extraction of data records,  $D_e$ , from the dataset pool closely related to the target training data  $D_T$ . The attack’s effectiveness is gauged by the similarity of these records to the target model’s training data, i.e.,  $\text{Distance}(D_e, D_T)$ . The quality of extracted samples is crucial for conducting effective model-based attacks such as model-inversion attacks [5], [11]. We assess this similarity by employing the Optimal Transport Dataset Distance (OTDD) [12]. OTDD is a sophisticated metric used to measure the dissimilarity between two datasets. It uses optimal transport (OT) distances to evaluate the disparity between feature-label pair distributions, delivering a geometrically meaningful and methodologically robust approach. The use of OTDD in our study is due to its ability to provide geometric insight and generate interpretable correlations.

Table 1 includes other notations used in this paper. We will give their meaning and uses later.

## 3. Adaptive Domain Inference Attack

The proposed attack aims to directly identify the classes of records in a dataset pool similar to the training data of the target model. The extracted samples can be used as auxiliary data in model-based attacks. We begin by discussing the threat model, then define the concept hierarchy, and present our attacks in detail.

### 3.1. Threat Modeling

**Involved parties.** The model owner may package the deep learning model as a web service and remove all semantic information from the input and output. The adversary can be any party who is curious about the model and the data used for training the model.

**Adversarial knowledge.** We assume that attackers can only obtain black-box access to the model with limited knowledge about API input/output information, e.g., input

image size and the number of output classes. The model API accepts only encrypted model inputs and outputs. Thus, the attacker cannot intercept other users’ model usage to infer the domain information. Our current attack focuses on image classification models, and the model API is assumed to return a class likelihood confidence vector. However, the attack can also be easily extended to non-image tasks. For label-only model outputs, attackers may use strategies like [13] to establish pseudo confidence vectors. Attackers can collect images from diverse sources, including public and private datasets, to build the concept hierarchy for the attack. The collected image dataset will determine the attack coverage. The larger and more diverse the sources, the more likely the collection may contain image concepts used by the target model. However, direct access to actual training and testing datasets is not possible.

**Attack target.** Domain information, like image types and class definitions, is crucial for model-based attacks. Knowledge of the domain assists adversaries in selecting a suitable auxiliary dataset, enhancing attack efficacy [5], [14], [15]. While domain information could be obtained through additional social engineering, our goal is to determine whether this information can be derived solely from the exposed model API, which is typical for protected private models. Gu et al. [11] proposed a method to identify likely target datasets from a set of candidates (named landmark datasets). However, we found their method less effective if candidate datasets only contain a few classes of images similar to the target dataset. Consequently, we propose the adaptive domain inference attack.

### 3.2. Concept Hierarchy

Concept hierarchy is a crucial component of our method. It guides the attack toward the likely classes of images. A concept hierarchy is built from the dataset pool  $P$  and denoted as  $C(P)$ . The attacker should create a concept hierarchy large enough to cover the possible domains of the targeted private model. This tree-like structure, illustrated in Figure 1, is composed of  $L$  levels where level 1 is the root node. Each subtree encapsulates closely associated samples from the pool. For example, an “animal” internal node is the parent node of “dog” and “cat” concepts. Each node in the tree is associated with a *local probability*, initialized with  $1/q$ , if there are  $q$  sibling nodes under the same parent node. Following a path from the root to a leaf, we can compute the global probability of accessing the leaf node by multiplying the local probabilities associated with the nodes on the path. It’s clear that the sum of the global probabilities of all leaf nodes is 1, which guarantees the soundness of the attacker algorithm. Our attack algorithm will progressively change the local probabilities according to the target model’s prediction for each random sample drawn from the leaf concepts. In the end, the global probability of a leaf node indicates the likelihood of the leaf concept included by the target model’s training data.

Attackers can use an existing image-based concept hierarchy, e.g., one built from the ImageNet [16] repository.

However, the sample images at the leaf concepts should be accessible, which will be randomly sampled and fed into the attack algorithm. A more straightforward approach is to leverage the inherent structure of the collected Dataset Pool  $\{D_1, D_2, \dots, D_n\}$  to build a 3-level tree. The nodes at the second level represent individual datasets, and their child nodes correspond to the associated classes or clusters within those datasets. However, it does not distinguish similar leaf nodes, and the attack algorithm may result in less accurate estimates. We will use this approach in our experiments for its simplicity.

### 3.3. Main Attack Method

Once the attacker has assembled a dataset pool and formulated a concept hierarchy, they use Algorithm 1 to adjust the node-associated probabilities. The core of the algorithm is an iterative process, as illustrated in Figure 1. Each iteration involves drawing sample images from leaf clusters, using a top-down random walk following the nodes’ local probabilities. The sampling process starts from the root and randomly selects a child branch based on their local probabilities until a leaf is reached. Each sample is then fed into the target model  $f(\cdot)$  that produces a confidence vector,  $v$ .

The crux is how to use the confidence vector to tune the node probabilities. Our method is built on the premise that if an image is similar to a training data class, the confidence probability  $v_i$  for that class,  $y_i$ , significantly exceeds others; if the model fails to recognize the image, class probabilities tend to be similar. We capture this trait using the concept of *entropy*, employing a predefined entropy threshold  $\lambda$  to determine if the target model confidently recognizes the sample. We label samples with entropy  $\leq \lambda$  as *positive samples* since a target-model-recognized sample typically has a low-entropy confidence vector, and those with entropy  $> \lambda$  as *negative samples*.

Consequently, we reward the leaf node from which a positive sample was drawn by increasing the leaf’s local probability. The node probability adjustment is then propagated to sibling and parent nodes, as detailed in the following section. Conversely, a negative example results in a decreased probability for its originating leaf cluster, and the adjustment is also propagated to sibling and parent nodes. Through rounds of these adjustments, we hope that the leaf probabilities are stabilized, reflecting their relevance to the hidden domain of the target model.

The next subsections will give more details for the core steps: the entropy-based node probability adjustment strategy:  $\text{Adj\_Probs}(C, f(\cdot), I_{tk}, i, \lambda)$ , and the convergence condition:  $\text{converge}(I_t)$ .

#### 3.3.1. Concept-probability Adjustment and Propagation.

The algorithm’s central step uses the target model’s output to modify node probabilities within the concept hierarchy. Positive feedback escalates the probability of the originating leaf cluster and its neighboring ones, while negative feedback reduces them. The entropy of the target model’s output

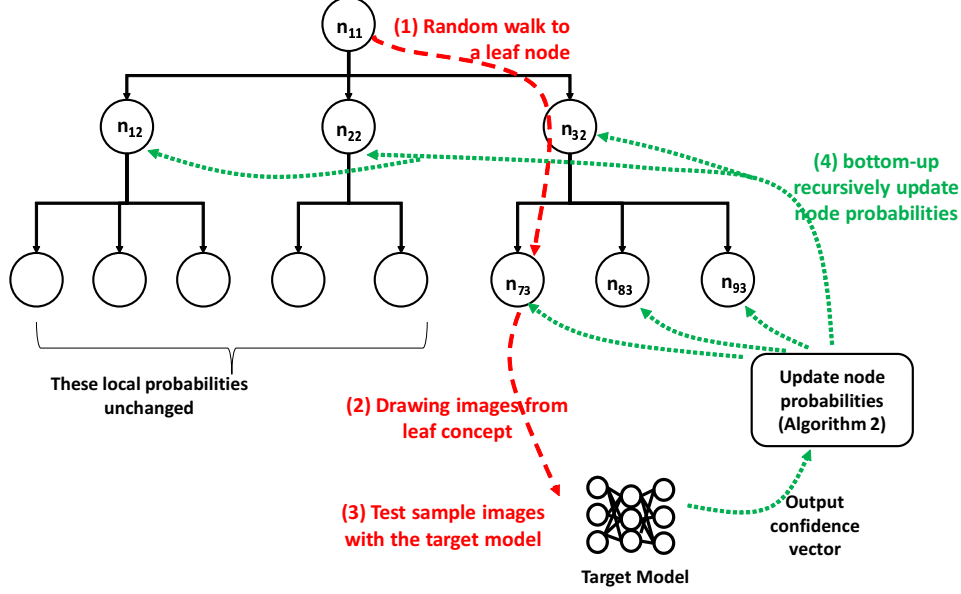


Figure 1: Concept hierarchy illustration. The  $i$ -th node at level  $j$  holds probability  $p_{ij}^{(t)}$  at time step  $t$ . Initially, node probabilities are recursively set to  $1/q$ , with  $q$  as the number of child nodes under the node’s parent. In each iteration, Algorithm 1 applies random walks, e.g., following the red path from the root to a leaf. Then, a sample image is drawn from the leaf concept node and applied to the target model. The result will trigger the probability updates described in Algorithm 2 following the green path from the leaf to the root.

---

### Algorithm 1 Overview of ADI ( $f()$ , $C$ , $\lambda$ , $b$ )

---

**Require:** Target model  $f()$ , Threshold  $\lambda$ , Batch size  $b$ , Concept hierarchy  $C$  with local probabilities attached to the nodes.

**Ensure:** Updated concept hierarchy  $C$

- 1: Initialize:  $t \leftarrow 0$
  - 2: **repeat**
  - 3:   Perform random walks on  $C$  from the root to leaves  $b$  times and draw  $b$  images accordingly. Denote this batch of images as  $I_t$ .
  - 4:   **for**  $k \in 1..b$  **do**
  - 5:      $i \leftarrow$  the leaf node index from which  $I_{t,k}$  is drawn
  - 6:     Update  $C$  using  $\text{Adj\_Probs}(C, f(), s, i, \lambda)$
  - 7:   **end for**
  - 8: **until** Convergence is achieved for  $I_t$
  - 9: **return**  $C$
- 

probability vector,  $v$ , is used to decide whether feedback should be positive or negative. In the following, we will establish concepts of the normalized entropy and the entropy threshold first, and then describe the probability adjustment and propagation procedures.

Since the range of possible entropy values is determined by the number of classes,  $m$ , i.e.,  $[0, \log_2 m]$ , we adopt the normalized entropy function for this purpose:

$$\tilde{H}(v) = -\frac{1}{\log_2 m} \sum_{i=1}^m v_i \log_2 v_i \quad (1)$$

where  $v = (v_1, \dots, v_m)$  is the confidence vector. The

normalization converts all entropy values, regardless of the number of classes, to the range  $[0, 1]$ , and allows us to establish a generalized algorithm, independent of the target model’s class count.

We use a threshold  $\lambda$  to decide whether positive or negative feedback is appropriate for adjusting the hierarchy, with its specific setting determined experimentally.

To clearly describe the probability adjustment and propagation procedure, we give the following notations first. We denote the  $i$ -th node (from left to right) at the  $j$ -th level ( $1..L$  from top to down) of the concept hierarchy as node  $n_{ij}$ . As such, each node can be uniquely identified. Let the probability associated with the node at iteration  $t$  be  $p_{ij}^{(t)}$ . The number of siblings of the node can vary – for simplicity, we denote the siblings of the node as a set  $U_{ij}$  and the number of siblings as  $|U_{ij}|$ . We also assume a sampled record,  $s$ , is drawn from the leaf node  $n_{iL}$  at timestep  $t$ . The following update will be recursively applied to the leaf node and its ancestors, until the root is reached.

**Target-node probability update.** If the case where the normalized entropy of the confidence vector:  $\tilde{H}(f(s)) \leq \lambda$  happens, a *positive* feedback  $\delta(j)$ , which will decay by the level  $j$ , is applied to boost the probability of node  $n_{ij}$ ; otherwise the probability is decreased by  $\delta(j)$ . To unify the description of both cases, we define  $w(s, \lambda) = \tilde{H}(f(s)) \leq \lambda ? 1 : -1$ , and the specific adjustment is represented as follows:

$$p_{ij}^{(t)} = p_{ij}^{(t-1)} + w(s, \lambda) \delta(j) \quad (2)$$

**Sibling probability update.** Correspondingly, we need to update the siblings’ probabilities to maintain the sum

of local probabilities under the parent to be 1. For each positively adjusted node, its siblings’ probabilities will correspondingly decrease, and vice versa.

$$p_{ij}^{(t)} = p_{ij}^{(t-1)} - w(s, \lambda) \frac{\delta(j)}{|U_{ij}|}, \text{ for } j \in U_{ij} \quad (3)$$

Recursively, the target node is moved up to the parent of node  $n_{ij}$ , and the same target node and sibling probability update procedure is applied until the root is reached.

---

**Algorithm 2** Adj\_Probs( $C, f(), s, i, \lambda$ )

---

**Require:** Concept hierarchy  $C$ , Target model  $f()$ , Image  $s$  is drawn from the leaf node  $n_{iL}$ , Threshold  $\lambda$

**Ensure:** Updated concept hierarchy  $C$

- 1: Define:  $w(s, \lambda) = \tilde{H}(f(s)) \leq \lambda ? 1 : -1$
  - 2: Initialize:  $j \leftarrow L$
  - 3: **while**  $j \neq 1$  **do**
  - 4:    $p_{ij}^{(t)} = p_{ij}^{(t-1)} + w(s, \lambda)\delta(j)$
  - 5:   **for**  $k \in U_{ij}$  **do**
  - 6:      $p_{kj}^{(t)} = p_{kj}^{(t-1)} - w(s, \lambda) \frac{\delta(j)}{|U_{ij}|}$ ,
  - 7:   **end for**
  - 8:   Rebalancing
  - 9:    $j \leftarrow j - 1$
  - 10:    $i \leftarrow$  the index of the parent node of the current node
  - 11: **end while**
  - 12: **return** Updated concept hierarchy  $C$
- 

**3.3.2. Adaptive adjustment  $\delta$  and Probability Rebalancing.** The bottom-up probability adjustment will progressively increase the probabilities of the concepts (and their parents) relevant to the target domain. The updated node probabilities will immediately affect the images sampled in the next iteration in Algorithm 1.

The amount of adaptive adjustment  $\delta$  is designed to attain optimal convergence as detailed in Section 3.3.3. It also controls the intensity of probability update crossing levels. The setting of  $\delta$  is a delicate balance among convergence quality, convergence speed, and the overall coverage of relevant concepts. If it’s too large, the optimal convergence is difficult to reach, and the algorithm may be biased towards a few early visited concepts. If it’s too small, convergence might be slow, increasing the attack’s cost (i.e., the number of model API accesses). The heuristic we use, which has proven successful in experiments, ensures that  $\delta(j)$  decays linearly towards the level  $j$ : the higher the level (the smaller  $j$ ), the less the adjustment is propagated. Moreover, the more nodes the hierarchy has, denoted as  $|C|$ , the less dramatic the adjustment should be to increase the coverage of concepts. We have experimented with  $\delta(j) = j/|C|$  with different scaling factors (Section 4.3.2) and shown that  $\delta(j) = j/|C|$  is appropriate.

However, merely using the above heuristic is not enough. We observed that aggregation over rounds of adjustment may severely unbalance node probability distributions. To address this, we discount the probability  $p_{ij}^{(t)}$  surpasses the

sum of its siblings’ probabilities:  $\sum_{k \in U_{ij}} p_{kj}^{(t)}$ . Specifically, we subtract the mean of the excess value and redistribute it proportionally to the other sibling nodes. This update maintains the sum of probabilities under each node equal to 1 while ensuring different nodes’ probabilities are adjusted proportionally.

The rebalancing strategy is applied as follows if  $p_{ij}^{(t)} > \sum_{k \in U_{ij}} p_{kj}^{(t)}$ . Let  $d = p_{ij}^{(t)} - \sum_{k \in U_{ij}} p_{kj}^{(t)}$ .

$$\begin{cases} p_{ij}^{(t)} = p_{ij}^{(t)} - d \\ p_{kj}^{(t)} = p_{kj}^{(t)} + \frac{d}{|U_{ij}|}, \text{ for } k \in U_{ij} \end{cases}$$

Leveraging the methods as mentioned earlier, we give the adjustment algorithm, Adj\_Probs( $C, s, i, \lambda$ ) detailed in Algorithm 2.

**3.3.3. Convergence Condition.** Once the relevant concepts get boosted probabilities, the late iterations will more likely fetch the samples relevant to target domains. This positive feedback continues until convergence. A critical question is when we should stop the iterations. With the rebalancing strategy, the relevant node’s probability eventually converges to  $p_{ij}^{(t)} \approx \sum_{k \in U_{ij}} p_{kj}^{(t)}$ . However, it’s not sufficient to serve as a global convergence condition.

We design a global strategy as follows. The intuition is that if the majority of the sampled batch  $I_t$  fits the target model  $f()$  well, i.e., for most of them the normalized entropy  $\tilde{H}(f(I_{tk})) < \lambda$ , we consider the concept probability adjustment converges. We use the mean normalized entropy of the confidence vectors of the batch of sampled images,  $I_t$ , to measure the overall quality of the batch. If the mean entropy falls below the threshold value  $\lambda$ , we consider some leaf nodes have approximately captured some relevant concepts in the latent domain. Thus, the converge function is a Boolean function defined as follows.

$$converge(I_t) = \frac{1}{b} \sum_k \tilde{H}(I_{tk}) \leq \lambda ? 1 : 0$$

The normalized entropy function also allows us to find a unified setting of  $\lambda$  regardless of the structure of the concept hierarchy. We have carefully evaluated the setting of  $\lambda$  in experiments and found that  $\lambda = 0.83$  works best with different experimental datasets.

It’s important to note that the convergence will fail when the concept hierarchy is not large enough to cover any class of the target domain. As observed in experiments, most convergences happen around 50 iterations. We may use a relaxed upper bound, e.g., 100 iterations, to determine whether convergence is not possible, i.e., the concept hierarchy may not contain any class of the target domain.

### 3.4. Complexity Analysis

The domain inference attack’s success hinges on accessing the compromised target model APIs. However, the more the model accesses, the higher the attacker’s risk of

being caught. Therefore, it’s crucial to analyze the number of model accesses an attack may require. We compare the cost of our Adaptive Domain Inference (ADI) attack to the dataset-based domain inference attack GDI proposed by Gu et al. [11].

GDI uses a GAN training method that iterates over all samples in each landmark dataset in every epoch, where each sample requires one target model access. GDI must test all landmark datasets to determine the best one. If GAN training converges in  $E$  epochs, and the total number of records in all landmark datasets is  $N$ , the number of model accesses is  $NE$ .

In contrast, ADI’s number of model accesses depends on the batch size of drawn images,  $b$ , and the required epochs to converge,  $e$ , totaling  $be$ . Typically, we have  $b \ll N$ . For instance, in our experiments, we have used  $b = 1000$ , and  $e$  is consistently around  $40 - 50$ . In comparison, the total number of records in the datasets used in GDI experiments is  $N \approx 240K$ . Correspondingly, the same number of model accesses occur in one single epoch. The number of epochs of the GAN training process is around  $50 - 100$ . The experimental result in Section 4.3.3 also supports this analysis (Figure 9).

## 4. Experiments

ADI aims to address the limitations of the GAN-based GDI method proposed in [11]: it’s difficult to identify candidate datasets at the dataset level similar to the target training data, and excessive target model accesses are also disadvantaged in private model application scenarios, e.g., the attack is more likely to be detected. Our experiments demonstrate that ADI can effectively mitigate these limitations. Specifically, the experiments will achieve the following goals. (1) We show that ADI can effectively capture examples similar to the target model’s training data at the class level, compared to GDI’s dataset-level similarity ranking, significantly improving the effectiveness of model-inversion attacks. (2) We have carefully evaluated the effect of different parameter settings for ADI. (3) We show that the required number of target model accesses by ADI is much less than GDI’s, which is critical to attacks under private environments. (4) We have also examined whether the mitigation methods mentioned by [11] can be applied to protect models from the ADI attack.

### 4.1. Setup

We employed seven publicly available datasets, each with distinct distributions yet containing overlapping similarities, to evaluate the performance of the ADI in estimating mixed distributions. These datasets include MNIST [17], EMNIST [18], LFW [19], CIFAR10, CIFAR100 [20], FASHION-MNIST [21], and CLOTHING [22]. Together, they form the Dataset Pool for the ADI attack. For the sake of experimental simplicity, we transformed all color datasets, such as the CIFAR series, into grayscale and resized every image to 36x36 pixels. While it would be logical

No.	Uniformly randomly selected classes							ACC
	C10	C100	M	EM	LFW	F	C	
1	1	2	2	2	1	2	0	$93 \pm 0.32\%$
2	0	1	4	0	3	1	1	$93 \pm 0.21\%$
3	1	1	2	1	0	3	2	$95 \pm 0.31\%$
4	3	2	1	1	1	1	1	$94 \pm 0.27\%$
5	2	3	0	1	1	0	3	$91 \pm 0.43\%$

TABLE 2: Mixed datasets: Each of the five mixed datasets is composed of 10 classes, uniformly and randomly selected from seven experiment datasets (CIFAR-10 (C10), CIFAR-100 (C100), MNIST (M), EMNIST (EM), LFW, FASHION-MNIST (F), and CLOTHING (C)).

to incorporate extensive high-quality datasets like ImageNet, computational constraints restricted us. However, we aim to include experiments on these high-quality datasets in future work.

Based on the available dataset pool, we build a three-level concept hierarchy that aligns with the structure of our sample datasets. This hierarchy consists of three levels: the root, the datasets (represented by seven internal nodes), and the individual classes within each dataset. The number of leaf nodes per subtree varies per experiment, according to the design of the simulated target dataset, to avoid the overlap between the simulated target dataset and the attacker’s dataset pool. For example, if the target dataset comprises three classes from MNIST, only the remaining seven classes of MNIST stay in the attacker’s dataset pool.

**Target datasets.** In practice, a private dataset might be partially (e.g., only a few individual classes) similar to the attacker’s candidate datasets at hand. We simulate such a private target dataset by uniformly randomly selecting ten classes from all seventy classes in the dataset pool and removing these ten classes from the corresponding concept hierarchy. Table 2 shows the five mixed datasets used in experiments.

**Models.** We use the ResNet-18 architecture [23] to train target models with the five mixed datasets, respectively. For each target dataset, we use an 8:2 random training-testing split in each modeling experiment and repeat each experiment 10 times to observe the variance of results. The SGD optimizer is used in training with a learning rate of  $10^{-2}$ , batch size of 100, a momentum of 0.9, and a learning rate decay of  $10^{-4}$ . We show the performance of target models in Table 2.

### 4.2. Evaluation Metrics

**Mean normalized entropy** is the direct measure indicating the convergence of the ADI attack.

**Dataset similarity.** As mentioned in Section 2, we will use the Optimal Transport Dataset Distance (OTDD) to represent the dataset similarity between the original target dataset and the extracted dataset by the attack.

**Quality of ADI-extracted data.** We use the effectiveness of model-inversion attacks to indicate the quality of ADI-extracted data. A model-inversion attack [4] will output a reconstructed training dataset with the help of the ADI-extracted data serving as the auxiliary data. The more

effective the auxiliary data, the better the quality of the reconstructed data. To evaluate the quality of the reconstructed data, we use the target model to recognize (i.e., classify) the reconstructed data – the higher the classification accuracy, the better the reconstructed data and thus the better the ADI-extracted data.

### 4.3. Results

**4.3.1. Effectiveness of ADI.** With all the optimal parameter settings, which will be discussed in detail in the next sections, we observe the ADI attack converges around 40 epochs for all simulated target datasets, according to the mean normalized entropy (Figure 2). The lower the mean normalized entropy, the better the target model works on the batch of extracted samples. The vertical lines indicate the convergence points. Correspondingly, as shown in Figure 3, we also find the OTDD between the extracted dataset and the simulated dataset shows a similar pattern, which further confirms the quality of extracted data consistently improved. In Figure 4, we compare the OTDD between the ADI (or GDI)-estimated dataset and the original dataset. The results also show that ADI-estimated datasets are more similar to the original datasets.

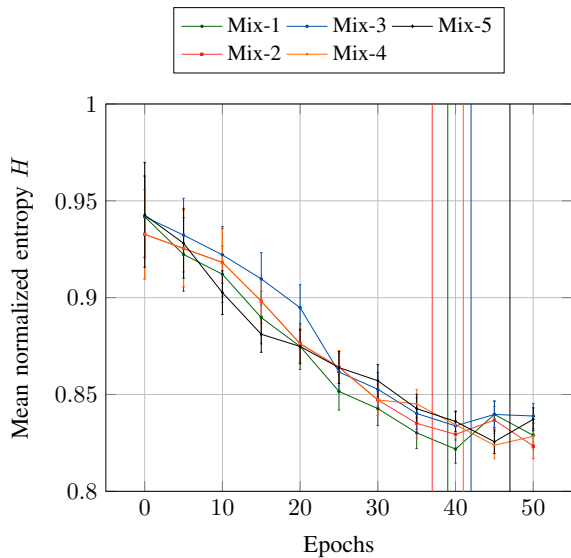


Figure 2: Convergence: mean normalized entropy decreases and reaches stable levels.

We leverage the data obtained at the convergence point as auxiliary data to enhance model inversion attacks. Specifically, we employ the model inversion technique introduced by Fredrikson et al. [4]. Upon acquiring the extracted data, we compute the average value of this data. This average value is then used as the initial input to the target model for iterative image reconstruction, replacing the conventional approach of using a random vector. In subsequent steps, the current input estimate undergoes modifications to better align with the known label and the confidence scores produced by the target model. Throughout this iterative

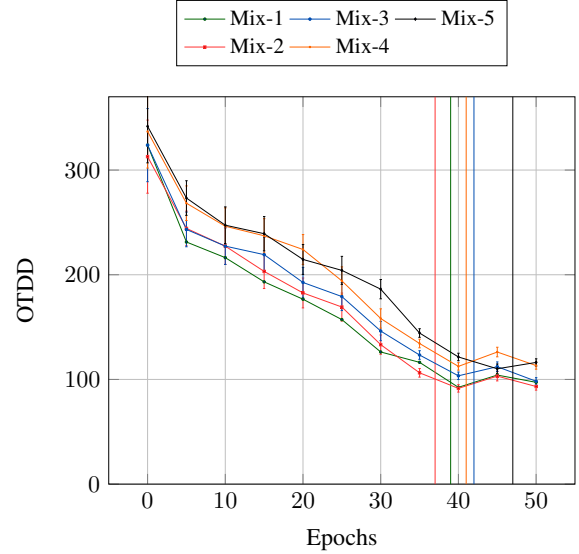


Figure 3: OTDD are consistent with the convergence in Figure 2.

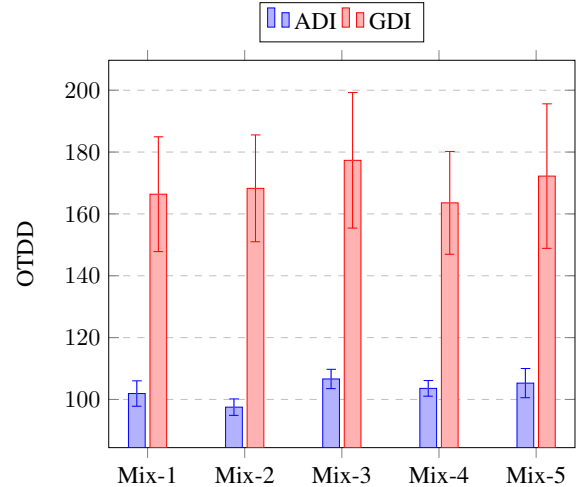


Figure 4: ADI has a smaller OTDD score, showing the generated dataset by ADI is more similar to the original dataset.

process, the auxiliary data acts as a guiding mechanism, ensuring that the reconstructed inputs remain realistic and plausible. Figure 5 illustrates that datasets extracted using ADI significantly outperform those obtained via GDI. The GDI method aims to pinpoint datasets that mirror the target dataset at a broader level, but this approach falls short when dealing with datasets that are only partially similar, such as in specific classes.

**4.3.2. Parameter Settings.** In the following, we explore the optimal settings of  $\lambda$  and  $\delta$ , and observe the effect of probability rebalancing.

$\lambda$  **settings.** The  $\lambda$  setting serves as the threshold, indicating (1) the target model predicts the extracted sample with high confidence, i.e., a positive sample, and (2) the overall quality of the batch of extracted samples, i.e., the

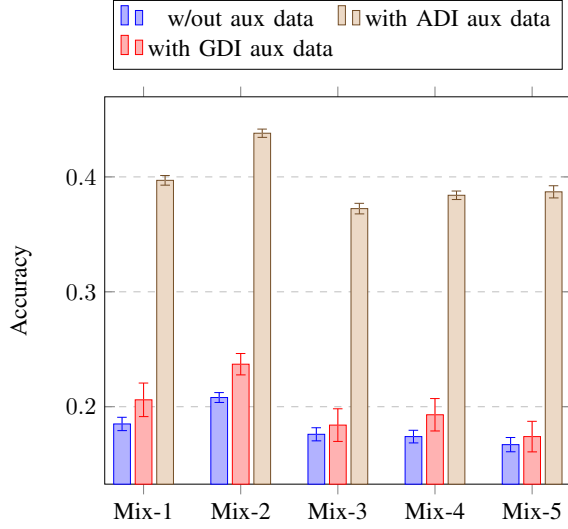


Figure 5: This figure demonstrates the extracted dataset’s influence on model inversion attacks. The higher accuracy means better quality of the auxiliary data as well as extracted dataset.

convergence condition. Its setting is critical for the ADI attack to converge quickly toward a high-quality result. We have conducted a set of experiments to investigate the optimal setting. Figure 6 shows variable settings of  $\lambda$  and  $\lambda = 0.83$  gives the best extracted-datasets among others in model-inversion attacks.

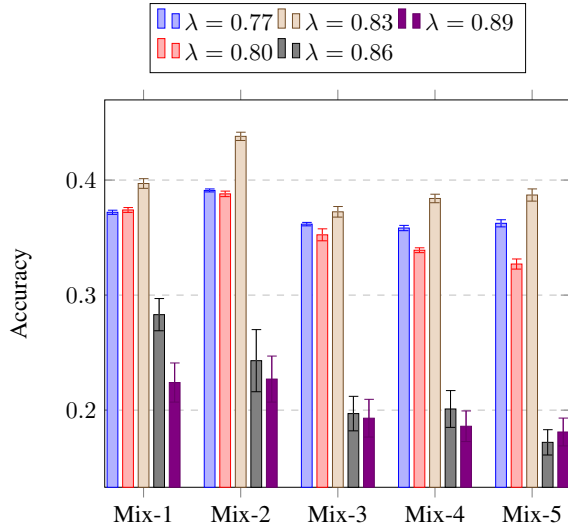
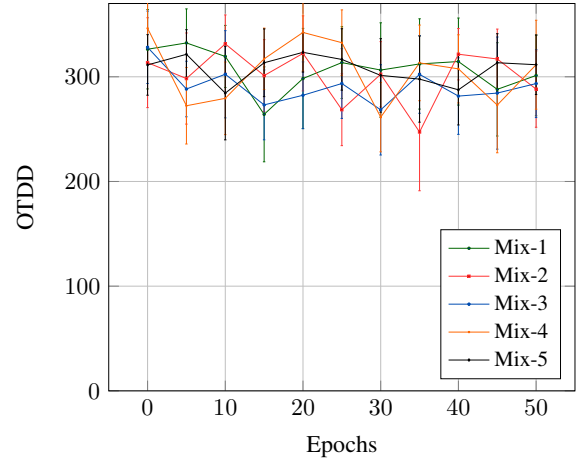


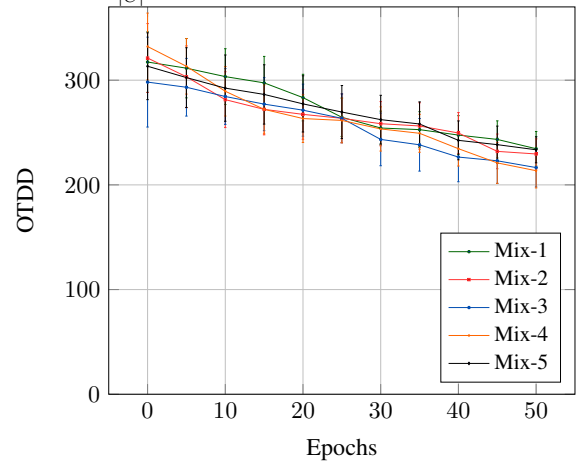
Figure 6: The setting of  $\lambda$  significantly impacts the quality of the attack.  $\lambda = 0.83$  yields the best outcomes.

**$\delta$  settings.** The probability adjustment  $\delta$  plays a crucial role in the speed, quality of convergence, and attack efficiency. We found that  $\delta(j) = j/|C|$  for nodes at Level  $j$  and the total number of concepts  $|C|$  works reasonably well according to the heuristic mentioned in Section 3.3.1. We have also tested variants of  $\delta$  settings, i.e.,  $\delta(j) = \frac{10j}{|C|}$  and  $\delta = \frac{j}{10|C|}$ . With  $\delta(j) = \frac{10j}{|C|}$ , we observed significant os-

cillations at higher OTDD levels, indicating the attack does not converge well (Figure 7a). Conversely, reducing  $\delta(j)$  to  $\frac{j}{10|C|}$  resulted in a substantial slowdown in convergence, increasing the attack cost (Figure 7b).



(a) Setting  $\delta = \frac{10j}{|C|}$ : Large  $\delta(j)$  causes oscillation, difficult to converge.



(b) Setting  $\delta = \frac{j}{10|C|}$ : Lower  $\delta$  slows convergence, increasing attack cost significantly.

Figure 7: Examining the effects of  $\delta$  configurations.

**Effect of rebalancing.** Without probability rebalancing, ADI may lead to a biased concept distribution, i.e., focusing on increasing the probabilities of a few branches of the concept hierarchy due to their already high probabilities. Other methods, such as fine-tuning the  $\delta$  values, may help address this issue. However, we find the rebalancing method works extremely well. Experiments show a remarkable improvement in ADI performance with rebalancing, as shown in Figure 8. Without rebalancing, the ADI gives unstable results. Due to the biased adjustment towards certain branches, the results are often stuck at suboptimal levels.

**4.3.3. Cost of Attack.** In this section, we examine the number of accesses required by the ADI and GDI methods to attack the target model. Figure 9 shows that ADI needs much fewer accesses to generate high-quality auxiliary data



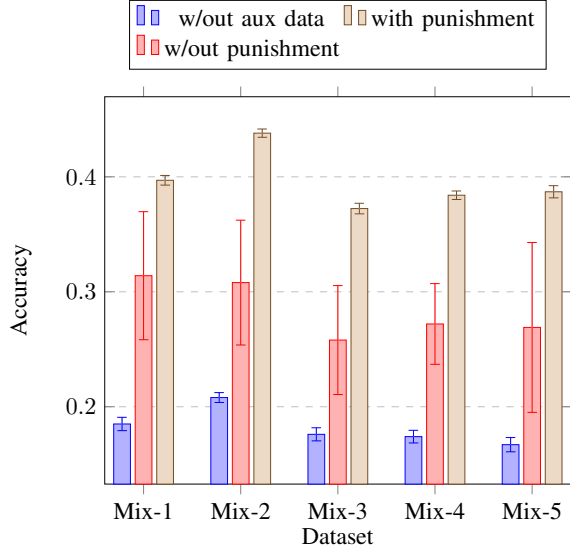


Figure 8: ADI-extracted datasets, with and without the rebalancing, impact model-inversion attacks.

to enhance the model inversion attack. We attribute this advantage of ADI to its class-based exploration approach, which allows it to quickly identify similar classes and ignore irrelevant ones. In contrast, GDI must test the whole dataset regardless of how irrelevant some of the classes in the dataset are.

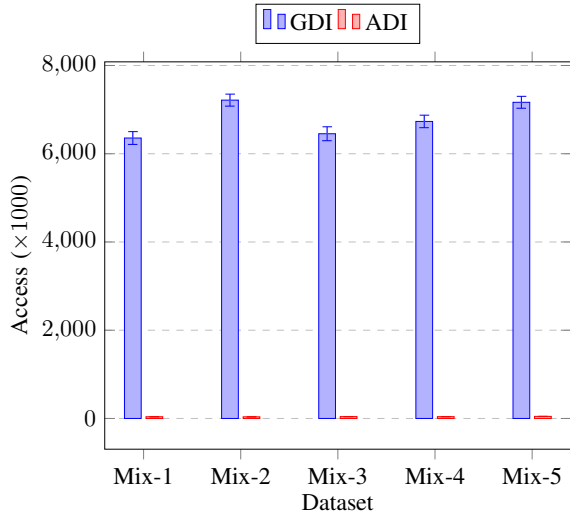


Figure 9: The required number of target-model accesses shows that ADI requires much fewer accesses.

#### 4.4. Discussion on Possible Mitigation Methods

A well-known approach to protecting models in model inference is secure inference with encrypted models, e.g., CryptoNet [24], MUSE [25], and SIMC [26]. However, these methods are still too expensive to deploy. Furthermore, since large models heavily depend on GPUs, secure infer-

ence with GPUs is still a significant challenge [27]–[29]. Trusted execution environments (TEEs) [30]–[32] can be a promising direction. However, side-channel attacks [33], [34] can be a severe concern, and TEEs for GPUs [35]–[37] are only available most recently (e.g., Nvidia H100 [37], whose protection mechanism is to be validated). Another popular approach is implementing differential privacy during the model training phase, for example, using differentially private optimizers as suggested by [38], [39]. These methods provide rigorous mathematical privacy protection. However, it’s important to note that models trained with differential privacy often experience a significant drop in model quality, as highlighted in studies by Zhao et al. [40] and Uniyal et al. [41]. Additionally, the practical implications of the privacy budget parameter  $\epsilon$  are not yet fully understood. Choosing a larger or smaller  $\epsilon$  value can lead to varying degrees of information protection. For instance, Apple has implemented a differential privacy system with  $\epsilon$  values ranging from 2 to 16 [42], while Google’s Community Mobility Reports adopted an  $\epsilon$  value of 2.64 [43].

An interesting approach is private instance encoding. Gu et al. [11] demonstrated that dataset encoding methods, such as RMT [44] and InstaHide [45], can effectively protect models from the GDI attack. As InstaHide is too weak to protect data [46], we also experimented with NeuraCrypt [47], a supposedly more robust protection method that uses a randomized auto-encoder to encode images. RMT and NeuraCrypt transform the image training data before training the target model; thus, the target model only works on the transformed data. RMT divides a training image into blocks and applies block-wise random projection. At the same time, NeuraCrypt constructs a disguised image by encoding the images using an auto-encoder. Both methods are provably secure if the attacker knows only the disguised images and models.

We are curious about if these two private instance encoding methods can protect models from ADI attacks. We have conducted experiments with simulated target datasets, and trained the target models on these encoded datasets while leaving the dataset pool unencoded as the encoding scheme keeps secret to attackers. The rationale is that the ADI attack will be ineffective as the dataset pool lacks any secretly transformed datasets guided by the target model. As expected, in Figure 10, both RMT and NeuraCrypt can significantly reduce the ADI effectiveness. However, these disguising methods also decreased the model quality by up to 7%, as shown in Figure 11.

## 5. Related Work

Machine learning models are increasingly deployed in sensitive areas such as intrusion detection and healthcare [1], [2]. These models, often packaged as authorized API services, are frequently targeted by attackers [48]–[52]. Among these attacks, model inversion and inference attacks are significant recent developments that breach the privacy of training data.

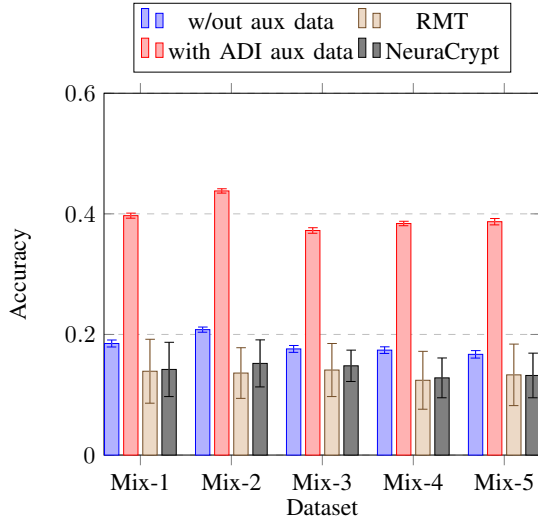


Figure 10: Accuracy deduction of MI shows that RMT and NeuraCrypt can effectively mitigate the ADI attack.

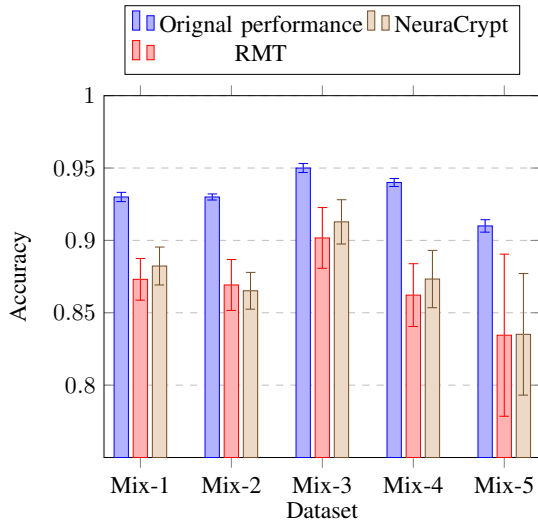


Figure 11: Accuracy deduction of target model shows that NeuraCrypt and RMT reduce the model’s quality.

Model inversion attacks attempt to reconstruct training examples from a target model with the help of auxiliary data that captures the training data distributions [4]. The attack has been constantly improved over the past few years [5], [53]. The recent GAN-based inversion attacks [5] can produce surprisingly high-quality reconstructed training data.

Membership inference attacks aim to determine if a specific record was used in training the target model [54]–[56]. Property inference attack shows that leakage of population-level properties, such as the overall distribution of certain features or the average values of attributes, is possible even when a sensitive attribute is irrelevant to the task [6], [7]. Both types of attacks involve creating shadow models and attack datasets and heavily depend on the prior knowledge

of the target model’s training data distribution, often termed “auxiliary data”.

While the assumption of an attacker knowing the auxiliary data or the domain distribution will not hold if all model information is removed from a model API service, it’s interesting to consider if a trained model inherently has contained the domain information. A recent method called GAN-based model-domain inference (GDI) attack aims to address this problem [11]. However, it does not work satisfactorily when only a few parts of candidate datasets are relevant to the target dataset. This attack also requires an excessive number of model accesses. Our proposed ADI attack directly addresses these two problems with the GDI attack.

As hiding domain information from model API accesses is insufficient to protect models from the proposed ADI attack, more sophisticated and thus possibly expensive model protection methods have to be applied. We have included a discussion about possible model defense mechanisms in Section 4.4. However, there is no clear consensus which model defense mechanism is the best.

## 6. Conclusion

Existing model-targeted attacks all assume the model domain information is available – without this information, these attacks’ performance is significantly reduced, and some may not even work anymore. One may wonder removing domain information from model APIs can be an effective countermeasure to these attacks. However, domain inference attacks show that we can still (partially) re-establish the domain information by accessing model APIs only. We have developed the Adaptive Domain Inference (ADI) attack that effectively identifies relevant sample classes within a large dataset pool by accessing the target model’s API only. Utilizing the concept hierarchy built from the dataset pool, ADI employs an iterative method to adjust the probability of concepts being part of the target model’s domain. ADI not only enhances the quality of the datasets extracted for domain inference but also reduces the number of required accesses to the model, compared to the first domain inference attack: GDI. We have conducted extensive experiments to examine the optimal parameter settings for ADI and demonstrate the unique strengths of ADI.

## References

- [1] S. Romero-Brufau, D. Whitford, M. G. Johnson, J. Hickman, B. W. Morlan, T. Therneau, J. Naessens, and J. M. Huddleston, “Using machine learning to improve the accuracy of patient deterioration predictions: Mayo clinic early warning score (mc-ews),” *Journal of the American Medical Informatics Association*, vol. 28, no. 6, pp. 1207–1215, 2021.
- [2] K. Shameer, K. W. Johnson, A. Yahi, R. Miotto, L. Li, D. Ricks, J. Jebakaran, P. Kovatch, P. P. Sengupta, S. Gelijns *et al.*, “Predictive modeling of hospital readmission rates using electronic medical record-wide machine learning: a case-study using mount sinai heart failure cohort,” in *Pacific Symposium on Biocomputing 2017*. World Scientific, 2017, pp. 276–287.

- [3] J. Alspector and T. Dietterich, "Darpa's role in machine learning," *AI Magazine*, vol. 41, no. 2, pp. 36–48, Jun. 2020. [Online]. Available: <https://ojs.aaai.org/aimagazine/index.php/aimagazine/article/view/5298>
- [4] M. Fredrikson, S. Jha, and T. Ristenpart, "Model inversion attacks that exploit confidence information and basic countermeasures," in *Proceedings of the 22nd ACM SIGSAC conference on computer and communications security*, 2015, pp. 1322–1333.
- [5] Y. Zhang, R. Jia, H. Pei, W. Wang, B. Li, and D. Song, "The secret revealer: Generative model-inversion attacks against deep neural networks," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 253–261.
- [6] K. Ganju, Q. Wang, W. Yang, C. A. Gunter, and N. Borisov, "Property inference attacks on fully connected neural networks using permutation invariant representations," in *Proceedings of the 2018 ACM SIGSAC conference on computer and communications security*, 2018, pp. 619–633.
- [7] W. Zhang, S. Tople, and O. Ohrimenko, "Leakage of dataset properties in {Multi-Party} machine learning," in *30th USENIX security symposium (USENIX Security 21)*, 2021, pp. 2687–2704.
- [8] R. Shokri, M. Stronati, C. Song, and V. Shmatikov, "Membership inference attacks against machine learning models," in *2017 IEEE symposium on security and privacy (SP)*. IEEE, 2017, pp. 3–18.
- [9] H. Hu, Z. Salicic, L. Sun, G. Dobbie, P. S. Yu, and X. Zhang, "Membership inference attacks on machine learning: A survey," *ACM Computing Surveys (CSUR)*, vol. 54, no. 11s, pp. 1–37, 2022.
- [10] Y. Zhang, R. Jia, H. Pei, W. Wang, B. Li, and D. Song, "The secret revealer: Generative model-inversion attacks against deep neural networks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 253–261.
- [11] Y. Gu and K. Chen, "Gan-based domain inference attack," *arXiv preprint arXiv:2212.11810*, 2022.
- [12] D. Alvarez-Melis and N. Fusi, "Geometric dataset distances via optimal transport," 2020.
- [13] C. A. Choquette-Choo, F. Tramèr, N. Carlini, and N. Papernot, "Label-only membership inference attacks," in *International conference on machine learning*. PMLR, 2021, pp. 1964–1974.
- [14] T. Matsumoto, T. Miura, and N. Yanai, "Membership inference attacks against diffusion models," *arXiv preprint arXiv:2302.03262*, 2023.
- [15] J. Hayes, L. Melis, G. Danezis, and E. De Cristofaro, "Logan: Membership inference attacks against generative models," *arXiv preprint arXiv:1705.07663*, 2017.
- [16] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.
- [17] L. Deng, "The mnist database of handwritten digit images for machine learning research," *IEEE Signal Processing Magazine*, vol. 29, no. 6, pp. 141–142, 2012.
- [18] G. Cohen, S. Afshar, J. Tapson, and A. van Schaik, "Emnist: an extension of mnist to handwritten letters," *arXiv preprint arXiv:1702.05373*, 2017.
- [19] G. B. Huang, M. Mattar, T. Berg, and E. Learned-Miller, "Labeled faces in the wild: A database for studying face recognition in unconstrained environments," in *Workshop on faces in 'Real-Life' Images: detection, alignment, and recognition*, 2008.
- [20] A. Krizhevsky, G. Hinton *et al.*, "Learning multiple layers of features from tiny images," 2009.
- [21] H. Xiao, K. Rasul, and R. Vollgraf, "Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms," 2017. [Online]. Available: <https://arxiv.org/abs/1708.07747>
- [22] A. Grigorev, "Clothing dataset (full, high resolution)," <https://www.kaggle.com/datasets/agrigorev/clothing-dataset-full>, 10 2020, accessed: 2020-10-21.
- [23] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [24] R. Gilad-Bachrach, N. Dowlin, K. Laine, K. Lauter, M. Naehrig, and J. Wernsing, "Cryptonets: Applying neural networks to encrypted data with high throughput and accuracy," in *Proceedings of The 33rd International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, M. F. Balcan and K. Q. Weinberger, Eds., vol. 48, 2016, pp. 201–210.
- [25] R. Lehmkuhl, P. Mishra, A. Srinivasan, and R. A. Popa, "Muse: Secure inference resilient to malicious clients," in *30th USENIX Security Symposium (USENIX Security 21)*. USENIX Association, Aug. 2021, pp. 2201–2218. [Online]. Available: <https://www.usenix.org/conference/usenixsecurity21/presentation/lehmkuhl>
- [26] N. Chandran, D. Gupta, S. L. B. Obbattu, and A. Shah, "SIMC: ML inference secure against malicious clients at Semi-Honest cost," in *31st USENIX Security Symposium (USENIX Security 22)*. Boston, MA: USENIX Association, Aug. 2022, pp. 1361–1378. [Online]. Available: <https://www.usenix.org/conference/usenixsecurity22/presentation/chandran>
- [27] F. Tramèr and D. Boneh, "Slalom: Fast, verifiable and private execution of neural networks in trusted hardware," in *International Conference on Learning Representations*, 2019. [Online]. Available: <https://openreview.net/forum?id=rjV0rjCcKQ>
- [28] T. Hunt, Z. Jia, V. Miller, A. Szekely, Y. Hu, C. J. Rossbach, and E. Witchel, "Telekine: Secure computing with cloud gpus," in *17th USENIX Symposium on Networked Systems Design and Implementation (NSDI 20)*. Santa Clara, CA: USENIX Association, Feb. 2020, pp. 817–833. [Online]. Available: <https://www.usenix.org/conference/nsdi20/presentation/hunt>
- [29] L. K. L. Ng, S. S. M. Chow, A. P. Y. Woo, D. P. H. Wong, and Y. Zhao, "Goten: Gpu-outsourcing trusted execution of neural network training," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 17, pp. 14 876–14 883, May 2021. [Online]. Available: <https://ojs.aaai.org/index.php/AAAI/article/view/17746>
- [30] V. Costan and S. Devadas, "Intel sgx explained," *IACR Cryptology ePrint Archive*, vol. 2016, p. 86, 2016.
- [31] D. Kaplan, J. Powell, and T. Woller, "Amd memory encryption," <https://www.amd.com/content/dam/amd/en/documents/epyc-business-docs/white-papers/memory-encryption-white-paper.pdf>, 2021.
- [32] S. Thornton, "Arm trustzone explained," <http://alturl.com/icptx>, 2017.
- [33] S. Fei, Z. Yan, W. Ding, and H. Xie, "Security vulnerabilities of sgx and countermeasures: A survey," *ACM Comput. Surv.*, vol. 54, no. 6, 2021.
- [34] X. Lou, T. Zhang, J. Jiang, and Y. Zhang, "A survey of microarchitectural side-channel vulnerabilities, attacks, and defenses in cryptography," *ACM Comput. Surv.*, vol. 54, no. 6, jul 2021.
- [35] S. Volos, K. Vaswani, and R. Bruno, "Graviton: Trusted execution environments on gpus," in *13th USENIX Symposium on Operating Systems Design and Implementation (OSDI 18)*. Carlsbad, CA: USENIX Association, Oct. 2018, pp. 681–696. [Online]. Available: <https://www.usenix.org/conference/osdi18/presentation/volos>
- [36] I. Jang, A. Tang, T. Kim, S. Sethumadhavan, and J. Huh, "Heterogeneous isolated execution for commodity gpus," in *Proceedings of the Twenty-Fourth International Conference on Architectural Support for Programming Languages and Operating Systems*, ser. ASPLOS '19. New York, NY, USA: Association for Computing Machinery, 2019, pp. 455–468. [Online]. Available: <https://doi.org/10.1145/3297858.3304021>
- [37] Nvidia, "Confidential compute on nvidia hopper h100," <https://images.nvidia.com/aem-dam/en-zz/Solutions/data-center/HCC-Whitepaper-v1.0.pdf>, 2023.
- [38] M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, and L. Zhang, "Deep learning with differential privacy," in *Proceedings of the 2016 ACM SIGSAC conference on computer and communications security*, 2016, pp. 308–318.

- [39] C. Dupuy, R. Arava, R. Gupta, and A. Rumshisky, "An efficient dp-sgd mechanism for large scale nlu models," in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022, pp. 4118–4122.
- [40] J. Zhao, Y. Chen, and W. Zhang, "Differential privacy preservation in deep learning: Challenges, opportunities and solutions," *IEEE Access*, vol. 7, pp. 48 901–48 911, 2019.
- [41] A. Uniyal, R. Naidu, S. Kotti, S. Singh, P. J. Kenfack, F. Mireshghallah, and A. Trask, "Dp-sgd vs pate: Which has less disparate impact on model accuracy?" *arXiv preprint arXiv:2106.12576*, 2021.
- [42] A. Differential Privacy Team, "Learning with privacy at scale differential," 2017. [Online]. Available: <https://api.semanticscholar.org/CorpusID:43986173>
- [43] A. Aktay, S. Bavadekar, G. Cossoul, J. Davis, D. Desfontaines, A. Fabrikant, E. Gabrilovich, K. Gadepalli, B. Gipson, M. Guevara, C. Kamath, M. Kansal, A. Lange, C. Mandayam, A. Oplinger, C. Pluntke, T. Roessler, A. Schlosberg, T. Shekel, S. Vispute, M. Vu, G. Wellenius, B. Williams, and R. J. Wilson, "Google covid-19 community mobility reports: Anonymization process description (version 1.1)," 2020.
- [44] S. Sharma, A. M. Alam, and K. Chen, "Image disguising for protecting data and model confidentiality in outsourced deep learning," in *IEEE Conference on Cloud Computing*, 2021.
- [45] Y. Huang, Z. Song, K. Li, and S. Arora, "Instahide: Instance-hiding schemes for private distributed learning," *CoRR*, vol. abs/2010.02772, 2020. [Online]. Available: <https://arxiv.org/abs/2010.02772>
- [46] N. Carlini, S. Deng, S. Garg, S. Jha, S. Mahloujifar, M. Mahmoody, A. Thakurta, and F. Tramèr, "Is private learning possible with instance encoding?" in *2021 IEEE Symposium on Security and Privacy (SP)*. IEEE, 2021, pp. 410–427.
- [47] A. Yala, H. Esfahanizadeh, R. G. L. D. Oliveira, K. R. Duffy, M. Ghobadi, T. S. Jaakkola, V. Vaikuntanathan, R. Barzilay, and M. Medard, "Neuracrypt: Hiding private health data via random neural networks for public training," 2021.
- [48] Z. Alkhalil, C. Hewage, L. Nawaf, and I. Khan, "Phishing attacks: A recent comprehensive study and a new anatomy," *Frontiers in Computer Science*, vol. 3, p. 563060, 2021.
- [49] B. B. Gupta, A. Tewari, A. K. Jain, and D. P. Agrawal, "Fighting against phishing attacks: state of the art and future challenges," *Neural Computing and Applications*, vol. 28, pp. 3629–3654, 2017.
- [50] B. Pingle, A. Mairaj, and A. Y. Javaid, "Real-world man-in-the-middle (mitm) attack implementation using open source tools for instructional use," in *2018 IEEE International Conference on Electro/Information Technology (EIT)*. IEEE, 2018, pp. 0192–0197.
- [51] M. A. Al-Shareeda, M. Anbar, S. Manickam, and I. H. Hasbullah, "Review of prevention schemes for man-in-the-middle (mitm) attack in vehicular ad hoc networks," *International Journal of Engineering and Management Research*, vol. 10, 2020.
- [52] B. Bhushan, G. Sahoo, and A. K. Rai, "Man-in-the-middle attack in wireless and computer networking—a review," in *2017 3rd International Conference on Advances in Computing, Communication & Automation (ICACCA)(Fall)*. IEEE, 2017, pp. 1–6.
- [53] S. Hidano, T. Murakami, S. Katsumata, S. Kiyomoto, and G. Hanaoka, "Model inversion attacks for prediction systems: Without knowledge of non-sensitive attributes," in *2017 15th Annual Conference on Privacy, Security and Trust (PST)*. IEEE, 2017, pp. 115–11509.
- [54] R. Shokri, M. Stronati, C. Song, and V. Shmatikov, "Membership inference attacks against machine learning models," in *2017 IEEE symposium on security and privacy (SP)*. IEEE, 2017, pp. 3–18.
- [55] B. Hui, Y. Yang, H. Yuan, P. Burlina, N. Z. Gong, and Y. Cao, "Practical blind membership inference attack via differential comparisons," *arXiv preprint arXiv:2101.01341*, 2021.
- [56] C. A. Choquette-Choo, F. Tramèr, N. Carlini, and N. Papernot, "Label-only membership inference attacks," in *International conference on machine learning*. PMLR, 2021, pp. 1964–1974.