# Do Not Trust Prediction Scores for Membership Inference Attacks

Dominik Hintersdorf[*,1], Lukas Struppek[*,1], and Kristian Kersting[1,2]

[1]Department of Computer Science, TU Darmstadt, Darmstadt, Germany
[2]Centre for Cognitive Science, TU Darmstadt, and Hessian Center for AI (hessian.AI)
*firstname.lastname@tu-darmstadt.de*

## Abstract

Membership inference attacks (MIAs) aim to determine whether a specific sample was used to train a predictive model. Knowing this may indeed lead to a privacy breach. Arguably, most MIAs, however, make use of the model's prediction scores—the probability of each output given some input—following the intuition that the trained model tends to behave differently on its training data. We argue that this is a fallacy for many modern deep network architectures, e.g., ReLU type neural networks produce almost always high prediction scores far away from the training data. Consequently, MIAs will miserably fail since this behavior leads to high false-positive rates not only on known domains but also on out-of-distribution data and implicitly acts as a defense against MIAs. Specifically, using generative adversarial networks, we are able to produce a potentially infinite number of samples falsely classified as part of the training data. In other words, the threat of MIAs is overestimated and less information is leaked than previously assumed. Moreover, there is actually a trade-off between the overconfidence of classifiers and their susceptibility to MIAs: the more classifiers know when they do not know, making low confidence predictions far away from the training data, the more they reveal the training data.

## 1 Introduction

Deep learning models achieve state-of-the-art performances in various tasks such as computer vision, language modeling, healthcare, and autonomous driving, among others. However, large amounts of data are needed to train these models appropriately (from scratch) and to generalize well on complex tasks. Collecting and, in particular, cleaning and labeling data is expensive. Hence, users may look for alternative data sources, which may not always be
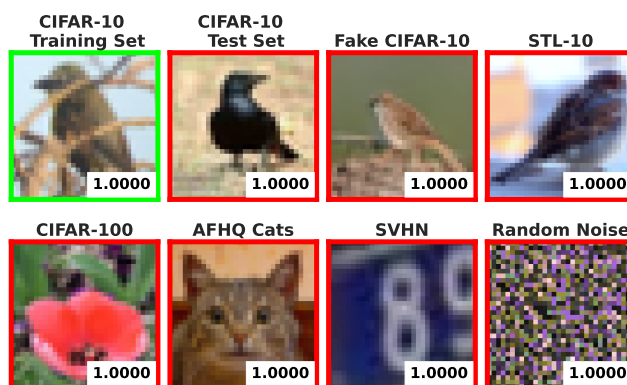
---
*equal contribution



Figure 1: False-positive membership inference attacks (red frames) against a ResNet-18 trained on CIFAR-10 and their assigned maximum prediction scores.

legal ones. To detect data abuse, it would be desirable to prove whether a model was trained on leaked or unauthorized retrieved data. However, to prove that a specific data point was part of the training set is difficult since neural networks do not store plain training data like lazy learners. Instead, the learned knowledge is encoded into the network's weights.

One way to distinguish between unseen data and data points used for training the neural networks is through membership inference attacks (MIAs). They attempt to identify training samples in a large set of possible inputs. Besides malicious intentions, MIAs might be used to prove illegal data abuse in deep learning settings. To use membership inference results of such attacks as evidence in court, high accuracy, as well as robustness to different data types and network architectures, is required. Previous work on MIAs, see e.g. [26, 25], state high accuracy in distinguishing between training and test data. However, the evaluation of MIAs reported in the literature is usually done in a cross-validation setting, i.e., on samples from the same data distribution.
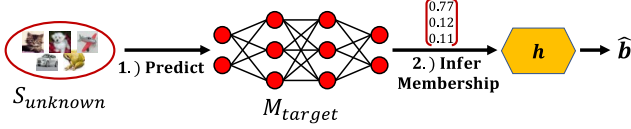
Figure 2: Application of inference model $h$ against a target model $M_{target}$. The adversary first queries $M_{target}$ to collect prediction scores, and the inference model $h$ is then used to make a prediction $\hat{b}$ on the membership status.

Specifically, we argue that MIAs, in particular attacks based on a model's prediction scores, are not robust and not very meaningful in realistic settings due to their high false-positive rates, also criticized by [24]. We take, however, a broader view and do not restrict evaluation on the target model's training data distribution. In a specific domain, there is a possibly infinite number of samples, and hence the number of false-positives can be increased arbitrarily. This leads to reduced informative value and low reliability of the attacks in realistic scenarios. Fig. 1 shows samples from various datasets for which all three MIAs studied in this paper make false-positive predictions, even if the inputs are nothing similar to the training data or do not contain any meaningful information at all. We practically demonstrate the theoretically unlimited number of false-positive member classifications by using a GAN-based approach to generate images following the training data distribution.

Our argumentation relies on the often reported overconfidence of modern deep neural architectures, see e.g. [20, 8, 5, 15], and our experimental results indicate that mitigating the overconfidence of neural networks using calibration techniques increases privacy leakage. We explore the trade-off between well-calibrated models and the predictive power of MIAs. Moreover, previous works considered generalization only in terms of the difference between training and test accuracy. In contrast, we interpret generalization as the ability of a model to know when it does not know something and consequently to express its uncertainty through its prediction scores. We argue that attack evaluations in previous works distorted the actual effectiveness of the attacks by using only data from the target model's training data distribution. We strongly believe that in reality, the attacker has only limited information about the real data distribution to effectively perform MIAs. There might not exist any meaningful MIA at all since the number of false-positive samples cannot be limited because an attacker will never be able to precisely approximate real data distributions.

To summarize, we make the following contributions:

1. We demonstrate that the practical relevance arising from MIAs has been overestimated due to their high false-positive rates on data from the training distribution as well as on out-of-distribution samples.
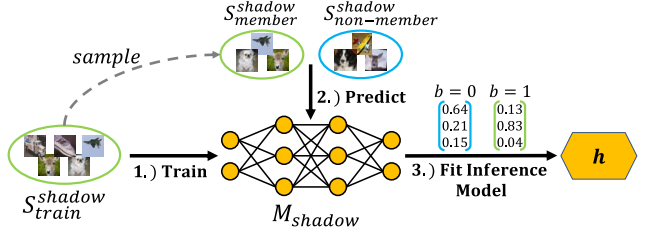


Figure 3: Membership inference preparation process. The adversary first trains a shadow model $M_{shadow}$ in order to then optimize the parameters of inference model $h$.

2. We empirically identify the overconfidence of neural networks as the main reason for the high false-positive rate of these attacks.

3. We draw attention to the fact that there is a trade-off between keeping a model secure against MIAs and mitigating overconfidence.

We proceed as follows. We start off by reviewing MIAs and (mitigating) overconfidence of neural networks. Afterward, we introduce the theoretical background, our experimental setup and present our experimental evidence. Before concluding, we discuss our results.

## 2 Membership Inference Attacks

Membership inference attacks (MIAs) on neural networks were first introduced by Shokri et al. [26]. In a general MIA setting, as usually assumed in the literature, an adversary is given an input $x$ following distribution $D$ and a target model $M_{target}$ which was trained on a training set $S_{train}^{target} \sim D^n$ with size $n$. The adversary is then facing the problem to identify whether a given $x \sim D$ was part of the training set $S_{train}^{target}$. To predict the membership of $x$, the adversary creates an inference model $h$. In score-based MIAs, the input to $h$ is the prediction score vector produced by $M_{target}$ on sample $x$. Fig. 2 visualizes such an attack exploiting prediction score vectors. Since MIAs are binary classification problems, precision, recall and false-positive rate (FPR) are used as attack evaluation metrics.

All MIAs exploit a difference in the behavior of $M_{target}$ on seen and unseen data. Most attacks in the literature follow Shokri et al. [26] and train so-called shadow models $M_{shadow}$ on a disjoint dataset $S_{train}^{shadow}$ drawn from the same distribution $D$ as $S_{train}^{target}$. $M_{shadow}$ is used to mimic the behavior of $M_{target}$ and adjust parameters of $h$, such as threshold values or model weights. Note that the membership status for inputs to $M_{shadow}$ are known to the adversary. Fig. 3 visualizes the attack preparation process.

In recent years, various MIAs have been proposed. Shokri et al. [26] trained multiple shadow models and

queried each of the shadow models with its training data (members), as well as unseen data (non-members) to retrieve the prediction scores of the shadow models. Multiple binary classifiers were then trained for each class label to predict the membership status.

Salem et al. [25] also exploited prediction scores and trained a single class-agnostic neural network to distinguish between members and non-members. In contrast to Shokri et al. [26], their approach relies on a single shadow model. The input of $h$ consists of the $k$ highest prediction scores in descending order.

Instead of focusing solely on the scores, Yeom et al. [33] took advantage of the fact that the loss of a model is lower on members than on non-members and fit a threshold to the loss values. More recent approaches [3, 16] focused on label-only attacks where only the predicted label for a known input is observed.

Most defense strategies either try to decrease the informative value of prediction scores or reduce overfitting. The informative value can be decreased by adding a large temperature to the softmax function to increase its entropy [26], adding carefully crafted noise to the predictions [9] or outputting only the predicted label without any score [26]. Various regularization techniques were proposed to reduce overfitting and thus the accuracy gap, e.g., L2 regularization [26] and dropout [26, 25].

## 3 Overconfidence of Neural Networks

Neural networks usually output prediction scores, e.g., by applying a softmax function. To take model uncertainty into account, it is generally desired that the prediction scores represent the probability of a correct prediction, which is usually not the case. This problem is generally referred to as model calibration. Guo et al. [5] demonstrated that modern networks tend to be overconfident in their predictions.

Generally, as Hein et al. [7] noted, there have been many cases reported where high prediction scores are made far away from the training data by neural networks, e.g., on fooling images, for out-of-distribution (OOD) images, in a medical diagnosis task, but also on the original task. Hein *et al.* then proved that ReLU networks are overconfident even on samples far away from the training data.

Scaling the inputs to ReLU network actually allows one to produce arbitrarily high prediction scores. Existing approaches to mitigate overconfidence can be grouped into two categories: post-processing methods applied on top of trained models and regularization methods modifying the training process.

As a post-processing method, Guo et al. [5] proposed temperature scaling using a single temperature parameter $T$ for scaling down the pre-softmax logits for all classes. The

larger $T$ is, the more the resulting scores approach a uniform distribution while its entropy increases.

Kristiadi et al. [13] proposed a Bayesian approach. They fixed the weights for all layers of a trained network except the last one and used a Kronecker-factored Laplace approximation (LA) on the weights of the final layer. Müller et al. [17] demonstrated that label smoothing regularization [28] not only improves the generalization of a model but also implicitly leads to better model calibration. It reduces the difference between the highest and the other logit values, thus reducing overconfident predictions. The calibration of a model can be measured by the expected calibration error (ECE) [18]. It computes a weighted average over the absolute difference between test accuracy and prediction scores.

## 4 Do Not Trust Prediction Scores for MIAs

In this section, we will show that predictions scores for MIAs cannot be trusted because score-based MIAs make membership decisions based mainly on the maximum prediction score. As a first step, we introduce our proposition and then verify our claims empirically.

Formally, a neural network $f(x)$ using ReLU activations decomposes the unrestricted input space $\mathbb{R}^m$ into a finite set of polytopes (linear regions). We can then interpret $f(x)$ as a piecewise affine function that is affine in any polytope. Due to the limited number of polytopes, the outer polytopes extend to infinity which allows to arbitrarily increase the prediction scores through scaling inputs by a large constant $\delta$ [7]. Applying these findings to MIAs results in the following proposition:

**Proposition 1.** *Given a ReLU-classifier, we can force almost any non-member input to be classified as a member by score-based MIAs simply through scaling it by a large constant.*

*Proof.* Let $f : \mathbb{R}^m \to \mathbb{R}^d$ be a piecewise affine ReLU-classifier. We define a score-based MIA inference model $h : \mathbb{R}^d \to \{0, 1\}$ with 1 indicating a classification as member. For almost any input $x \in \mathbb{R}^m$ and a sufficiently small $\epsilon > 0$ if $\max_{i=1,...,d} f(x)_i \geq 1 - \epsilon$ it follows that $h(f(x)) = 1$. Since $lim_{\delta \to \infty} \max_{i=1,...,d} f(\delta x)_i = 1$, then $lim_{\delta \to \infty} h(f(\delta x)) = 1$ already holds. $\square$

By scaling the whole non-member dataset, one can force the FPR to be close to 100%. Indeed, the proposition holds only for ReLU-networks and unbounded inputs, which are not restricted to the range of $[0, 1]^m$. Next, we empirically show that one cannot trust predictions scores for MIAs in more general settings without input scaling required and using other activation functions.

## 4.1 Experimental Protocol

We provide our source code on github [1] and state further information for reproducibility in Appx. A.

**Threat Model.** Most demonstrations of MIAs [26, 33, 25, 27] make assumptions on the adversary such as using shadow models, knowledge of the target model structure, and having a set of samples from the same distribution as the target model's training data. Here, we followed the setting of [26] for our inference attacks but only trained a single shadow model for each attack following [25]. We even aimed at simulating a worst-case scenario, i.e., the adversary further knows the exact architecture and training procedure of the target model. Therefore, a strong shadow model could be trained, following the procedure depicted in Fig. 3. In this scenario, the adversary only has access to the target model's prediction score vector for membership inference.

**Datasets**; more details on the datasets and preprocessing steps in Appx. A.6. We evaluated the attacks on models trained on the CIFAR-10 [14] and Stanford Dogs [11] datasets for image classification. To improve generalization and classification accuracy on the Stanford Dogs dataset, we performed usual data augmentation techniques during training.

We created two disjoint training datasets for the target and shadow models, each containing 12,500 (CIFAR-10) and 8,232 (Stanford Dogs) samples. We then randomly drew 2,500 and 2,058 samples, respectively, from the training and test sets to create the member and non-member datasets.

We used various datasets to demonstrate the susceptibility of prediction score-based MIAs to high scores on samples from neighboring distributions and samples far away from the training data—a kind of OOD setting. For CIFAR-10 models, we used STL-10 [4], CIFAR-100 [14] and Animal Faces-HQ (AFHQ) Cats dataset [2] as neighboring distributions containing samples of similar classes or visual styles. For distributions further away, we used samples from the SVHN [19] dataset.

Additionally, we used a class-conditional pre-trained StyleGAN2 [10] to generate synthetic CIFAR-10 samples, referred to as Fake CIFAR-10. To push our approach to the extreme, we created two additional datasets based on CIFAR-10 test images by randomly permuting the images' pixels to create random noise samples and scaling pixel values by factor 255 to imitate possible mistakes during preprocessing. In the following, we refer to these two datasets as Permuted and Scaled, respectively.

For evaluation on the Stanford Dogs models, we used the dog samples from AFHQ as neighboring distribution and the cat samples as OOD data. We also created Permuted

and Scaled datasets based on Stanford Dogs test images and generated synthetic dog images by using a StyleGAN2 trained on the AFHQ dataset, referred to as Fake Dogs.

**Neural Network Architectures**; more details on the architectures and hyperparameters in Appx. A.1. We trained ResNet-18 [6], EfficientNetB0 [29] and a simple convolutional neural network following Salem et al. [25], referred to as SalemCNN, on CIFAR-10. For the Stanford dogs dataset, we used a larger ResNet-50 architecture pre-trained on ImageNet. ResNets and SalemCNN are ReLU networks and can be interpreted as piecewise linear functions [1]. EfficientNetB0 uses Swish activation functions [23], which are not piecewise linear and, therefore, our proposition does not hold. Nevertheless, we demonstrate that also non-ReLU networks suffer from overconfidence, and MIAs are also not robust.

**Prediction Score-Based Attacks**; more details on the attacks in Appx. A.2. We base our analysis on three different MIAs exploiting the maximum value, the entropy, and the top-3 values of the prediction score vector, mainly following Salem et al. [25]. For the top-3 prediction score attack, we trained a small neural network with a single hidden layer with 64 neurons as inference model. It uses the three highest scores of $M_{target}$ in descending order as inputs.

The maximum prediction score attack only relies on the highest score, while the entropy attack computes the entropy on the whole score vector. An input sample is then classified as a member if the maximum value is higher or if the entropy is lower than a threshold, respectively. We fit both thresholds on the shadow models' outputs.

## 4.2 Experimental Results

Our intention is to investigate the following questions using the protocol: **(Q1)** How robust are prediction score-based MIAs? **(Q2)** Does overconfidence interact with MIAs? **(Q3)** How does calibrating neural networks influence the success of MIAs? **(Q4)** Are defenses contrary to calibration?

**(Q1) MIAs Are Not Robust.** Tab. 1 summarizes the prediction accuracy of the CIFAR-10 target models, as well as the attack metrics, and Tab. 2 for the Stanford Dogs models. As one can see, the attacks performed quite similarly while the recall is always significantly higher than the precision, indicating the problem of many false-positive predictions.

To examine the robustness of the attacks, we then used the remaining datasets as non-member inputs and measured the FPRs. Fig. 4 shows the FPR of the attacks against the CIFAR-10 models. The results demonstrate that the attacks not only tend to falsely classify samples from the test data as members but also samples from other distributions. The attacks misclassified STL-10 samples as members in more

---

[1]https://github.com/ml-research/Do-Not-Trust-Prediction-Scores-for-Membership-Inference-Attacks

|  | SalemCNN | ResNet-18 | EfficientNetB0 |
|---|---|---|---|
| Train Accuracy | 100.00% | 100.00% | 99.03% |
| Test Accuracy | 59.04% | 69.38% | 71.06% |
| Entropy Pre | 65.51% | 67.35% | 61.36% |
| Entropy Rec | 88.52% | 92.32% | 79.96% |
| Entropy FPR | 46.60% | 44.76% | 50.36% |
| Max. Score Pre | 65.34% | 67.35% | 61.43% |
| Max. Score Rec | 87.48% | 92.32% | 79.64% |
| Max. Score FPR | 46.40% | 44.76% | 50.00% |
| Top-3 Scores Pre | 62.48% | 63.84% | 60.74% |
| Top-3 Scores Rec | 100.00% | 98.04% | 82.60% |
| Top-3 Scores FPR | 60.04% | 55.52% | 53.40% |

Table 1: Training and attack metrics for the target models trained on CIFAR-10. We measure the attacks' precision and recall on equally-sized member and non-member subsets from CIFAR-10.

| ResNet-50 | Standard | LS | LA | Temp | L2 |
|---|---|---|---|---|---|
| Train Accuracy | 98.48% | 99.62% | 98.48% | 98.48% | 74.05% |
| Test Accuracy | 59.69% | 64.65% | 59.62% | 59.69% | 48.15% |
| ECE | 25.09% | 5.80% | 9.92% | 51.03% | 11.86% |
| Entropy Pre | 68.22% | 76.33% | 66.12% | 59.45% | 60.50% |
| Entropy Rec | 84.50% | 82.56% | 87.61% | 47.38% | 50.68% |
| Entropy FPR | 39.36% | 25.61% | 44.90% | 32.31% | 33.09% |
| Max. Score Pre | 68.30% | 77.32% | 67.42% | 63.96% | 59.13% |
| Max. Score Rec | 83.97% | 81.83% | 86.49% | 65.55% | 56.66% |
| Max. Score FPR | 38.97% | 24.00% | 41.79% | 36.93% | 39.16% |
| Top-3 Scores Pre | 68.52% | 76.57% | 66.64% | 67.65% | 58.61% |
| Top-3 Scores Rec | 83.87% | 85.42% | 88.44% | 85.96% | 59.23% |
| Top-3 Scores FPR | 38.53% | 26.14% | 44.27% | 41.11% | 41.84% |

Table 2: Training and attack metrics for ResNet-50 target models trained on Stanford Dogs. We compare the results for the standard model to models trained with label smoothing (LS) and Laplace approximation (LA) as calibration techniques and temperature scaling (Temp) and L2 regularization as defense techniques.

than a third and samples from AFHQ Cats in allmost 60% of the cases. For the OOD SVHN dataset, the attacks even showed an FPR of up to 28%.

Moreover, the Permuted and Scaled CIFAR-10 test samples resulted in the highest FPR, empirically confirming our proposition and demonstrating that neural networks are not able to recognize when they are operating on unknown inputs and still produce high prediction scores. This leads to less robust MIAs. The fact that the FPR on the Fake CIFAR-10 samples is even higher than on the real CIFAR-10 samples clearly shows that there exists a potentially infinite number of false-positive samples.

This behavior is not limited to ReLU networks. The FPR of the EfficientNetB0 on the datasets is quite similar to the FPR of the ResNet-18. This indicates that the problem of high FPR is affecting modern deep architectures in general. We repeated the experiments on the more complex Stanford Dogs dataset using a ResNet-50. Training and attack results are summarized in Tab. 2. The attacks achieved a precision of about 68% and a recall >80% but still lacked the ability to recognize unknown domains.

The transparent bars in Figs. 5a and 5c show the FPR for this setting. The attacks still showed false-positive predictions on about a third of the AFHQ Dogs and up to 13% on the AFHQ Cats samples. By generating synthetic dog images, however, we could easily fool the attacks in more than a quarter of the cases. Scaling Stanford Dogs samples resulted in a 98% FPR with maximum scores close to 1.0, confirming our proposition.

**(Q2) High Prediction Scores May Lower Privacy Risks.** To shed light on the connection between overconfidence and high FPR of the MIAs, we analyzed the mean maximum prediction scores (MMPS) of the target models' predictions.

Tab. 3 shows the MMPS values measured on a standard ResNet-50 among others and underlines our assumption that all score-based MIAs against models trained with standard procedure mainly rely on the maximum score since there is a clear difference between the MMPS of false-positive and true-negative predictions. We have obtained similar results for the CIFAR-10 models and present these in Appx. B.5. For the MMPS values of all attacks on the ResNet-50 using all datasets see Appx. B.1.

It seems that the non-maximum scores are not providing significant information on the membership status since the MMPS values of the false-positive predicted samples using the maximum score attack and the top-3 attack differ only slightly. Modifying the top-3 attack to use a larger part of the score vector for inferring membership of the samples did not significantly improve the membership inference either.

So on one side, neural networks are overconfident in their predictions, even on inputs without any meaningful information. It prevents a reasonable interpretation regarding a model's probability of being correct in its predictions. During MIAs, on the other side, this behavior implicitly protects the training data since the information content of the prediction score is rather low. Consequently, there is a trade-off between a model's ability to react to unknown inputs and its privacy leakage. We explore this trade-off in Q3.

We further argue that any adversarial example maximizing the target model's scores in an arbitrary class would also be classified as a member in almost all cases. So it is possible to hide members in a larger dataset of non-members that are altered by adversarial attacks to maximize the target model's scores in the true class label.

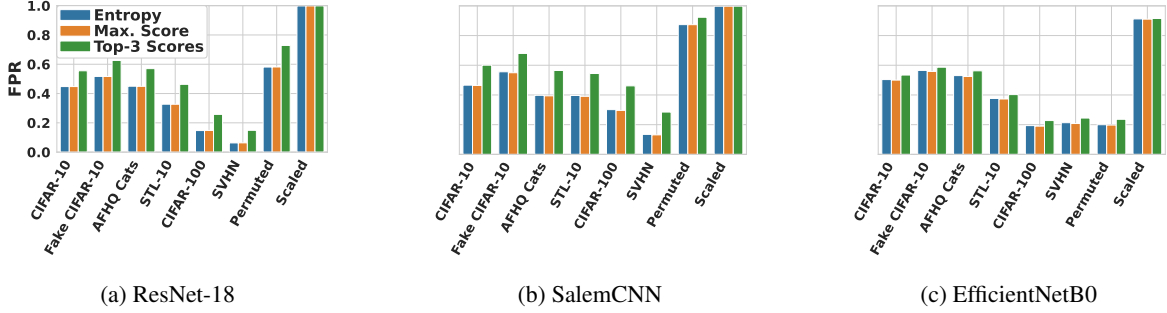(a) ResNet-18       (b) SalemCNN       (c) EfficientNetB0

Figure 4: False-positive rates (FPR) for prediction score-based membership inference attacks against CIFAR-10 models. All attacks yield high FPR, even on generated CIFAR-10 images and samples from unknown distributions.

| Dataset | ResNet-50 Model | FP MMPS | TN MMPS |
|---|---|---|---|
| Stanford Dogs | Standard | 0.9986 | 0.7596 |
| | Label Smoothing | 0.9098 | 0.4819 |
| | L2 Regularization | 0.8492 | 0.4167 |
| | Temperature | 0.1367 | 0.0539 |
| Fake Dogs | Standard | 0.9980 | 0.7737 |
| | Label Smoothing | 0.8890 | 0.4508 |
| | L2 Regularization | 0.8376 | 0.4548 |
| | Temperature | 0.1481 | 0.0749 |
| AFHQ Cats | Standard | 0.9969 | 0.7193 |
| | Label Smoothing | 0.8930 | 0.3519 |
| | L2 Regularization | 0.8223 | 0.3736 |
| | Temperature | 0.0840 | 0.0444 |

Table 3: MMPS for false-positive (FP) and true-negative (TN) predictions on the various ResNet-50 models and the Top-3 scores attack on selected datasets. Label smoothing and L2 regularization both increase the MMPS gap, while label smoothing increases the model's vulnerability, and L2 and temperature reduce it.



(a) ResNet-50 (LS)       (b) ResNet-18 (LS)
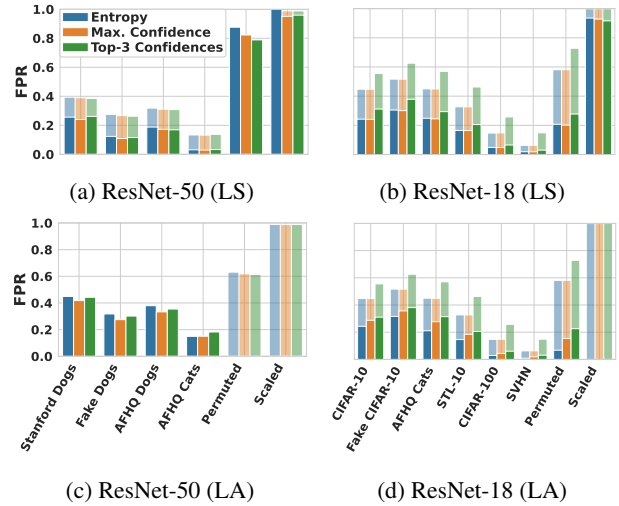
(c) ResNet-50 (LA)       (d) ResNet-18 (LA)

Figure 5: False-positive rates (FPR) of membership inference attacks against ResNet-18 and ResNet-50. The transparent bars represent the FPR of the standard models while the solid bars represent the FPR of the models with the respective modification given in parentheses - label smoothing (LS) and Laplace approximation (LA). Both calibration methods lead to a reduction in FPR for almost all inputs.

**(Q3) Calibration May Increase Privacy Risks.** Ideally, neural networks are properly calibrated and their prediction scores represent the probabilities of correct predictions. To calibrate the models and to reduce the overconfidence, we retrained the ResNet-18 and ResNet-50 with label smoothing. Smoothing factors were set to $\alpha = 0.1$ (ResNet-50) and $\alpha = \frac{0.1}{12}$ (ResNet-18). We performed the same calibration method on both the target and the shadow models, which reflects a worst-case scenario, with an adversary knowing the exact calibration method and hyperparameters.

Label smoothing not only calibrates a model but may also improve its test accuracy, as shown in Tab. 2 for ResNet-50; detailed results for ResNet-18 are given in Appx. B.3. The expected calibration error (ECE) computed on the Stanford Dogs test data with 15 bins drops from 25.09% for the standard ResNet-50 down to 5.8% with label smoothing. Similarly, for the ResNet-18 trained on CIFAR-10, the ECE drops from 24.02% to 13.33%. In both cases, it demonstrates a strong calibration effect when using label smoothing.

Previous works on MIAs suggested that minimizing the accuracy gap between the training and test accuracy on the same architecture leads to weaker attacks and, therefore, to lower privacy risks. However, as demonstrated by the results summarized in Tab. 2, label smoothing improves the test accuracy and still yields higher attack precision values for all three attacks on both architectures. Figs. 5a and 5b further illustrate that label smoothing reduces the number of false-positive membership predictions.

(a) ResNet-50    (b) ResNet-50 (LS)    (c) ResNet-50 (Temperature)    (d) ResNet-50 (L2)
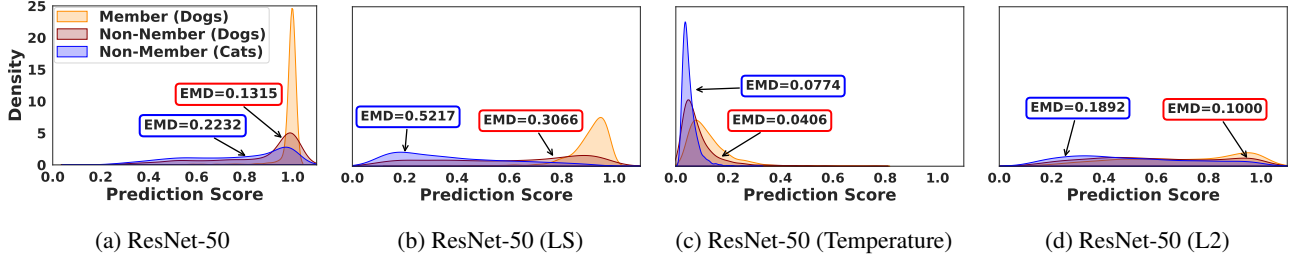
Figure 6: Kernel density estimation applying Gaussian kernels on the top prediction scores values of ResNet-50 target models. We use equally-sized member and non-member subsets of Stanford Dogs and AFHQ Cats. We further state the earth mover's distance (EMD) between each dataset and the member dataset. Label smoothing (LS) moves the non-member distributions further away, and consequently, the members become easier to separate. Temperature scaling and L2 regularization show an inverse effect and increase the overlapping, making membership inference attacks harder.



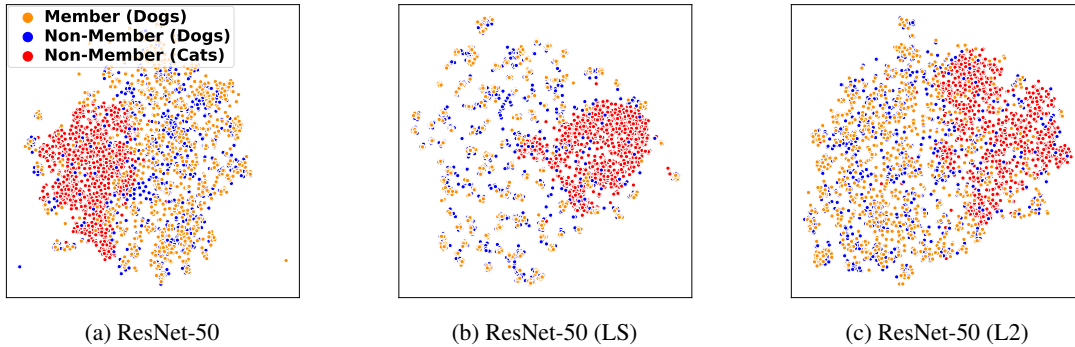(a) ResNet-50    (b) ResNet-50 (LS)    (c) ResNet-50 (L2)

Figure 7: T-SNE visualization of the penultimate ResNet-50 layer activations on training samples (green), test samples (blue), and OOD samples (red) as model inputs. Label smoothing (LS) creates much tighter clusters of training samples, while the OOD cat samples also form a tighter cluster and are easier to separate from the members. L2 regularization shows a reverted effect and increases overlapping, making it harder to separate members from non-members.

While the FPR on the Permuted samples is drastically reduced for ResNet-18, the FPR of the ResNet-50 on the Permuted samples even increases when using label smoothing. We note that this effect does only occur in some training runs. In other cases, the FPR for Permuted data drops similar to the ResNet-18 results. On all datasets, the reductions in the FPR are comparable between the ResNet-18 and ResNet-50. The FPR also decreases for inputs similar to the training data. For comparison, we also apply Laplace approximation (LA) on the weights of the final layers to mitigate overconfidence. As shown in Figs. 5c and 5d, LA is better suited to avoid high prediction scores on the Permuted and Scaled samples.

The attack results demonstrate that if a model shows reduced prediction scores on unseen inputs, the samples of the training data are easier to separate. It reduces the protection induced by overconfident predictions (on unseen inputs) and increases vulnerability to MIAs. We applied a kernel density estimation (KDE) to visualize the distribution of the maximum prediction scores of the ResNet-50 target models on member and non-member data from the Stanford

Dogs and the AFHQ Cats. Figs. 6a and 6b show the estimated density functions. Without label smoothing, all three distributions have their mode around prediction scores of 1.0. This leads to a large overlap of the distributions. Samples with prediction scores this high are most likely classified as false-positive members as the MMPS values in Tab. 3 suggest. We also state the earth mover's distance (EMD) in the KDE plots to quantify the distance between the member and non-member distributions. Label smoothing separates the three distributions clearly and doubles both EMD values. The label smoothing model tends to be less overconfident in its predictions on unknown input data, and hence the member samples are easier to separate from non-members. This increases the potential privacy leakage of MIAs.

As depicted in Fig. 7, we further used t-SNE [30] to plot the penultimate layer activations on samples from the same datasets as used for the KDE plots. While the standard model in Fig. 7a shows an overlapping between the activations of the three datasets, label smoothing in Fig. 7b creates tighter clusters of dog samples and separates the OOD cat images more clearly.

**(Q4) Defenses Are Contrary to Calibration.** While calibration tries to maximize the informative value of the prediction scores, many defenses against MIAs aim to reduce the informative value and to align the score distributions of members and non-members. They reduce the generalization of a model in terms of its ability to distinguish between samples from known and unknown inputs and express meaningful scores. First, we evaluated temperature scaling with $T = 10$ on top of the trained ResNet-50 standard model without calibration. Fig. 6c shows the estimated maximum prediction score distributions. The score vectors converge to a uniform distribution, and the distributions of the top scores are much more similar. This can be seen by the significantly lower EMD values. With an ECE of 51% using temperature scaling, the information content of the actual prediction score is greatly reduced. As seen in Tab. 2, temperature scaling is only effective against the maximum score and entropy attack. We suspect it is due to the softmax being a monotone transformation in the softmax space, not removing information encoded in the top-3 score patterns.

We also investigated L2 regularization as a stronger defense applied during training. We retrained the ResNets with a weight decay of $\lambda = 0.001$ for ResNet-50 and $\lambda = 0.0003$ for ResNet-18. L2 regularization effectively reduces the vulnerability to MIAs. For all attacks, both precision and recall drop significantly at the cost of reduced test accuracy, as Tab. 2 states. The distribution of the highest prediction scores can be seen in Fig. 6d. Similar to temperature scaling, L2 regularization aligns the distributions of members and non-members but distributes the maximum scores more equally instead of pushing it towards a single value. Fig. 7c shows a similar effect of overlapping distributions in the penultimate layer activations, making it harder to separate members from non-members and OOD data.

## 5    Discussion, Defense, and other Attacks

Our results demonstrate that previous works on prediction score-based MIAs overestimate the threat to privacy. Attack evaluations are usually done using cross-validation within the same data distribution not considering other data distributions. This introduces a strong evaluation bias since false-positive predictions are not taken into account. Our experiments point out that standard neural networks are not inherently able to identify inputs from unknown domains and cannot adapt their behavior in terms of reducing the prediction scores. This is why the expressiveness of MIAs in practical applications is reduced, and the associated risks to privacy are thus lower than previously assumed. In all our analyses, we assumed a strong adversary with full knowledge about the architecture and training procedure of the target model and with access to data from the target's train-

ing data distribution to fit the attacks. Loosening these assumptions even further decreases the power of MIAs.

One way to mitigate the problem of false-positive predictions on unseen data is to first identify and remove all OOD samples. This would indeed prevent false-positive predictions on completely different data distributions. However, we demonstrated that the problem of high FPR also occurs on more meaningful datasets with similar classes. Moreover, by generating synthetic images, we showed that there is a potentially unlimited number of samples that follow the training data distribution and result in false-positive predictions.

In this paper, we only considered prediction score-based MIAs but we expect our results to be similar for other kinds of attacks. Doing so provides an interesting avenue for future work. We also assume a high FPR for loss-based attacks since the loss is similarly based on prediction scores. We further expect a significant amount of false-positive predictions by label-only attacks, utilizing the distance of a sample to the decision boundary. The t-SNE plots in Fig. 7 show that the activations of members and non-members from different distributions are closely intertwined, and distances in the high dimensional space to the boundaries are probably also not clearly distinguishable. Still, we leave the empirical proof open for future work.

## 6    Conclusions

We have shown that MIAs produce high false-positive rates due to overconfident predictions of neural networks for in- and out-of-distribution data. In stark contrast to previous works stating strong attack results on standard neural networks, we demonstrate that the inference results are actually not reliable in realistic scenarios. Our results suggest that there is a trade-off between reducing a model's overconfidence and its susceptibility to MIAs. Mitigating overconfidence using calibration techniques leads to lower false-positive rates while precision increases and recall stays roughly the same. Therefore, the informative value of MIAs increases on calibrated models. Previously proposed defense methods against MIAs like temperature scaling and L2 regularization have an opposite effect and reduce the expressiveness of a model's prediction scores by aligning their distribution. Even though our experiments focus on score-based attacks, we want to draw attention to the fact that MIAs are not necessarily as powerful as previously thought and are at odds with the meaning of neural networks' prediction scores.

# References

[1] Raman Arora, Amitabh Basu, Poorya Mianjy, and Anirbit Mukherjee. Understanding deep neural networks with rectified linear units. In *ICLR*, 2018.

[2] Yunjey Choi, Youngjung Uh, Jaejun Yoo, and Jung-Woo Ha. Stargan v2: Diverse image synthesis for multiple domains. In *CVPR*, pages 8185–8194, 2020. doi: 10.1109/CVPR42600.2020.00821.

[3] Christopher A. Choquette-Choo, Florian Tramer, Nicholas Carlini, and Nicolas Papernot. Label-only membership inference attacks. In *ICML*, volume 139, pages 1964–1974, 18–24 Jul 2021.

[4] Adam Coates, Andrew Ng, and Honglak Lee. An analysis of single-layer networks in unsupervised feature learning. In *AISTATS*, volume 15, pages 215–223, 2011.

[5] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. On calibration of modern neural networks. In *ICML*, volume 70, pages 1321–1330, 2017.

[6] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016. doi: 10.1109/CVPR.2016.90.

[7] Matthias Hein, Maksym Andriushchenko, and Julian Bitterwolf. Why relu networks yield high-confidence predictions far away from the training data and how to mitigate the problem. In *CVPR*, pages 41–50, 2019. doi: 10.1109/CVPR.2019.00013.

[8] Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. In *ICLR*, 2017.

[9] Jinyuan Jia, Ahmed Salem, Michael Backes, Yang Zhang, and Neil Zhenqiang Gong. Memguard: Defending against black-box membership inference attacks via adversarial examples. In *CCS*, pages 259–274, 2019. doi: 10.1145/3319535.3363201.

[10] Tero Karras, Miika Aittala, Janne Hellsten, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Training generative adversarial networks with limited data. In *NeurIPS*, 2020.

[11] Aditya Khosla, Nityananda Jayadevaprakash, Bangpeng Yao, and Li Fei-Fei. Novel dataset for fine-grained image categorization. In *CVPR Workshop*, June 2011.

[12] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.

[13] Agustinus Kristiadi, Matthias Hein, and Philipp Hennig. Being bayesian, even just a bit, fixes overconfidence in relu networks. In *ICML*, volume 119, pages 5436–5446, 2020.

[14] Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, Master's thesis, Department of Computer Science, University of Toronto, 2009.

[15] C. Leibig, V. Allken, M. S. Ayhan, P. Berens, and S. Wahl. Leveraging uncertainty information from deep neural networks for disease detection. *Scientific Reports*, 7(0), 2019.

[16] Zheng Li and Yang Zhang. Membership leakage in label-only exposures, 2021.

[17] Rafael Müller, Simon Kornblith, and Geoffrey E. Hinton. When does label smoothing help? In *NeurIPS*, pages 4696–4705, 2019.

[18] Mahdi Pakdaman Naeini, Gregory F. Cooper, and Milos Hauskrecht. Obtaining well calibrated probabilities using bayesian binning. In *AAAI*, pages 2901–2907, 2015.

[19] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y. Ng. Reading digits in natural images with unsupervised feature learning. In *NeurIPS Workshop*, 2011.

[20] Anh Mai Nguyen, Jason Yosinski, and Jeff Clune. Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. In *CVPR*, pages 427–436, 2015. doi: 10.1109/CVPR.2015.7298640.

[21] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *NeurIPS*, pages 8024–8035, 2019.

[22] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.

[23] Prajit Ramachandran, Barret Zoph, and Quoc V. Le. Searching for activation functions. In *ICLR*, 2018.

[24] Shahbaz Rezaei and Xin Liu. On the difficulty of membership inference attacks. In *CVPR*, pages 7892–7900, 2021.

[25] Ahmed Salem, Yang Zhang, Mathias Humbert, Pascal Berrang, Mario Fritz, and Michael Backes. Ml-leaks: Model and data independent membership inference attacks and defenses on machine learning models. In *NDSS Symposium*, 2019.

[26] Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference attacks against machine learning models. In *IEEE Symposium on Security and Privacy*, pages 3–18, 2017. doi: 10.1109/SP.2017.41.

[27] Liwei Song and Prateek Mittal. Systematic evaluation of privacy risks of machine learning models. In *USENIX Security Symposium*, 2021.

[28] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *CVPR*, pages 2818–2826, 2016. doi: 10.1109/CVPR.2016.308.

[29] Mingxing Tan and Quoc V. Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *ICML*, volume 97, pages 6105–6114, 2019.

[30] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *JMLR*, 9(86):2579–2605, 2008.

[31] Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan J. van der Walt, Matthew Brett, Joshua Wilson, K. Jarrod Millman, Nikolay Mayorov, Andrew R. J. Nelson, Eric Jones, Robert Kern, Eric Larson, C J Carey, İlhan Polat, Yu Feng, Eric W. Moore, Jake VanderPlas, Denis Laxalde, Josef Perktold, Robert Cimrman, Ian Henriksen, E. A. Quintero, Charles R. Harris, Anne M. Archibald, Antônio H. Ribeiro, Fabian Pedregosa, Paul van Mulbregt, and SciPy 1.0 Contributors. SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. *Nature Methods*, 17:261–272, 2020. doi: 10.1038/s41592-019-0686-2.

[32] Michael L. Waskom. seaborn: statistical data visualization. *Journal of Open Source Software*, 6(60):3021, 2021. doi: 10.21105/joss.03021.

[33] Samuel Yeom, Irene Giacomelli, Matt Fredrikson, and Somesh Jha. Privacy risk in machine learning: Analyzing the connection to overfitting. In *IEEE Computer Security Foundations Symposium*, pages 268–282, 2018.

# A  Experimental Setup Details

We performed all our experiments on NVIDIA DGX machines running NVIDIA DGX Server Version 4.4.0 and Ubuntu 18.04 LTS. The machines have 1.6TB of RAM and contain Tesla V100-SXM3-32GB-H GPUs and Intel Xeon Platinum 8174 CPUs. We further relied on Python 3.8.8 and PyTorch 1.8.1 with Torchvision 0.9.1 [21] for the implementation and training of the neural networks. We provide a dockerfile together with our code to facilitate execution and reproducibility. We performed a single experimental run and set the seed for all experiments to 42 to allow reproducibility.

## A.1  Architectures

We use ResNet-50, ResNet-18, EfficientNetB0 and a custom CNN (SalemCNN) for our experiments.

**ResNet-50**: We use the ResNet-50 implementation and the ImageNet weights provided by PyTorch.

**ResNet-18 and EfficientNetB0**: For ResNet-18 and EfficientNetB0, we rely on the implementations provided at `https://github.com/kuangliu/pytorch-cifar` under MIT License. Note that this ResNet-18 and EfficientNetB0 implementations slightly differ from the official architectures since we train on CIFAR-10 instead of ImageNet. Differences occur mainly in early layers and show up in smaller kernel sizes and strides to avoid a large reduction in feature map sizes.

**SalemCNN**: Following Salem et al. [25], the model consists of two convolutional layers, each containing 32 filters of size 5 and a padding of 2 to maintain the spatial ratio. Each layer is followed by a ReLU activation and a $2 \times 2$ max pooling layer with stride 2 for downsampling. After that, two fully-connected layers further process the extracted features. The first fully-connected layer contains 128 neurons, while the number of neurons of the second one corresponds to the number of classes on the training set. Note that in its original version, the model uses a tanh activation on the first fully-connected layer. We change it to a ReLU to keep the network piecewise linear. We do not notice any significant performance differences between both variants during our experiments.

## A.2  Attacks

**Top-3 Score Attack** We use a simple neural network as an inference model for membership inference. The model consists of a neural network with one hidden layer containing 64 neurons and ReLU activations. Unlike Salem et al. [25], our inference model only uses a single output neuron, followed by a sigmoid function. During training, we first query the shadow model with samples with known membership status and collect the predicted scores. The values of each score vector are then sorted in descending order and the three highest values are used together with the membership status to train the inference model. The inference model is then trained on the membership dataset gathered by querying the shadow model and collecting the prediction score vectors. We use Adam optimizer [12] with learning rate 0.01, optimizing a binary cross-entropy loss. The training uses a batch size of 16 and is stopped if the loss is not decreasing by at least $5e^{-4}$ for 15 epochs.

**Maximum Prediction Score Attack** To find the threshold for the maximum prediction score attack a receiver operating characteristic (ROC) curve is created with the maximum values of each prediction score vector. We then choose the best threshold that maximizes the true-positive rate while minimizing the false-positive rate.

**Entropy Attack** For each of the prediction score vectors the entropy is calculated. To find the threshold a linear search is performed.

## A.3  Training Hyperparemters

To set the hyperparameters, we perform a small grid search for each model and dataset. Hence, we do not aim to achieve maximum test accuracy but to keep the training procedure simple. We roughly optimized the number of epochs {50, 100, 200}, learning rate {0.1, 0.05, 0.01, 0.005, 0.001, 0.0005, 0.0001} and optimizer type {Adam, SGD}. Since we usually expect the MIAs to perform worse on models with higher test accuracy, we argue that further hyperparameter optimization would degrade the attack metrics and strengthen our statements. We finally choose Adam optimizer [12] with default parameters ($\beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 1e - 08$) to train all our models. We use a batch size of 64 and a seed of 42 for training and experiments. All models are trained for 100 epochs.

**CIFAR-10 models**: For the SalemCNN, we set the learning rate to 0.001, for ResNet-18 and EfficientNetB0 we use a higher learning rate of 0.01.

**Stanford Dogs models**: We use pre-trained ImageNet weights for ResNet-50 and set the learning rate to 0.001. We replace the final fully-connected layer to match the number of classes.

## A.4  KDE Plots

We compute the earth mover's distance using the Wasserstein metric as implemented in SciPy 1.6.3 [31] with default parameters. Kernel density estimations are created with Seaborn 0.11.1 [32] and default parameters on 2,058 samples for each dataset.
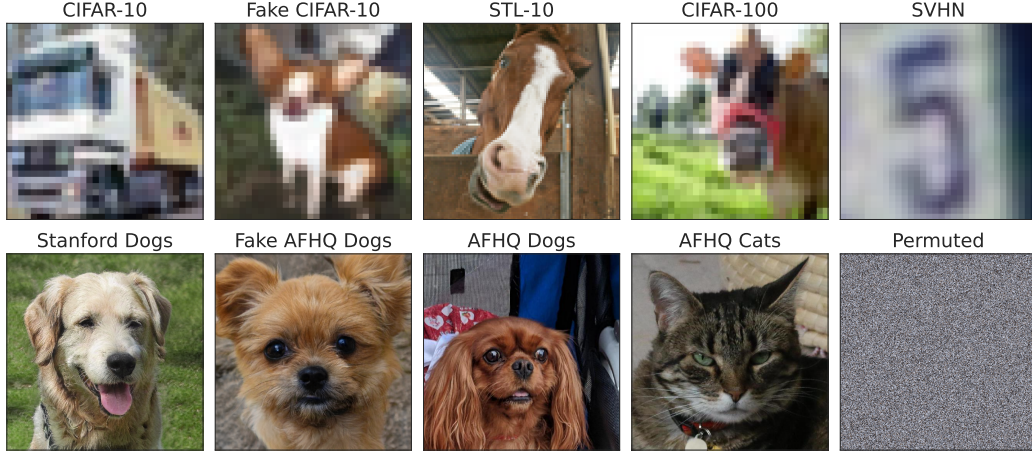
Figure 8: Randomly selected samples from the datasets without any preprocessing. At inference time, we scale all images to 32x32 for models trained on CIFAR-10 and 224x224 for models trained on Stanford Dogs. We further use StyleGAN2 to generate fake samples of CIFAR-10 and AFHQ Dogs, demonstrating a potentially infinite number of samples following a similar distribution.

## A.5 t-SNE Plots

We create all three t-SNE plots using scikit-learn 0.24.2 [22] with the same hyperparameters. Each embedding is initialized with a PCA. We set the perplexity to 30, the learning rate to 100, and the maximum number of iterations to 1,000. The random state is set to 13.

## A.6 Datasets

We normalize all input images on the statistics of the target model's training data by computing the standard score $z = \frac{x-\mu}{\sigma}$ of each input. We state exact parameters in table 4. For inference on CIFAR-trained models, we downsize samples from other distributions to 32x32 pixels. For models trained on Stanford Dogs, we resize all inputs to 224x224 pixels.

| Dataset | Mean | Std |
|---|---|---|
| CIFAR-10 | (0.4914, 0.4822, 0.4465) | (0.2470, 0.2435, 0.2616) |
| Stanford Dogs | (0.485, 0.456, 0.406) | (0.229, 0.224, 0.225) |

Table 4: Statistics for dataset normalization.

**CIFAR-10/ CIFAR-100** [14]: The CIFAR-10 and CIFAR-100 datasets each consist of 60,000 color images of size 32x32. Both training and test splits contain 50,000 and 10,000 samples, respectively. CIFAR-10 samples are grouped into 10 classes, CIFAR-100 correspondingly into 100 classes. The number of samples per class is completely balanced. CIFAR-10 contains samples from the classes airplane, automobile, bird, cat, deer, dog, frog, horse, ship and truck. More information is available at `https://www.cs.toronto.edu/~kriz/cifar.html`.

**Stanford Dogs Dataset** [11]: The Stanford Dogs dataset contains 20,580 images of 120 different dog breeds. Number of samples per class is not balanced. The dataset is built on ImageNet samples and has a significantly higher image resolution than the other datasets we use, except AFHQ. We do not rely on the official dataset split to increase the number of training samples for the target and shadow models. We use 80% of the data as training set, resulting in 16,464 samples for training and 4,116 samples as test data. We evaluate both, target and shadow models, on the full test split to keep it as large as possible. To improve generalization, we apply random rotation in a range of 20 degrees and resize the images so that the smaller side has a length of 230 pixels. We then randomly crop out a square image with size 224 pixels and flip it horizontally with 50% probability. Besides training, we only resize the inputs to 224 pixels on the shorter side and center crop to obtain a square image with size 224 pixels. The dataset and a list of its classes are available at `http://vision.stanford.edu/aditya86/ImageNetDogs/`.

**Animal Faces-HQ (AFHQ)** [2]: The AFHQ dataset contains 16,130 images of animal faces of size 512x512 and split into 14,630 training and 1,500 test samples. We use the training split for our experiments, providing 4,739 dog, 5,153 cat and 4,738 wild animal samples.

**STL-10** [4]: The STL-10 dataset is inspired by CIFAR-10 and contains 96x96 color images of ten different classes. The classes are identical to CIFAR-10, except class *monkey*. We therefore remove all samples containing monkeys. The full dataset contains a total of 5,000 labeled training samples, 8,000 test images and 100,000 unlabeled images from similar distributions for unsupervised learning. We use the training set for out experiments.

**SVHN** [19]: The Street View House Numbers (SVHN) dataset provides over 600,000 digit images of cropped house numbers in natural scene images. The dataset consists of 73,257 training and 26,032 test images. We use the training set for our experiments.

**Fake CIFAR-10**: We rely on a pre-trained class-conditional StyleGAN2 [10] to generate synthetic CIFAR-10 samples. We create a balanced dataset of 2,500 synthetic images, 250 for each class. Pre-trained StyleGAN models are available at `https://github.com/NVlabs/stylegan2-ada-pytorch`.

**Fake AFHQ Dogs**: Similarly to Fake CIFAR-10, we also create 2,500 synthetic dog images using another pre-trained StyleGAN2 trained on AFHQ Dog images using adaptive discriminator augmentation. Note that since AFHQ does not provide fine-granular labels for dog breeds, images are generated randomly without defining the dog breeds.

**Permuted**: We generate noisy images by randomly permutating the pixels of the non-members from the CIFAR-10 and Stanford Dog test sets. The resulting images do no longer contain any structural information.

**Scaled**: We scale samples from the non-members test set after normalization by factor 255 to imitate mistakes in the preprocessing steps. Samples of this dataset follow the theorem of [7] and correspond to a scaling factor $\delta = 255$.

### A.7 Evaluation Metrics

For evaluating our experiments precision, recall, false-positive rate (FPR), and mean maximum prediction scores (MMPS) are used. The first three metrics are based on the count of true-positive (TP), false-positive (FP) and false-negative (FN) predictions made by MIAs. The MMPS relies on the prediction score vector $f(x)$ produced by a neural network $f$ given input $x$. $f(x)$ includes the application of a softmax function to compute the prediction scores. We chose precision, recall and FPR since MIAs can be interpreted as a binary classification task. We further computed the MMPS to examine the influence of the maximum prediction scores on MIA classification decisions. The formulas we used to calculate the evaluation metrics can be seen below:

- $Precision = \frac{TP}{TP+FP}$

- $Recall = \frac{TP}{TP+FN}$

- $FPR = \frac{FP}{FP+TN}$

- $MMPS = \frac{1}{N} \sum_{n=1}^{N} \max_{i=1,...,d} f(x_n)_i$

## B  Additional Experimental Results

We here state additional results for our paper, that due to size restrictions, we were not able to include into the main part.

### B.1  MMPS for Stanford Dogs Models

Table 5 states the mean maximum prediction scores (MMPS) of false-positive and true-negative membership predictions for the default Stanford Dogs ResNet-50 model.

| Dataset | Attack | FP MMPS | TN MMPS |
|---|---|---|---|
| Stanford Dogs | Entropy | 0.9984 | 0.7565 |
| | Max. Score | 0.9985 | 0.7580 |
| | Top-3 Scores | 0.9986 | 0.7596 |
| Fake Dogs | Entropy | 0.9977 | 0.7700 |
| | Max. Score | 0.9979 | 0.7724 |
| | Top-3 Scores | 0.9980 | 0.7737 |
| AFHQ Dogs | Entropy | 0.9978 | 0.7636 |
| | Max. Score | 0.9980 | 0.7661 |
| | Top-3 Scores | 0.9980 | 0.7670 |
| AFHQ Cats | Entropy | 0.9972 | 0.7205 |
| | Max. Score | 0.9972 | 0.7208 |
| | Top-3 Scores | 0.9969 | 0.7193 |
| Permuted | Entropy | 0.9989 | 0.8237 |
| | Max. Score | 0.9991 | 0.8292 |
| | Top-3 Scores | 0.9992 | 0.8313 |
| Scaled | Entropy | 1.0000 | 0.8744 |
| | Max. Score | 1.0000 | 0.8744 |
| | Top-3 Scores | 1.0000 | 0.8744 |

Table 5: MMPS for false-positive (FP) and true-negative (TN) predictions on the standard ResNet-50 model.

13

## B.2 MMPS for Regularized Stanford Dogs Models

Table 6 states the MMPS of false-positive and true-negative membership predictions of the top-3 score attack against the various Stanford Dogs ResNet-50 models.

| Dataset | ResNet-50 | FP MMPS | TN MMPS |
|---|---|---|---|
| Stanford Dogs | Standard | 0.9986 | 0.7596 |
| | Label Smoothing | 0.9098 | 0.4819 |
| | L2 Regularization | 0.8492 | 0.4167 |
| | Temperature | 0.1367 | 0.0539 |
| Fake Dogs | Standard | 0.9980 | 0.7737 |
| | Label Smoothing | 0.8890 | 0.4508 |
| | L2 Regularization | 0.8376 | 0.4548 |
| | Temperature | 0.1481 | 0.0749 |
| AFHQ Dogs | Standard | 0.9980 | 0.7670 |
| | Label Smoothing | 0.9115 | 0.4613 |
| | L2 Regularization | 0.8359 | 0.4466 |
| | Temperature | 0.1479 | 0.0704 |
| AFHQ Cats | Standard | 0.9969 | 0.7193 |
| | Label Smoothing | 0.8930 | 0.3519 |
| | L2 Regularization | 0.8223 | 0.3736 |
| | Temperature | 0.0840 | 0.0444 |
| Permuted | Standard | 0.9992 | 0.8313 |
| | Label Smoothing | 0.9764 | 0.6280 |
| | L2 Regularization | 0.8449 | 0.4482 |
| | Temperature | 0.2257 | 0.1068 |
| Scaled | Standard | 1.0000 | 0.8744 |
| | Label Smoothing | 0.9885 | 0.8100 |
| | L2 Regularization | 0.9987 | 0.5486 |
| | Temperature | 0.6385 | 0.9924 |

Table 6: MMPS for false-positive and true-negative predictions the Top-3 scores attack against ResNet-50.

## B.3 Attack Results on ResNet-18

Table 7 states additional training and attack metrics for ResNet-18 trained on CIFAR-10.

| ResNet-18 | Standard | LS | LA | Temp | L2 |
|---|---|---|---|---|---|
| Train Accuracy | 100.00% | 100.00% | 100.00% | 100.00% | 68.48% |
| Test Accuracy | 69.38% | 71.94% | 69.08% | 69.38% | 58.04% |
| ECE | 24.02% | 13.33% | 8.65% | 22.04% | 16.95% |
| Entropy Pre | 67.35% | 79.60% | 68.02% | 63.47% | 51.43% |
| Entropy Rec | 92.32% | 94.56% | 51.72% | 84.44% | 28.72% |
| Entropy FPR | 44.76% | 24.24% | 24.32% | 48.60% | 27.12% |
| Max. Score Pre | 67.35% | 79.72% | 68.87% | 64.96% | 51.40% |
| Max. Score Rec | 92.32% | 94.48% | 63.64% | 90.76% | 28.64% |
| Max. Score FPR | 44.76% | 24.04% | 28.76% | 48.96% | 27.08% |
| Top-3 Scores Pre | 63.84% | 76.14% | 69.21% | 65.91% | 50.96% |
| Top-3 Scores Rec | 98.04% | 99.44% | 69.76% | 96.12% | 43.52% |
| Top-3 Scores FPR | 55.52% | 31.16% | 31.04% | 49.72% | 41.88% |

Table 7: Training and attack metrics for ResNet-18 target models trained on CIFAR-10.

## B.4 False-positive rates for CIFAR-10 Models

Table 8 states the underlying numerical FPR results for the models trained on CIFAR-10. All three attacks tend to predict unknown samples falsely as members even on data different from the training data distribution. Figure 9 further plots the FPR for the CIFAR-10 models and the effect of Laplace approximation and label smoothing, respectively.

## B.5 MMPS for CIFAR-10 Models

Table 9 states the mean maximum prediction scores (MMPS) of false-positive and true-negative membership predictions for CIFAR-10 models.



(a) ResNet-18 (LS)  (b) ResNet-18 (LA)

(c) SalemCNN (LS)  (d) SalemCNN (LA)

(e) EfficientNetB0 (LS)  (f) EfficientNetB0 (LA)
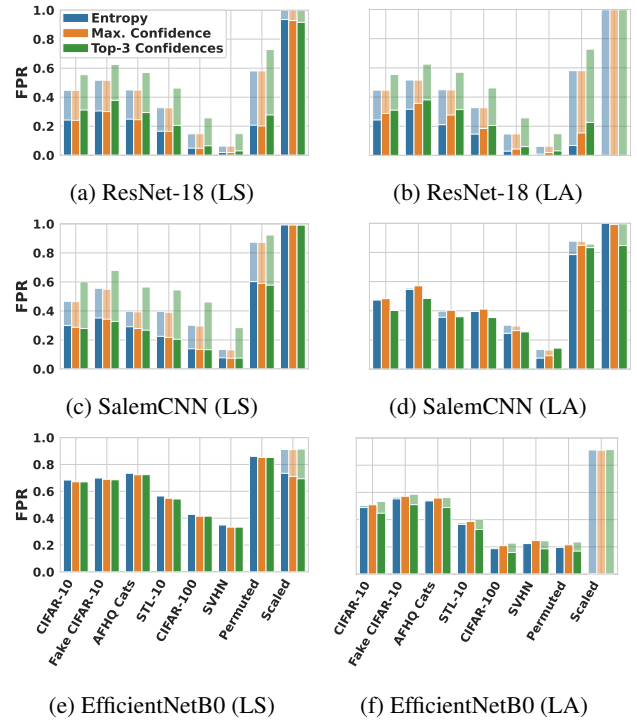
Figure 9: False-positive rates of score-based membership inference attacks against ResNet-18, SalemCNN, and the EfficientNetB0 model trained on the CIFAR-10 dataset. The transparent bars represent the false-positive rate of the standard models while the solid bars represent the false-positive rate of the models with Laplace approximation (LA) or label smoothing (LS), respectively.

| Attack | Architecture | CIFAR-10 | Fake CIFAR-10 | AFHQ Cats | STL-10 | CIFAR-100 | SVHN | Permuted | Scaled |
|---|---|---|---|---|---|---|---|---|---|
| | SalemCNN | 46.60% | 55.60% | 39.76% | 39.64% | 30.12% | 13.44% | 87.36% | 99.76% |
| Entropy | ResNet-18 | 44.76% | 51.68% | 45.00% | 32.72% | 14.72% | 6.28% | 58.08% | 99.84% |
| | EfficientNetB0 | 50.36% | 56.52% | 53.04% | 37.64% | 19.36% | 21.32% | 19.96% | 91.20% |
| | SalemCNN | 46.40% | 54.96% | 39.40% | 39.08% | 29.60% | 13.12% | 87.32% | 99.76% |
| Max. Score | ResNet-18 | 44.76% | 51.64% | 44.92% | 32.68% | 14.80% | 6.32% | 58.12% | 99.84% |
| | EfficientNetB0 | 50.00% | 55.84% | 52.44% | 37.28% | 19.00% | 20.76% | 19.72% | 91.04% |
| | SalemCNN | 60.04% | 67.96% | 56.48% | 54.44% | 46.12% | 28.48% | 92.32% | 99.80% |
| Top-3 Scores | ResNet-18 | 55.52% | 62.60% | 57.00% | 46.28% | 25.76% | 14.88% | 72.84% | 99.88% |
| | EfficientNetB0 | 53.40% | 58.68% | 56.24% | 40.24% | 22.76% | 24.44% | 23.60% | 91.48% |

Table 8: False-positive rates (FPR) for score-based membership inference attacks against CIFAR-10 target models.

| | | Entropy | | Max. Score | | Top-3 Scores | |
|---|---|---|---|---|---|---|---|
| | | FP MMPS | TN MMPS | FP MMPS | TN MMPS | FP MMPS | TN MMPS |
| | SalemCNN | 1.0000 | 0.9242 | 1.0000 | 0.9245 | 1.0000 | 0.8987 |
| CIFAR-10 | ResNet-18 | 1.0000 | 0.8873 | 1.0000 | 0.8873 | 0.9999 | 0.8601 |
| | EfficientNetB0 | 0.9999 | 0.8575 | 0.9999 | 0.8585 | 0.9998 | 0.8483 |
| | SalemCNN | 1.0000 | 0.9212 | 1.0000 | 0.9224 | 1.0000 | 0.8909 |
| Fake CIFAR-10 | ResNet-18 | 1.0000 | 0.8988 | 1.0000 | 0.8989 | 0.9999 | 0.8694 |
| | EfficientNetB0 | 0.9999 | 0.8649 | 0.9999 | 0.8670 | 0.9998 | 0.8579 |
| | SalemCNN | 1.0000 | 0.9178 | 1.0000 | 0.9183 | 1.0000 | 0.8862 |
| AFHQ Cats | ResNet-18 | 1.0000 | 0.9051 | 1.0000 | 0.9053 | 0.9999 | 0.8787 |
| | EfficientNetB0 | 0.9999 | 0.8866 | 0.9999 | 0.8880 | 0.9998 | 0.8784 |
| | SalemCNN | 1.0000 | 0.9144 | 1.0000 | 0.9152 | 1.0000 | 0.8866 |
| STL-10 | ResNet-18 | 1.0000 | 0.8841 | 1.0000 | 0.8842 | 0.9999 | 0.8549 |
| | EfficientNetB0 | 0.9999 | 0.8428 | 0.9999 | 0.8437 | 0.9998 | 0.8360 |
| | SalemCNN | 1.0000 | 0.9089 | 1.0000 | 0.9095 | 1.0000 | 0.8819 |
| CIFAR-100 | ResNet-18 | 1.0000 | 0.8560 | 1.0000 | 0.8559 | 0.9999 | 0.8346 |
| | EfficientNetB0 | 0.9997 | 0.8258 | 0.9998 | 0.8266 | 0.9995 | 0.8182 |
| | SalemCNN | 1.0000 | 0.8926 | 1.0000 | 0.8930 | 0.9999 | 0.8700 |
| SVHN | ResNet-18 | 1.0000 | 0.8394 | 1.0000 | 0.8393 | 0.9998 | 0.8232 |
| | EfficientNetB0 | 0.9997 | 0.8301 | 0.9998 | 0.8313 | 0.9996 | 0.8231 |
| | SalemCNN | 1.0000 | 0.9405 | 1.0000 | 0.9407 | 1.0000 | 0.9022 |
| Permuted | ResNet-18 | 1.0000 | 0.9323 | 1.0000 | 0.9322 | 0.9999 | 0.8956 |
| | EfficientNetB0 | 0.9997 | 0.8086 | 0.9997 | 0.8092 | 0.9995 | 0.7995 |
| | SalemCNN | 1.0000 | 0.8329 | 1.0000 | 0.8329 | 1.0000 | 0.7995 |
| Scaled | ResNet-18 | 1.0000 | 0.8926 | 1.0000 | 0.8926 | 1.0000 | 0.8568 |
| | EfficientNetB0 | 1.0000 | 0.8922 | 1.0000 | 0.8941 | 1.0000 | 0.8887 |

Table 9: Mean maximum prediction scores (MMPS) for false-positive (FP) and true-negative (TN) member predictions for CIFAR-10 models.