# Practical Black-Box Attacks against Deep Learning Systems using Adversarial Examples

**Nicolas Papernot** - *The Pennsylvania State University*            ngp5056@cse.psu.edu
**Patrick McDaniel** - *The Pennsylvania State University*            mcdaniel@cse.psu.edu
**Ian Goodfellow** - *Google Inc.*            goodfellow@google.com
**Somesh Jha** - *University of Wisconsin-Madison*            jha@cs.wisc.edu
**Z. Berkay Celik** - *The Pennsylvania State University*            zbc102@cse.psu.edu
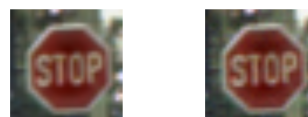**Ananthram Swami** - *US Army Research Laboratory*            ananthram.swami.civ@mail.mil

**Abstract** - *Advances in deep learning have led to the broad adoption of Deep Neural Networks (DNNs) to a range of important machine learning problems, e.g., guiding autonomous vehicles, speech recognition, malware detection. Yet, machine learning models, including DNNs, were shown to be vulnerable to adversarial samples—subtly (and often humanly indistinguishably) modified malicious inputs crafted to compromise the integrity of their outputs. Adversarial examples thus enable adversaries to manipulate system behaviors. Potential attacks include attempts to control the behavior of vehicles, have spam content identified as legitimate content, or have malware identified as legitimate software. Adversarial examples are known to transfer from one model to another, even if the second model has a different architecture or was trained on a different set. We introduce the first practical demonstration that this cross-model transfer phenomenon enables attackers to control a remotely hosted DNN with no access to the model, its parameters, or its training data. In our demonstration, we only assume that the adversary can observe outputs from the target DNN given inputs chosen by the adversary. We introduce the attack strategy of fitting a substitute model to the input-output pairs in this manner, then crafting adversarial examples based on this auxiliary model. We evaluate the approach on existing DNN datasets and real-world settings. In one experiment, we force a DNN supported by MetaMind (one of the online APIs for DNN classifiers) to mis-classify inputs at a rate of 84.24%. We conclude with experiments exploring why adversarial samples transfer between DNNs, and a discussion on the applicability of our attack when targeting machine learning algorithms distinct from DNNs.*

## 1   Introduction

*Deep Neural Networks* transformed machine learning by outperforming existing approaches at classification [17] and reinforcement learning [27], as well as allowing for breakthroughs in unsupervised learning [35]. These advances have generated a wide interest in *deep learning*, the machine learning technique which uses DNNs to represent a hierarchy of increasingly complex concepts [14]. Thus, deep learning has been broadly adapted to automotive systems [29], finance [22], healthcare [1], computer vision [23], speech recognition [18], natural language processing [32], and security [11, 36].

It has been found that a wide variety of machine learning algorithms are vulnerable to integrity attacks [38, 15, 30]. Namely, *adversarial samples* are inputs crafted by adding small, often imperceptible, perturbations to force a trained machine learning model to misclassify the resulting inputs, while remaining correctly classified by a human observer. To illustrate, consider the following images, potentially consumed by an autonomous vehicle:



To humans, these two images appear to be the same— our biological classifiers (vision) identify each image as a stop sign. The image on the left [37] is indeed an ordinary image of a stop sign. We produced the image on the right by adding a small, precise perturbation that forces a particular image-classification DNN to classify it as a yield sign. Here, the adversary could potentially use the altered image to cause the car to behave dangerously, if the car did not have failsafes such as maps of known stop sign locations. This hypothetical attack would require modifying the image used internally by the car—it is not yet known whether adversarial example generation requires greater control over the image than can be obtained by physically altering the sign. It is conceivable that such an image could be generated by maliciously modifying the sign itself, but it is also conceivable that low bit depth in the camera, compression, lens artifacts, and the variety of viewing angles could make it harder to give the sign a physical appearance deceiving the model.

To evaluate perturbations required to craft adversarial samples, previous approaches used detailed knowledge of the DNN architecture and parameters [15, 30, 38], or an independently collected training set to fit an auxiliary model [15, 38]. This limited their applicability to strong adversaries with the capability of gaining insider knowledge of the targeted deep learning based system or of collecting and labeling large training sets. Moreover, previous work studied cross-model transfer from the point of view of what it implies about the operation of machine learning models. To study it from the perspective of security, it is important to carry out correctly blinded experiments, in which the experimenters truly do not know the architecture of the model being attacked.

In this paper, we consider a weak adversary whose capabilities are limited to accessing the targeted model's output. We construct attacks against a remotely hosted machine learning system whose architecture is unknown to us, thus providing the first correctly blinded experiment concerning adversarial examples as a security risk. Specifically, the adversary has access to an *oracle* returning the DNN output for any chosen input. The adversarial goal is to force the oracle to misclassify inputs.

Our approach builds on previous attacks to craft adversarial samples under this more complex threat model. More precisely, we introduce a substitute DNN approximating the targeted oracle's learned model, thus skirting the lack of knowledge of the oracle's architecture and training set. Training a substitute DNN gives us the benefit of having full access to its inner details, enabling us to employ previous adversarial sample crafting methods. To train the substitute without access to the oracle training set, we introduce a Jacobian-based dataset augmentation technique. It allows us to select points in the input domain that are representative of the oracle model. We then leverage the *transferability* of adversarial samples [38]: adversarial samples crafted for a given machine learning model are also largely misclassified by other models trained on the same task. Thus, adversarial samples crafted using our substitute DNN are largely misclassified by the targeted oracle, as evaluated in Section 5. The contributions of this paper are the following:

- We introduce in Section 4 an attack targeting the integrity of machine learning oracles, which stems from training a substitute DNN using a Jacobian-based dataset augmentation. Adversaries can thus craft adversarial samples misclassified by the oracle without knowledge of its architecture or training set.
- Our dataset augmentation makes the attack tractable by limiting oracle querying. We show in Section 5 how to force the MetaMind API for DNN classifiers to mis-classify 84.24% of inputs altered with perturbations hard to distinguish for humans. We

also force a second oracle trained on a traffic sign dataset to misclassify 64.24% of inputs crafted.
- We calibrate our attack algorithm in Section 6. We show that the choice of substitute architecture does not significantly impact adversarial sample transferability. We fine-tune perturbations introduced in order to maximize adversarial sample transferability.
- We provide an intuition of why adversarial samples transfer across DNN architectures by empirically observing that substitute models have cost gradient sign matrices that are correlated to the oracle's.

## 2  About Deep Neural Networks

In this section, we provide preliminaries of deep learning to enable understanding of our threat model and attack strategy. We refer readers interested to the detailed presentation of deep neural networks in [14].

A *Deep Neural Network*, as illustrated in its simplest form in Figure 1, is a machine learning technique that uses a hierarchical composition of $n$ parametric functions to model a high dimensional input $\vec{x}$ [6, 14]. Each function $f_i$ for $i \in 1..n$ is modeled using a layer of neurons, which are essentially elementary computing units applying an *activation function* to the previous layer's weighted representation of the input to generate a new representation. Each neuron layer is parameterized by a weight vector $\theta_i$ impacting the activation of each neuron. Such weights essentially hold the knowledge of a DNN model $F$ and are evaluated during its training phase, as detailed below. Thus, the composition of functions modeled by deep neural networks can be formalized as:

$$F(\vec{x}) = f_n(\theta_n, f_{n-1}(\theta_{n-1}, \dots f_2(\theta_2, f_1(\theta_1, \vec{x})))) \quad (1)$$

The *training phase* of a neural network $F$ learns values for its set of parameters $\theta_F = \{\theta_1, ..., \theta_n\}$. In this paper, we focus on the case of supervised learning where the task of the DNN is to model a known input-output relation. One example of such a task is classification, where the goal is to assign inputs a label among a predefined set of labels. To do so, the network is given a large set of known input-output pairs $(\vec{x}, \vec{y})$ and it adjusts weight parameters to reduce a cost quantifying the prediction error between the model prediction $F(\vec{x})$ and the correct output $\vec{y}$. The adjustment is typically performed using techniques derived from the backpropagation algorithm [39]. Briefly speaking, such techniques successively propagate error gradients with respect to network parameters from the output layer of the network to its input layer.

During the *test phase*, the DNN is deployed with a fixed set of parameters $\theta_F$ and is used to make predictions on inputs unseen during training. For instance, in our paper we consider DNNs used as classifiers: for a
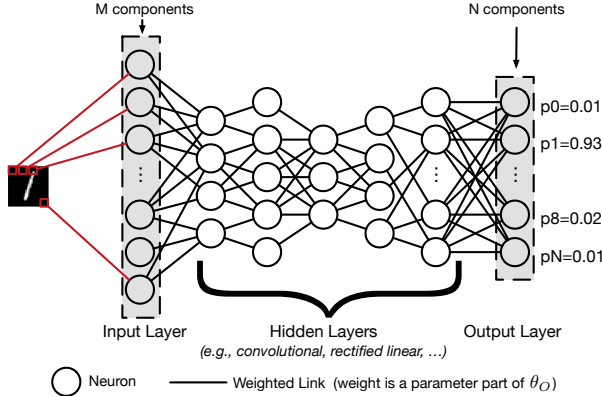
Figure 1: **DNN Classifier Architecture:** a DNN is a succession of neuron layers modeling representations of the input and building to an expected output. This model maps images of digits with probability vectors indicating the digit identified in the image (adapted from [31]).

given input $\vec{x}$, the neural network produces a probability vector $F(\vec{x})$ encoding its belief of input $\vec{x}$ being in each of the predefined classes (cf. Figure 1). The weight parameters $\theta_F$ hold the knowledge acquired by the model during training. Ideally, the learning algorithm used during the training phase should allow the model to generalize so as to make accurate predictions for inputs outside of the domain explored during training. However, this is not the case in practice as shown by previous attacks manipulating DNN outputs using adversarial samples [15, 30, 38].

## 3 Threat Model

We now describe our threat model for attacks against deep learning oracles. The adversary's goal is to force a targeted DNN to misbehave in a controlled manner: for instance, if the DNN is used for classification, the adversary seeks to force it to misclassify inputs in classes different from their correct classes. To achieve this goal, we consider a weak adversary only capable of accessing the targeted DNN's output. In other words, the adversary has no knowledge of the architectural choices made to design the DNN, which include the number, type, and size of layers, nor of the training data used to evaluate the parameters of the DNN.[1] Previous attacks depended on a knowledge of the DNN internals [15, 30, 38] or cross-model transferability from models learned using training sets hard to collect and expensive to label [15, 38]. Hence, the black-box threat model explored in this paper represents a more serious attack vector on deep learning.

A taxonomy of threat models for deep learning de-

---

[1] Such attacks are referred to as *black box attacks*, where adversaries need not know internal details of a targeted system to compromise it.

ployed in adversarial settings is introduced in [30]. In the present paper, we consider attackers targeting a DNN used as a multi-class classifier. Such a DNN produces outputs taking the form of probability vectors. Each vector component encodes the network's belief of the input being part of one of the predefined classes. To illustrate our attack scenario, we consider the ongoing example of a DNN used for classifying images, as illustrated in Figure 1. Such an architecture can be used to classify handwritten digits into classes corresponding to digits between 0 and 9, images of objects in a fixed number of categories, or images of traffic signs into classes corresponding to sign types (STOP, yield, ...).

**Adversarial Capabilities -** The targeted DNN is referred to as the *oracle* and denoted by $O$. The name oracle refers to the only capability of the adversary: accessing the DNN's label $\tilde{O}(\vec{x})$ for any input $\vec{x}$ by querying the oracle $O$. In this paper, we consider the output $\tilde{O}(\vec{x})$ to take the form of a label, which is defined as the index of the class assigned the largest probability by the DNN:

$$\tilde{O}(\vec{x}) = \arg \max_{j \in 0..N-1} O_j(\vec{x}) \qquad (2)$$

where $O_j(\vec{x})$ is the j-th component of the probability vector $O(\vec{x})$ output by network $O$ on input $\vec{x}$. The distinction between labels and probabilities is important as it makes adversaries more realistic–users more often have access to output labels than output vectors in real world applications–but weaker as labels encode much less information about the model's learned behavior than probability vectors. Accessing labels $\tilde{O}$ is the only capability of adversaries considered in our threat model. As such, we do not assume that the adversary has access to the oracle's architecture, parameters, or training data. Again, this weaker adversary makes our threat model more realistic but in turn makes the attack much harder to execute.

**Adversarial Goal -** The adversarial goal is to take any input $\vec{x}$ and produce a minimally altered version of $\vec{x}$, named *adversarial sample* and denoted by $\vec{x^*}$, that has the property of being misclassified by oracle $O$. This corresponds to an attack against the oracle's output integrity. Thus, the adversarial sample is such that $\tilde{O}(\vec{x^*}) \neq \tilde{O}(\vec{x})$, and solves the following optimization problem:

$$\vec{x^*} = \vec{x} + \arg \min \{\vec{z} : \tilde{O}(\vec{x} + \vec{z}) \neq \tilde{O}(\vec{x})\} = \vec{x} + \delta_{\vec{x}} \quad (3)$$

Some examples of adversarial samples can be found in Figure 2. The first row contains legitimate samples while the second row contains corresponding adversarial samples that are misclassified. This misclassification must be achieved by adding a minimal perturbation $\delta\vec{x}$ so as to evade human detection. Even with total knowledge of the network architecture used to train model $O$ as well as the parameters resulting from training, finding such a
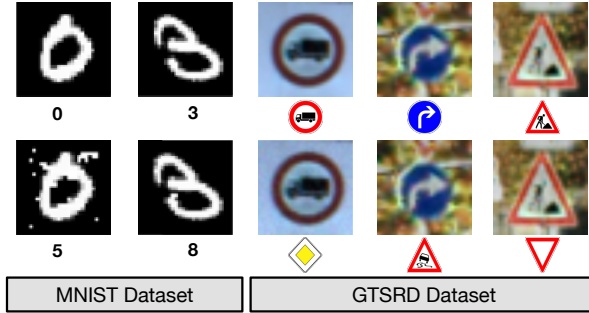
Figure 2: **Adversarial Samples:** the top row contains legitimate samples [25, 37] used to create the adversarial samples in the bottow row, which are misclassified by a DNN. The DNN output is identified below each sample.

minimal perturbation is not trivial, as properties of deep neural networks preclude the optimization problem from being linear or convex. This problem is exacerbated by our threat model: removing knowledge of model $O$'s architecture and training dataset makes it even harder to find a perturbation such that $\tilde{O}(\vec{x} + \delta\vec{x}) \neq \tilde{O}(\vec{x})$ holds.

## 4   Attack Methodology

We introduce here an attack against DNN oracles under the threat model described in Section 3. The adversary wants to craft adversarial samples misclassified by the oracle $\tilde{O}$ using his/her sole capability of accessing the label assigned by oracle $\tilde{O}$ for any chosen input. Briefly speaking, the attack strategy consists in learning a *substitute DNN* approximating the targeted oracle and then using this substitute model to craft adversarial samples. As adversarial samples transfer between architectures, we expect these samples to be misclassified by the oracle.

To understand the difficulty of conducting the attack under this threat model, recall Equation 3 formalizing the adversarial goal of finding a minimal perturbation $\delta\vec{x}$ resulting in an adversarial sample misclassified by the targeted DNN oracle $O$. Unfortunately for attackers but fortunately for defenders, Equation 3 cannot be solved in closed form. The basis for most adversarial attacks [15, 30, 38] is to approximate the solution to Equation 3 by using gradient-based optimization on functions defined by the DNN. Because evaluating these functions and their gradients requires knowledge of the DNN architecture and parameters, such an attack is not possible under our black-box scenario. In [38], it was shown that adversaries with access to an independently collected labeled training set from the same population distribution than the oracle could train a model with a different architecture and use it as a substitute: adversarial examples designed to manipulate the substitute's out-

put often manipulate the targeted model's output. However, many modern machine learning systems require extremely large and expensive training sets to fit reliably. For instance, in our experiments, we consider models trained with several tens of thousands of labeled examples. This makes attacks based on this paradigm unfeasible for adversaries without large labeled datasets.

In this paper, we demonstrate that black-box attacks can be accomplished at a much lower cost, without the construction of an independent labeled training set. In our approach, we use the target as an oracle to provide the labels used to train the substitute model. The attacker builds an accurate approximation $F$ of the model $O$ learned by the oracle, then uses this approximated model $F$, which we name the *substitute network*, to craft adversarial samples misclassified by $F$ that will in turn transfer and be misclassified by the oracle $O$.

Indeed, as adversaries have full knowledge of the DNN architecture and parameters of the substitute network, they can use one of the previously described attack techniques [15, 30] to craft adversarial samples misclassified by $F$. *As long as the transferability property holds between F and O, adversarial samples crafted for F will also be misclassified by O.* This leads us to propose the following two-fold attack strategy:

1. **Substitute Model Training:** the attacker queries the oracle with inputs selected by a Jacobian-based heuristic to build a model $F$ approximating the oracle model $O$'s decision boundaries.

2. **Adversarial Sample Crafting:** the attacker uses substitute network $F$ to craft adversarial samples, which are then misclassified by oracle $O$ due to the transferability of adversarial samples.

The most crucial step is of training the substitute model using the very limited knowledge available.

### 4.1   Substitute Model Training

Training a model $F$ approximating oracle $O$ is challenging because: (1) we must select an architecture for our substitute DNN without knowledge of the targeted oracle's architecture, and (2) we must limit the number of queries made to the oracle in order to ensure that the approach is tractable. Our approach overcomes these challenges mainly by introducing a dataset augmentation technique, which we name *Jacobian-based Dataset Augmentation*. Let us emphasize that *our augmentation technique is not aimed at improving the substitute DNN's accuracy but rather ensuring that it approximates the oracle's decision boundaries*. The training procedure described below is illustrated in Figure 3.

**Substitute Architecture -** This factor is not the most limiting as the adversary must at least have some partial
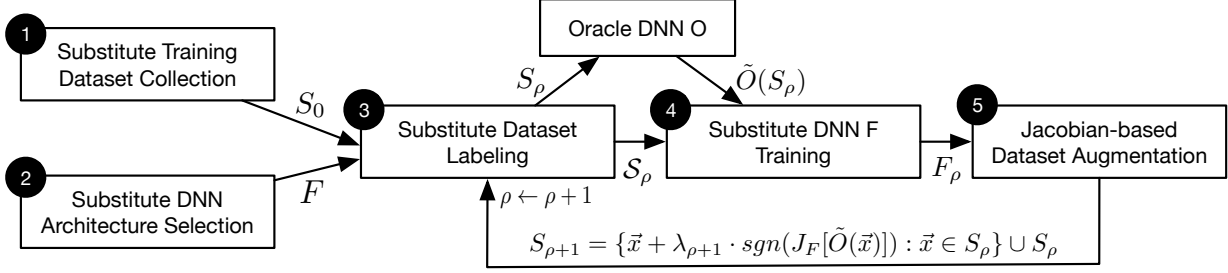
Figure 3: **Training of the Substitute DNN Architecture $F$**: the attacker (1) collects an initial substitute training set $S_0$ and (2) selects a substitute architecture $F$. Using the oracle $\tilde{O}$, the attacker (3) labels $S_0$ and (4) trains substitute DNN $F$. After (5) Jacobian-based dataset augmentation, steps (3) through (5) are repeated for several substitute epochs $\rho$.

knowledge of the oracle input (e.g., images, text, documents) and expected output (e.g., multi-class classification). The adversary can thus use an architecture adapted to the input-output relation. For instance, a convolutional neural network is suitable for image classification. Furthermore, we show in Section 6 that the type, number, and size of layers used in the substitute architecture has relatively little impact on the success of the attack. Adversaries can also consider performing an architecture exploration and train several substitute models before selecting the one yielding the highest attack success rate.

**Exploration of the Input Domain -** To better understand factor (2), note that we could potentially make an infinite number of queries to obtain the oracle's output $O(\vec{x})$ for any input $\vec{x}$ belonging to the input domain. This would provide us with a copy of the oracle. However, this is simply not tractable: consider a DNN with $M$ input components, each taking discrete values among a set of $K$ possible values, the number of possible inputs to be queried is $K^M$. The intractability is even more apparent for inputs in the continuous domain. Furthermore, making a large number of queries renders the adversarial behavior easy to detect by a defender. In order to address this issue, we therefore use a heuristic that performs an efficient exploration of the input domain and, as shown in Sections 5 and 6, drastically limits the number of oracle queries. Furthermore, our technique also ensures that the substitute network is an accurate approximation of the targeted oracle, i.e. it learns similar decision boundaries.

The heuristic is based on identifying directions in which the model's output is varying, around an initial set of training points. The intuition is that such directions will require more input-output pairs (produced by querying the oracle) to capture the output variations from model $O$. Therefore, to get a substitute network accurately approximating the oracle's decision boundaries, the heuristic prioritizes these samples for queries made to the oracle. To identify such directions, the substitute DNN's Jacobian matrix $J_F$ is evaluated at several input

points (the method used to choose these points is described later in this section). Directions, in which the substitute network $F$ varies for any input $\vec{x}$ are then found by computing the following:

$$\text{sgn}\left(J_F(\vec{x})[\tilde{O}(\vec{x})]\right) \qquad (4)$$

which is the sign of the Jacobian matrix's dimension corresponding to the oracle's label for input $\vec{x}$. We name this technique *Jacobian-based Dataset Augmentation*. We base our training algorithm on the idea of iteratively refining the model in directions identified using the Jacobian by making additional queries to the oracle.

---

**Algorithm 1 - Substitute DNN Training:** the algorithm takes as inputs the oracle $\tilde{O}$, the maximum number $max_\rho$ of substitute training epochs to be performed, a substitute architecture $F$, and an initial training set $S_0$.

---

**Input:** $\tilde{O}, max_\rho, S_0$
1: Define architecture $F$
2: **for** $\rho \in 0 .. max_\rho - 1$ **do**
3:     // *Label the substitute training set*
4:     $D \leftarrow \left\{ (\vec{x}, \tilde{O}(\vec{x})) : \vec{x} \in S_\rho \right\}$
5:     // *Train F on D to evaluate parameters $\theta_F$*
6:     $\theta_F \leftarrow \text{train}(F, D)$
7:     // *Perform Jacobian-based dataset augmentation*
8:     $S_{\rho+1} \leftarrow \{\vec{x} + \lambda_{\rho+1} \cdot \text{sgn}(J_F[\tilde{O}(\vec{x})]) : \vec{x} \in S_\rho\} \cup S_\rho$
9: **end for**
10: **return** $\theta_F$

---

**Substitute DNN Training Algorithm -** The resulting training procedure is a five step process outlined in Algorithm 1 and described more specifically below:

1. **Initial Collection:** The adversary first collects a very limited number of inputs representative of the input domain. For instance, if the targeted oracle $O$ classifies handwritten digits, the adversary collects 10 images of each digit 0 through 9. This constitutes the initial training set $S_0$ for the substitute DNN.

We show in Section 5 that the substitute training set does not necessarily have to come from the distribution from which the targeted oracle was trained.

2. **Architecture Selection:** The adversary selects a DNN architecture to be trained as the substitute model $F$. Again, this can be done using high-level knowledge of the classification task performed by the targeted oracle (e.g., convolutional neural networks offer good performance for vision tasks)

3. **Labeling:** By querying for the labels $\tilde{O}(\vec{x})$ output by oracle $O$, the adversary labels each sample $\vec{x} \in S_0$ in its initial substitute training set $S_0$.

4. **Training:** The adversary trains the architecture chosen at step (2) using the substitute training set $S_0$. This is performed using the classical training techniques used for DNNs.

5. **Augmentation:** The adversary performs Jacobian-based Dataset Augmentation on the initial substitute training set $S_0$ to produce a larger substitute training set $S_1$. This new training set better represents the model's decision boundaries. The adversary repeats steps (3) and (4) with the augmented set $S_1$.

Note that steps (3),(4), and (5) can be repeated several times to increase both the accuracy of substitute network $F$ and the similarity of its decision boundaries with those of the oracle. We introduce the notion of *substitute training epoch* to refer to each iteration performed, and index epochs with $\rho$. This leads to the following formalization of the Jacobian-based Dataset Augmentation performed at step (5) of our substitute DNN training algorithm:

$$S_{\rho+1} = \{\vec{x} + \lambda_{\rho+1} \cdot \text{sgn}(J_F[\tilde{O}(\vec{x})]) : \vec{x} \in S_\rho\} \cup S_\rho \quad (5)$$

where $\lambda_{\rho+1}$ is a parameter of the augmentation technique: it defines the size of the step taken in the sensitive direction identified by the sign of the Jacobian matrix to augment the dataset $S_\rho$ into $S_{\rho+1}$.

## 4.2 Adversarial Sample Crafting

Once the adversary has trained a substitute DNN, the next step in our strategy is to craft adversarial samples. This step is performed by implementing two previously introduced approaches described in [15, 30]. We now provide an overview of both adversarial sample crafting techniques, namely the *Goodfellow et al. algorithm* and the *Papernot et al. algorithm*. Both approaches share a similar intuition of evaluating the model's sensitivity to input components in order to select a small perturbation achieving the adversarial misclassification goal.

**Goodfellow et al. attack -** This algorithm is also known as the *fast gradient sign method* [15]. Given a model $F$

with an associated cost function $c(F, \vec{x}, y)$, the adversary crafts an adversarial sample $\vec{x^*} = \vec{x} + \delta_{\vec{x}}$ for a given legitimate sample $\vec{x}$ by computing the following perturbation:

$$\delta_{\vec{x}} = \varepsilon \, \text{sgn}(\nabla_{\vec{x}} c(F, \vec{x}, y)) \quad (6)$$

where perturbation $\text{sgn}(\nabla_{\vec{x}} c(F, \vec{x}, y))$ is the sign of the model's cost function[2] gradient. The cost gradient is computed with respect to $\vec{x}$ using sample $\vec{x}$ and label $y$ as inputs. The value of the *input variation parameter $\varepsilon$* factoring the sign matrix controls the perturbation's amplitude. Increasing its value increases the likelihood of $\vec{x^*}$ being misclassified by model $F$ but on the contrary makes adversarial samples easier to detect by humans. In Section 6, we evaluate the impact of parameter $\varepsilon$ on the successfulness of our DNN oracle attack.

**Papernot et al. attack -** This algorithm is suitable for source-target misclassification attacks where adversaries seek to take samples from any legitimate source class to any chosen target class [30]. The misclassification attack is a special case of the source-target misclassification attack where the target class can be any class different from the legitimate source class. Given a model $F$, the adversary crafts an adversarial sample $\vec{x^*} = \vec{x} + \delta_{\vec{x}}$ for a given legitimate sample $\vec{x}$ by adding a perturbation. Perturbation $\delta_{\vec{x}}$ is a subset of the input components $\vec{x}_i$.

To choose input components forming perturbation $\delta_{\vec{x}}$, components are sorted by decreasing adversarial saliency value. The adversarial saliency value $S(\vec{x}, t)[i]$ of component $i$ for an adversarial target class $t$ is defined as:

$$S(\vec{x}, t)[i] = \begin{cases} 0 \text{ if } \frac{\partial F_t}{\partial \vec{x}_i}(\vec{x}) < 0 \text{ or } \sum_{j \neq t} \frac{\partial F_j}{\partial \vec{x}_i}(\vec{x}) > 0 \\ \frac{\partial F_t}{\partial \vec{x}_i}(\vec{x}) \left| \sum_{j \neq t} \frac{\partial F_j}{\partial \vec{x}_i}(\vec{x}) \right| \text{ otherwise} \end{cases} \quad (7)$$

where matrix $J_F = \left[\frac{\partial F_j}{\partial \vec{x}_i}\right]_{ij}$ is the model's Jacobian matrix. Input components $i$ are added to perturbation $\delta_{\vec{x}}$ in order of decreasing adversarial saliency value $S(\vec{x}, t)[i]$ until the resulting adversarial sample $\vec{x^*} = \vec{x} + \delta_{\vec{x}}$ is misclassified by model $F$. The perturbation introduced for each selected input component can vary. Greater individual variations tend to reduce the number of components perturbed to achieve misclassification.

Each attack algorithm has its benefits and drawbacks. The Goodfellow algorithm is well suited for fast crafting of large amounts of adversarial samples with relatively large perturbations that are thus relatively easy for a defender to notice. The Papernot algorithm reduces the perturbation at the expense of a greater computing cost.

---

[2]As described here, the goal of the fast gradient sign method is to cause a generic model to make a mistake. The method has also been extended to make classifiers misclassify the input as belonging to a specific class by differentiating the probability of the desired class rather than differentiating the cost function.

# 5 Validation of the Attack

To test the effectiveness of our attack, we first apply it to target the oracle provided by MetaMind, a startup serving an API[3] that allows one to build classifiers using deep learning. The API returns labels produced by a DNN for any given input but does not provide access to the DNN. This fits the characteristics of the oracle DNN described in our threat model. We show that an adversary implementing our attack can reliably force an oracle trained using MetaMind on MNIST [25] to misclassify 84.24% of samples altered with a perturbation not affecting human recognition performance. We then target a different oracle trained on the German Traffic Signs Recognition Benchmark (GTSRB) [37], and show that we can force it to misclassify more than 64.24% of its inputs altered with a perturbation not affecting human recognition.

## 5.1 Attack against the MetaMind Oracle

### 5.1.1 Description of the Oracle

The dataset used to train the oracle is the MNIST handwritten digit dataset [25]. It is made of $60,000$ training samples, $10,000$ validation samples, and $10,000$ test samples. Each sample is an image of a handwritten digit. The task associated with this dataset is to train a classifier to identify the digit corresponding to each image. Each 28x28 grayscale sample is encoded as a vector of pixel intensities normalized in the interval $[0,1]$ and obtained by reading the image pixel matrix row-wise.

We registered for an API key on MetaMind's website, which gave us access to three functionalities: dataset upload, automated model training, and model prediction querying. We uploaded the $50,000$ samples included in the MNIST training set to MetaMind's servers and then used the API to train a classifier on the dataset. Again, we emphasize that training is automated: we have no access to the training algorithm, model architecture, or model parameters. All we are given is the accuracy of the resulting model, computed by MetaMind using a validation set created by isolating 10% of the samples we provided for training. Details can be found on MetaMind's website.

Training took approximatively 36 hours and returned a classifier yielding a 94.97% accuracy. This performance cannot be improved as we cannot access or modify the model's specifications and training algorithm. Once training is completed, we access the model predictions, for any input of our choice, through the API. Predictions take the form of a label indicating the class assigned the highest probability. This corresponds to adversarial capabilities of the threat model described in Section 3.

Now that the oracle is set up, we take the adversarial perspective and follow the attack described in Section 4.

---

[3]The MetaMind API can be accessed online at www.metamind.io

### 5.1.2 Initial Substitute Training Sets

First, the adversary collects an initial substitute training set. We describe two such sets used to attack the Meta-Mind oracle: a subset of MNIST and a handcrafted set.

**From the MNIST test set -** Our first initial substitute training set is made of 150 samples extracted from the MNIST test set. We use 100 samples for training and 50 for validation. Note that these samples are different from those used by the oracle for training as the test and training sets are distinct. We assume adversaries can collect such a limited amount of samples under the threat model described in Section 3 using minimal knowledge of the oracle task: in our case, handwritten digit classification.

**Handcrafted set -** To ensure our results do not stem from similarities between the MNIST test and training sets, we also consider a *handcrafted* initial substitute training set. We handcrafted 100 samples by handwriting 10 digits for each class between 0 and 9 with a laptop trackpad. We then adapt them to the MNIST format of 28x28 grayscale pixels. Handcrafted samples are shown below.



### 5.1.3 Substitute DNN Training

The adversary then uses the initial substitute training sets and the oracle to train two subsitute DNNs. Our substitute architecture A is described in Table 3 (cf. appendix). It is a standard architecture for computer vision classification. We train all DNNs in this paper with Theano [5, 7] and Lasagne [12]. The DNN is trained on our machine for 6 substitute epochs. During each of these 6 epochs, the DNN is trained for 10 epochs from scratch with a learning rate of $10^{-2}$ and a momentum of 0.9. Between substitute epochs, we perform a Jacobian-based dataset augmentation with a step of $\lambda_\rho = 0.1$.

The accuracy of the 2 substitute DNNs is reported in Table 1. It is computed with the MNIST test set (minus the 150 samples used in the first initial substitute training set). The adversary does *not* have access to this full test set: we solely use it to analyze our results. The two substitute DNNs respectively achieve a 81.20% and 67.00% accuracy on the MNIST test set after 6 substitute training epochs. These accuracies fall short of current state-of-the-art accuracies on this particular task. However one should keep in mind that the adversary has access to a limited number of samples (in this case $6,400 = 100 \times 2^6$ instead of $50,000$ for state-of-the-art models). Furthermore, the adversarial goal is to craft adversarial samples misclassified by the oracle. *Instead of training a substitute DNN matching the oracle accuracy, the adversary is interested in training a substitute DNN capable of mimicking the oracle's decision boundaries.*

| Substitute Epoch | Initial Substitute Training Set from | |
|---|---|---|
| | MNIST test set | Handcrafted digits |
| 0 | 24.86% | 18.70% |
| 1 | 41.37% | 19.89% |
| 2 | 65.38% | 29.79% |
| 3 | 74.86% | 36.87% |
| 4 | 80.36% | 40.64% |
| 5 | 79.18% | 56.95% |
| 6 | 81.20% | 67.00% |

Table 1: **Substitute DNNs Accuracies:** the accuracy is evaluated for two DNNs corresponding to each initial substitute training set: 150 examples from the MNIST test set, and handcrafted digits. The accuracy is reported on the remaining 9,850 test set examples.
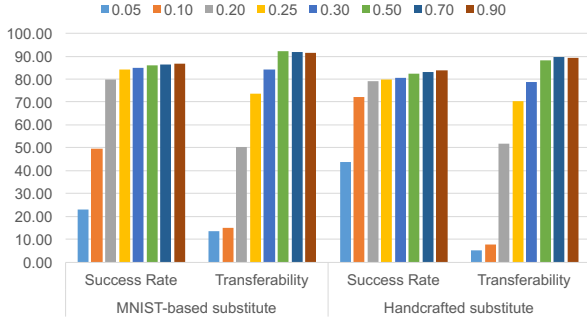


Figure 4: **Success Rate and Transferability of Adversarial Samples for the MetaMind attacks:** performed using the MNIST-based and the handcrafted substitutes.

### 5.1.4 Adversarial Sample Crafting

Using the substitute DNNs, the adversary then moves on to the next fold of our attack: crafting adversarial samples using Goodfellow's algorithm (cf. Section 4). We decided to use the 10,000 samples from the MNIST test set as our legitimate samples.[4] We evaluate sample crafting using two metrics: *success rate* and *transferability*. The *success rate* is the proportion of adversarial samples misclassified by the substitute DNN. Our goal is to verify whether these samples are also misclassified by the oracle or not. Therefore, the *transferability of adversarial samples* refers to the oracle misclassification rate of adversarial samples crafted using the substitute DNN.

Figure 4 details both metrics for each substitute DNN and for several values of the input variation $\varepsilon$ (cf. Equation 6). Transferability reaches 84.24% for the first substitute DNN and 78.72% for the second, with input variations of $\varepsilon = 0.3$. Our attack strategy is thus effec-

---

[4]Again, note that adversaries do not need access to the dataset and can use any legitimate sample of their choice to craft adversarial samples. We use the dataset in order to show that expected inputs can be misclassified on a large scale.
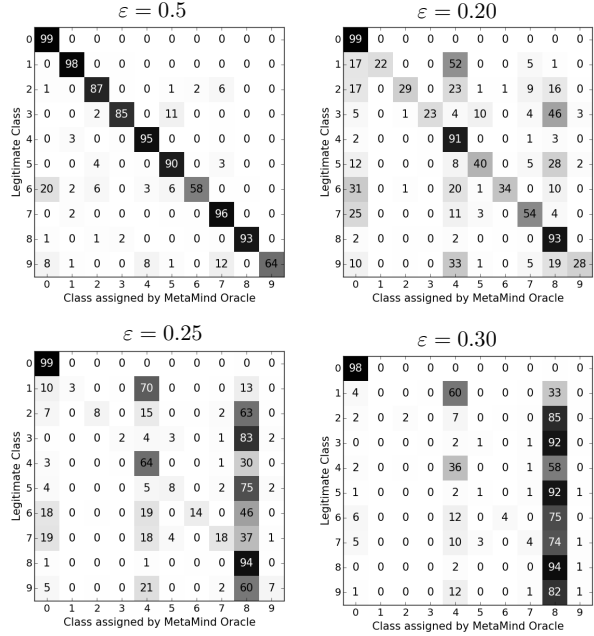


Figure 5: **MetaMind Oracle Confusion Matrices:** for 4 input variations $\varepsilon$. Cell $(x, y)$ indicates the percentage of instances of digit $y$ classified by the oracle as digit $x$.

tively able to severely damage the output integrity of the MetaMind oracle. Using the substitute training set handcrafted by the adversary limits the transferability of adversarial samples when compared to the substitute set extracted from MNIST data, for all input variations except $\varepsilon = 0.2$. Yet, the transferability of both substitutes is similar, corroborating that our attack can be executed without access to any of the oracle's training data.

To analyze the labels assigned by the MetaMind oracle, we plot confusion matrices for adversarial samples crafted using the first substitute DNN and with 4 values of input variation $\varepsilon$. In Figure 5, rates on the diagonal indicate the proportion of samples correctly classified by the oracle for each of the 10 classes. Off-diagonal rates indicate the proportion of samples misclassified in a wrong class. For instance, cell $(8, 3)$ in the third matrix indicates that 89% instances of a 3 are classified as a 8 by the oracle when samples are perturbed with an input variation of $\varepsilon = 0.25$. Interestingly, confusion matrices converge to most samples being classified as 4s and 8s as $\varepsilon$ increases. As observed in previous work on adversarial crafting [30], an explanation to this phenomenon could be that DNNs more easily classify adversarial samples in these classes. Additionally, it is possible that MetaMind augments the DNN with an ensemble of classifiers collectively assigning labels to samples. Adversarial samples might display properties characteristics of an 8 like for instance a significant proportion of white pixels.

## 5.2 Attacking an oracle for the GTSRB

We now validate our attack on a different dataset, using an oracle trained with our machine to recognize traffic signs on the GTSRB dataset. The attack achieves *higher transferability rates at lower distortions compared to the MNIST oracle*. We believe this is due to the higher dataset dimensionality, in terms of inputs and outputs.

**Description of the Oracle -** The GTSRB dataset is an image collection of 43 types of traffic signs [37]. Images vary in size and are RGB-encoded. To simplify our classifier, we resize images to 32x32 pixels, recenter them by subtracting the mean component value, and rescale them by factoring their standard deviations out. We keep $35,000$ images for our training set and $4,000$ for our validation set (out of the $39,209$ training set), and $10,000$ for our test set (out of $12,630$). We train the oracle on our machine, using architecture B from Table 3 (cf. appendix), for 50 epochs with a learning parameter of $10^{-2}$ and a momentum of 0.9 (both decayed by 0.5 every 10 epochs). Data is processed in batches of 100 samples.

**Substitute DNN Training -** The adversary uses two initial substitute training sets extracted from the GTSRB test set. The first one includes the first $1,000$ samples and the second one the first 500. The number of initial samples is higher than for MNIST substitutes as the inputs have a higher dimensionality. We train three substitute DNN architectures C, D, and E (cf. Table 3) using the oracle for 6 substitute training epochs with a Jacobian-based dataset augmentation parameter of $\lambda_\rho = 0.1$. Substitute C and E where trained with the $1,000$ sample initial substitute training set and achieve a 71.42% accuracy. Substitute D was trained with the initial set of 500 samples. Its accuracy of 60.12% is lower than C and E.

**Sample Crafting -** We use Goodfellow's algorithm with input variations $\varepsilon$ between 0.01 and 0.5 to craft adversarial samples from the test set. Results for all substitute models are shown in Figure 6. Adversarial samples crafted with input varitions $\varepsilon < 0.3$ are more transferable than those crafted with the same $\varepsilon$ for MNIST models. This is likely due to the higher input dimensionality— $3,072$ components instead of 784—which means almost 4 times more perturbation is applied with the same $\varepsilon$. Nevertheless, with success rates higher than 98.98% and transferability rates ranging from 64.24% to 69.03% for $\varepsilon = 0.3$, which is hard to distinguish for humans, *the attack is successful*. Interestingly, the transferability of adversarial samples crafted using substitute DNN D is comparable or higher than corresponding samples for DNNs C and E, even though it was trained with less samples and is less accurate. This emphasizes our intuition that adversaries should not primarily be focused with improving the accuracy of substitute DNNs.
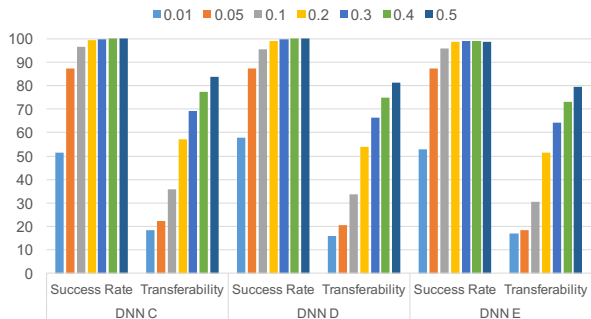


Figure 6: **Success Rate and Transferability of Adversarial Samples crafted on the GTRSRB dataset**

## 6 Attack Algorithm Calibration

Having shown in Section 5 that an adversary can force an MNIST oracle from MetaMind, as well as a GTSRB oracle trained locally, to misclassify inputs using the attack described in Section 4, we now perform an exploration of the parameter space of both attack steps–the substitute DNN training and the adversarial sample crafting–using the MNIST dataset. We explore the following questions: "(1) How can the substitute DNN training be fine-tuned to improve adversarial sample transferability?" and (2) "For each adversarial sample crafting strategies, which parameters optimize transferability?". Our findings are:

- The choice of substitute DNN architecture (number of layers, size, activation function, type) has a limited impact on adversarial sample transferability. Increasing the number of epochs, after the substitute DNN has reached an asymptotic accuracy, does not improve adversarial sample transferability.

- At comparable perturbation distortions, the Goodfellow and Papernot algorithms have similar transferability rate. The Goodfellow algorithm is however more computationally efficient.

In this section, we use a local oracle to limit querying of MetaMind. We train architecture A (cf. Table 3) for 50 epochs with a learning parameter $10^{-2}$ and a momentum 0.9 (both decayed by 0.5 every 10 epochs). Our oracle yields an accuracy of 99.50% on the MNIST test set.

## 6.1 Calibrating Substitute DNN Training

We first seek to quantify the impact of the substitute DNN training on adversarial sample transferability.

**Choosing an Architecture -** We train substitute architectures A, F, G, H, I, J, K, L, and M (cf. Table 3) using the substitute training set from Section 5, which holds 150 samples from the MNIST test set. During each of the 6 substitute training epochs, the network is trained for 5

| DNN ID | Accuracy ($\rho = 2$) | Accuracy ($\rho = 6$) | Transferability ($\rho = 6$) |
|--------|--------|--------|--------|
| A | 30.50% | 82.81% | 75.74% |
| F | 68.67% | 79.19% | 64.28% |
| G | 72.88% | 78.31% | 61.17% |
| H | 56.70% | 74.67% | 63.44% |
| I | 57.68% | 71.25% | 43.48% |
| J | 64.39% | 68.99% | 47.03% |
| K | 58.53% | 70.75% | 54.45% |
| L | 67.73% | 75.43% | 65.95% |
| M | 62.64% | 76.04 | 62.00% |

Table 2: **Substitute Accuracy** at $\rho = 2$ and $\rho = 6$ substitute epochs and **Transferability of Adversarial Samples:** crafted with $\varepsilon = 0.4$ after $\rho = 6$ substitute epochs.



Figure 7: **Impact of the input variation $\varepsilon$ in the Goodfellow crafting algorithm on the transferability of adversarial samples:** for architectures from Table 2.

epochs from scratch. Between epochs, the dataset is augmented using Jacobian-based dataset augmentations with step $\lambda_\rho = 0.1$. The substitute architectures differ from the oracle architecture by the type, number, and size of layers used. In Table 2, we report the accuracy of each architecture after 2 and 6 substitute training epochs, evaluated on the MNIST test set. After 6 epochs, all architectures reach comparable accuracies. Beyond 6 epochs, we are not able to improve the accuracies, suggesting that our dataset augmentation technique needs to be refined, or additional samples need to be added to the initial substitute training set. However, we are primarily interested in adversarial sample transferability.

Adversarial samples are crafted using the Goodfellow algorithm with an input variation of $\varepsilon = 0.4$ (which we justify later). The last column of Table 2 reports the transferability of adversarial samples for the substitute DNNs measured on the MNIST test set. These results show that the architecture choice has a limited impact on adversarial sample transferability, and therefore on the attack success. The most important transferability drop follows from removing all convolutional layers. Changing the hidden layer activation function from rectified linear to a sigmoid does not impact transferability significantly. Finally, there seem to be no apparent correlation between the accuracy of the substitute DNN and the transferability of adversarial samples.

**Choosing the number of substitute epochs -** An additional vector of improvement for substitute DNNs is the number of substitute training epochs for which it they are trained. Intuitively, one would hypothesize that the longer we train the substitute DNN, the more samples labeled using the oracle are included in the substitute training set, thus the higher the transferability of adversarial samples will be. This intuition is only partially confirmed by our experiments on substitute DNN A. We find that for for input variations $\varepsilon \leq 0.3$, the transferabil-
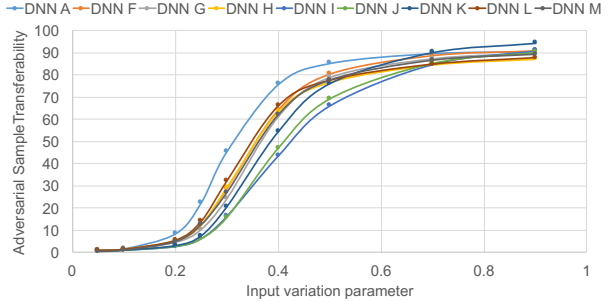
ity is slightly improved by a rate between $+3\%$ to $+9\%$, but for variations $\varepsilon \geq 0.4$, the transferability is slightly degraded by less than 1%.

**Setting the dataset augmentation step size -** We trained substitute A using different step size $\lambda_\rho$ for the Jacobian-based dataset augmentation. Increasing the step size from 0.1 (used in the rest of this paper) to 0.2 or 0.3 does not modify the substitute DNN's accuracy by more than 3% but leads to a less stable learning convergence. Reducing the step from 0.1 to 0.01 or 0.05 did not change the accuracy by more than 3% but made the learning convergence slower. Generally speaking, increasing step size $\lambda_\rho$ negatively impacts adversarial sample transferability : for instance with a step size of 0.3 compared to 0.1, the transferability rate for $\varepsilon = 0.25$ is 10.82% instead of 22.35% and for $\varepsilon = 0.5$, 82.07% instead of 85.22%.

## 6.2 Adversarial Sample Crafting

We now compare the transferability of adversarial samples produced by each of the two crafting algorithms introduced previously [15, 30]. We first calibrate both algorithms to improve the transferability of adversarial samples produced. We then compare the results to elect the strongest technique under our threat model.

**Goodfellow's algorithm -** Recall from earlier that perturbations for the Goodfellow attack are evaluated using:

$$\delta_{\vec{x}} = \varepsilon \, \mathrm{sgn}(\nabla_{\vec{x}} cost(F, \vec{x}, y)) \qquad (8)$$

Thus, the only parameter to optimize for this crafting algorithm is $\varepsilon$: the variation added in the direction indicated by the sign of the cost function's gradient for each input component. We use the same architecture collection as before to quantify the impact of $\varepsilon$ on adversarial sample transferability. As shown in Figure 7, architecture A outperforms all others because it is a copy of the oracle and acts as a baseline. Other architectures have
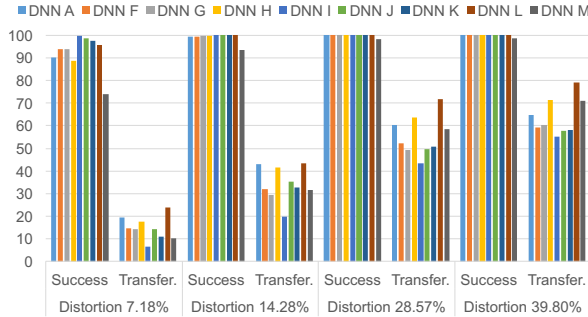
Figure 8: **Impact of the maximum distortion factor $\Upsilon$ in the Papernot crafting strategy on the success rate and transferability of adversarial samples:** although not required to achieve a close to a 100% success rate at adversarial sample crafting, increasing the distortion leads to higher transferability rates across all DNNs.



Figure 9: **Impact of the input variation $\varepsilon$ in the Papernot algorithm on the success rate and adversarial sample transferability** computed for $\varepsilon \in \{0.5, 0.7, 1\}$ on DNNs from Table 2 with a distortion of $\Upsilon = 39.80\%$.

asymptotic transferability rates ranging between 72.24% and 80.21%, confirming that *the substitute architecture choice has a limited impact on transferability*. Increasing the value of $\varepsilon$ above 0.4 yields little improvement in transferability and should be avoided so as to guarantee indistinguishability of adversarial samples to humans. Interestingly, the curves can be partitioned in two groups: one corresponding to architectures using convolutional layers (A, F, G, H, L, M) and a second for architectures only using fully connected layers (I, J, K).

**Papernot's algorithm -** This algorithm is fine-tuned by two parameters: the *maximum distortion* $\Upsilon$ and the *input variation* $\varepsilon$. The maximum distortion[5] defines the number of input components that are selected to construct perturbation $\delta_{\vec{x}}$. The input variation, similarly to the Goodfellow algorithm, controls the amount of change induced to input components making up perturbation $\delta_{\vec{x}}$.

We first evaluate the impact of the maximum distortion $\Upsilon$ on adversarial sample transferability. For now, components selected to be perturbed are increased by $\varepsilon = 1$. Intuitively, one expects that increasing the maximum distortion will make adversarial samples more transferable. Indeed, even though some adversarial samples would be misclassified by the substitute DNN with lower distortions, higher distortions increase the confidence of the substitute DNN making a misclassification, and also increases the likelihood of the oracle also misclassifying the sample. In Figure 8, we confirm this intuition with different values of the maximum distortion $\Upsilon$. Results show that increasing distortion $\Upsilon$ from 7.14% to 28.57% improves transferability: at a 7.14% distortion, the aver-

age transferability with the oracle across all architectures is 14.70% whereas at a 28.57% is averages at 55.53%.

We run a second experiment to quantify the impact of the input variation $\varepsilon$ introduced to each component selected to be part of $\delta_{\vec{x}}$. We find that reducing the input variation from $\varepsilon = 1$ to $\varepsilon = 0.7$ significantly degrades the transferability of adversarial samples, approximatively by a factor of 2, as reported in Figure 9. This is due to the fact that we kept the distortion parameter $\Upsilon$ fixed, thus the crafting algorithm is not able to compensate the lost of effectiveness induced by the decrease of $\varepsilon$ by increasing the number of components altered.

**Comparing Crafting Algorithm -** To compare both strategies, we compute the L1 norm of perturbation $\delta_{\vec{x}}$ introduced by each algorithm, defined as:

$$\|\delta_{\vec{x}}\|_1 = \varepsilon \cdot \|\delta_{\vec{x}}\|_0 \qquad (9)$$

where $\|\cdot\|_0$ is used to denote the number of non-zero entries[6] in a vector, and $\varepsilon$ is the input variation—both crafting algorithms introduce a constant variation to components perturbed. Note that for the Goodfellow algorithm, we always have $\|\delta_{\vec{x}}\|_0 = 1$, whereas for the Papernot algorithm, values vary for both factors $\varepsilon$ and $\|\delta_{\vec{x}}\|_0$. For instance, the case $\|\delta_{\vec{x}}\|_1 = 0.4$ corresponds to a Goodfellow algorithm with $\varepsilon = 0.4$ and a Papernot algorithm with $\varepsilon = 1$ and distortion $\Upsilon = 40\%$. Corresponding transferability rates can be found respectively in Table 2 and Figure 8 for our running set of architectures. Performances are comparable with some DNNs performing better with the Goodfellow algorithm and others with the Papernot algorithm. Thus, the choice of algorithm depends on the type of perturbation acceptable: e.g., all features perturbed a little vs. some features perturbed a lot. Indeed, the Goodfellow algorithm gives more control on $\varepsilon$ while the Papernot algorithm gives more control on $\Upsilon$.

---

[5]Note that in the original algorithm introduced by Papernot et al. [30], the algorithm stopped perturbing the input if it reached a target class different from the source class. Here, we force the algorithm to continue perturbing the input until it has perturbed $\Upsilon$ input components.
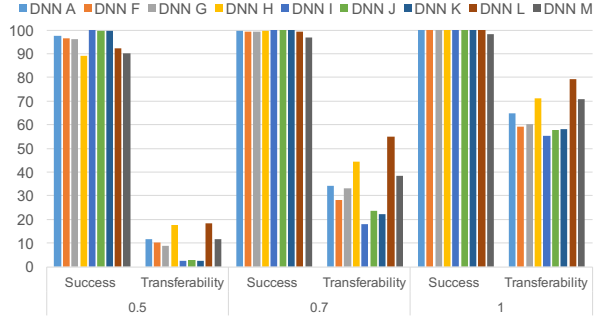
[6]It is not strictly a norm in the mathematical sense because it lacks the homogeneity and triangle inequality property

# 7 Intuition Behind Transferability

Previous work has started explaining why adversarial samples generalize and transfer between different architectures [15, 38]. Here, we build an intuition behind transferability based on *statistical hypothesis testing* [26] and an analysis of DNN cost gradient sign matrices. A formal treatment of transferability is left as future work.

Recall the perturbation added to craft adversarial samples in the Goodfellow algorithm. Inspecting Equation 6, it is clear that, given a sample $\vec{x}$, the noise added would be the same for two DNNs $F$ and $G$ if $\text{sgn}(\nabla_{\vec{x}}cost(F,\vec{x},y))$ and $\text{sgn}(\nabla_{\vec{x}}cost(G,\vec{x},y))$ were equal. Let us write the space of these matrices, which have entries in $\{+1,-1\}$, as $\text{Sgn}_{n \times m}$. Assume that samples $\vec{x}$ are generated from a population distribution $\mathscr{D}$ (e.g., in our case the distribution from which the images of digits are drawn). The formula $\text{sgn}(\nabla_{\vec{x}}cost(F,\vec{x},y))$ and $\mathscr{D}$ induce a distribution $\mathscr{D}_F$ over $\text{Sgn}_{n \times m}$ (i.e. randomly draw a sample from the distribution $\mathscr{D}$ and compute the quantity in Eq. 6). Similarly, DNN $G$ and distribution $\mathscr{D}$ induce a distribution $\mathscr{D}_G$ over matrices in $\text{Sgn}_{n \times m}$. Our main conjecture is:

> For two "similar" architectures $F$ and $G$ distributions $\mathscr{D}_F$ and $\mathscr{D}_G$ induced by a population distribution $\mathscr{D}$ are highly correlated.

Note that if the distributions $\mathscr{D}_F$ and $\mathscr{D}_G$ were independent, then the noise they add during adversarial sample crafting is independent. In this case, our intuition tells us that adversarial samples would not transfer (after all, in the two cases you are adding noise that are independent). The question is: how to verify our conjecture? The challenge is that the population distribution $\mathscr{D}$ is unknown.

We turn to statistical hypothesis testing. Note that we can empirically estimate the distributions $\mathscr{D}_F$ and $\mathscr{D}_G$ based on known samples. First, we generate two sequences of sign matrices $\sigma_1 = \langle M_1, M_2, \cdots \rangle$ and $\sigma_2 = \langle N_1, N_2, \cdots \rangle$ using the sample set (e.g. MNIST in our case) for a substitute DNN $F$ and oracle $G$. Next we pose the following *null hypothesis*:

> $H_N$: The sequences $\sigma_1$ and $\sigma_2$ are drawn from independent distributions.

We use standard tests from the statistical hypothesis testing literature to test the hypothesis $H_N$. If the hypothesis $H_N$ is *rejected*, then we know that the sign matrices corresponding to the two architectures $F$ and $G$ are correlated.

Next, we describe the test we use. There are several algorithms for hypothesis testing, and we picked a simple one based on chi-square test. A thorough investigation of other hypothesis-testing techniques, is left as future work. Let $p_{i,j}$ and $q_{i,j}$ be the frequency of $+1$ in the $(i,j)$-th entry of the matrices in the sequences $\sigma_1$ and $\sigma_2$, respectively. Let $r_{i,j}$ be the frequency of the $(i,j)$-th entry
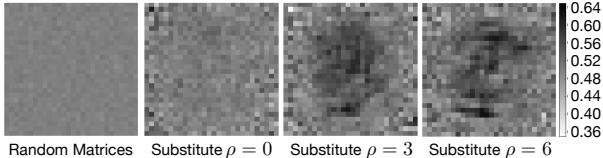


Figure 10: **Frequencies of cost gradient sign matrix components equal between substitute DNN A and the oracle** at three substitute training epochs $\rho \in \{0, 3, 6\}$, compared to a pair of random sign matrices (on the left).

being $+1$ in both sequences $\sigma_1$ and $\sigma_2$ simultaneosuly.[7] Note that if the distributions were independent then $r_{i,j} = p_{i,j}q_{i,j}$. However, if the distributions are correlated, then we expect $r_{i,j} \neq p_{i,j}q_{i,j}$. Consider quantity:

$$\chi^{2\star} = \sum_{i=1}^{m} \sum_{j=1}^{n} \frac{(r_{i,j}N - p_{i,j}q_{i,j}N)^2}{p_{i,j}q_{i,j}N}$$

where $N$ is the number of samples. In the $\chi$-square test, we compute the probability that $P(\chi^2 > \chi^{2\star})$, where $\chi^2$ has degrees of freedom $(m-1)(n-1) = 27 \times 27 = 729$ for the MNIST data. The $\chi^{2\star}$ scores for substitute DNNs from Table 2 range between $61,403$ for DNN A and $88,813$ for DNN G. Corresponding P-values are below $10^{-5}$ for all architectures, with confidence $p < 0.01$. Thus, for all substitute DNNs, the hypothesis $H_N$ is largely rejected: sequences $\sigma_1$ ans $\sigma_2$, and therefore sign matrices corresponding to pairs of a substitute DNN and the oracle, are highly correlated. As a baseline comparison, we generate 2 random sign matrices and compute the corresponding $\chi^{2\star}$ score: 596. We find a P-Value of 0.99 with a confidence of 0.01, meaning that these matrices were indeed drawn from independent distribution in contrast to previous results.

However, we must now complete our analysis to characterize the correlation suggested by the hypothesis testing. In Figure 10, we plot the frequency matrix $R = [r_{i,j}]$ for several pairs of matrices. The first on the left is a pair of random matrices of $\{+1, -1\}$. The other matrices correspond to the oracle and substitute DNN A at different substitute training epochs $\rho$. Frequencies are computed using the $10,000$ samples of the MNIST test set. Although all frequencies in the random pairs are very close to $1/2$, frequencies corresponding to pixels located in the center of the image are higher in the $(substitute, oracle)$ matrix pairs. Furthermore, the phenomenon amplifies as training progresses through the substitute epochs. We then compute the frequencies separately for each sample source class in Figure 11. For each class, it is clear that the sign matrices agree on pixels relevant for classi-

---

[7] We assume that the frequencies are normalized so they can be interpreted as probabilities, and also assume that all frequencies are $> 0$ to avoid division by zero, which can be achieved by rescaling.
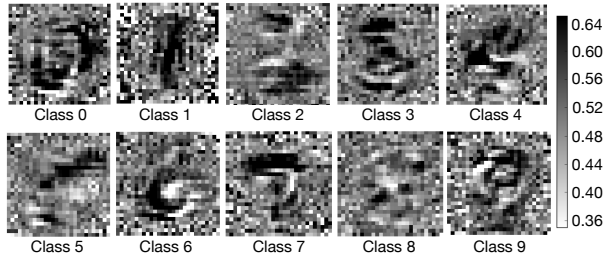
Figure 11: **Frequencies of cost gradient sign matrix components equal between substitute DNN A and the oracle** computed separately by sample source class.

fication in that class. Similar results for other substitute DNNs are summarized in Figure 12 of the appendix. One observation is that substitutes that yielded smaller transferability rates in our evaluation also have less components of their cost gradient sign matrix frequently equal to the oracle's. This suggests that *correlations between the respective sign matrices of the substitute DNN and of the oracle—for input components that are relevant to classification in each respective class—could explain cross-model adversarial sample transferability.*

## 8 Discussion

This paper contributes to a line of works on machine learning security [3, 4, 20, 30, 31]. Attacks were proposed for classifiers built using machine learning techniques distinct from DNNs like tree ensembles [21], random forests [40], and SVMs [40]. It has been demonstrated that attacks can be executed at training time [9] and test time [10]. A framework for securing binary classifiers (e.g., logistic regression, SVM) is described in [8].

Our black-box DNN threat model is challenging because it considers adversaries without knowledge of the oracle architecture and training set. Knowledge of either one of these reduces our threat model to previous work. Finally, we assumed attackers only had access to the label $\tilde{O}$ assigned by the oracle. Attackers with access to the probability vector $O$ could more accurately extract the oracle model's knowledge [2, 19]. However, this attack is outside of the scope of this paper.

In some cases, it may be difficult to collect the initial substitute training set. It requires knowledge of the task achieved by the oracle, like the input and output types. Encouragingly, we saw in Section 5 that our attack is successful even with samples not extracted from the oracle training distribution. Furthermore, targeted oracles were successfully attacked with 100 samples on MNIST and 500 samples on GTSRB. A second crucial point is the number of queries: our dataset augmentation limits oracle querying, thus making the attack tractable, but it may remind too high to go undetected in some environments. Both of these points leave opportunity for future work.

Discussing defenses against this attack is outside the scope of this paper. Previous work suggested the use of adversarial sample training [15], Jacobian-based regularization [16], and distillation [31] as means to make DNNs robust to adversarial samples. Such solutions apply to our settings. Additionally, a careful analysis of oracle queries (taking into account distributed queries) could help thwart attempts at substitute DNN training.

Finally, we only discussed in this paper oracles built using DNNs due to space constraints. However, we emphasize here that no part of the attack design constraints it to DNN-based oracles. In fact, we conducted preliminary experiments with a *k*-nearest neighbor (k-NN) oracle making predictions on unseen inputs by assigning the label of the closest[8] sample from the MNIST training set. Following the method described in Section 4, we trained a substitute DNN using this k-NN oracle and an initial substitute training set of 150 samples from the MNIST test set. The substitute DNN yielded an accuracy of 80.64% on the remaining 9,850 samples of the test set after 6 substitute epochs. We then used this substitute DNN to craft 1,000 adversarial samples. When crafted with an input variation of $\varepsilon = 0.25$ (respectively $\varepsilon = 0.5$), adversarial samples are misclassified at a rate of 27.3% (respectively 56.7%) by the k-NN oracle. This suggests that our attack is applicable to oracles using machine learning algorithms fully different from deep learning. An extended investigation is left as future work.

## 9 Conclusions

We introduced and evaluated a novel *black-box* attack against machine learning oracles. Our work constitutes a significant step towards relaxing strong assumptions, regarding adversarial capabilities, made by previous work. We demonstrated that practical *black-box* attacks can be conducted in properly blinded settings against remotely hosted oracles built using deep learning, a state-of-the-art machine learning technique. Indeed, we successfully targeted one of the online APIs providing deep learning as a service, namely MetaMind, and forced it to misclassify 84.24% of our adversarial samples. The proposed attack introduces a substitute DNN training method to craft adversarial samples misclassified by oracles with no access to both their architectures and training set. Finally, we empirically built an intuition for the transferability of adversarial samples across architectures.

Future work should deploy this attack against oracles built with machine learning algorithms distinct from deep learning to formalize the initial results mentioned in Section 8. Defenses, outside of the scope of this paper, should also be designed to mitigate this novel threat vector against machine learning oracles.

---

[8]Distances between samples are evaluated with the Frobenius norm.

# References

[1] ALIPANAHI, B., DELONG, A., WEIRAUCH, M. T., AND FREY, B. J. Predicting the sequence specificities of DNA-and RNA-binding proteins by deep learning. *Nature biotechnology* (2015).

[2] BA, J., AND CARUANA, R. Do deep nets really need to be deep? In *Advances in Neural Information Processing Systems* (2014), pp. 2654–2662.

[3] BARRENO, M., NELSON, B., JOSEPH, A. D., AND TYGAR, J. The security of machine learning. *Machine Learning 81*, 2 (2010), 121–148.

[4] BARRENO, M., NELSON, B., SEARS, R., JOSEPH, A. D., AND TYGAR, J. D. Can machine learning be secure? In *Proceedings of the 2006 ACM Symposium on Information, computer and communications security* (2006), ACM, pp. 16–25.

[5] BASTIEN, F., LAMBLIN, P., PASCANU, R., BERGSTRA, J., GOODFELLOW, I., BERGERON, A., BOUCHARD, N., WARDE-FARLEY, D., AND BENGIO, Y. Theano: new features and speed improvements. *arXiv preprint arXiv:1211.5590* (2012).

[6] BENGIO, Y. Learning deep architectures for AI. *Foundations and trends® in Machine Learning 2*, 1 (2009), 1–127.

[7] BERGSTRA, J., BREULEUX, O., BASTIEN, F., LAMBLIN, P., PASCANU, R., DESJARDINS, G., TURIAN, J., WARDE-FARLEY, D., AND BENGIO, Y. Theano: a CPU and GPU math expression compiler. In *Proceedings of the Python for scientific computing conference (SciPy)* (2010), vol. 4, Austin, TX, p. 3.

[8] BIGGIO, B., FUMERA, G., AND ROLI, F. Security evaluation of pattern classifiers under attack. *IEEE Transactions on Knowledge and Data Engineering 26*, 4 (2014), 984–996.

[9] BIGGIO, B., NELSON, B., AND LASKOV, P. Support vector machines under adversarial label noise. In *ACML* (2011), pp. 97–112.

[10] BIGGIO, B., NELSON, B., AND PAVEL, L. Poisoning attacks against support vector machines. In *Proceedings of the 29th International Conference on Machine Learning* (2012).

[11] DAHL, G. E., STOKES, J. W., DENG, L., AND YU, D. Large-scale malware classification using random projections and neural networks. In *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (2013), IEEE, pp. 3422–3426.

[12] DIELEMAN, S., SCHLTER, J., RAFFEL, C., OLSON, E., SNDERBY, S. K., ET AL. Lasagne: First release., Aug. 2015.

[13] FUKUSHIMA, K., AND MIYAKE, S. Neocognitron: A new algorithm for pattern recognition tolerant of deformations and shifts in position. *Pattern recognition 15*, 6 (1982), 455–469.

[14] GOODFELLOW, I., BENGIO, Y., AND COURVILLE, A. Deep learning. Book in preparation for MIT Press, 2016.

[15] GOODFELLOW, I. J., SHLENS, J., AND SZEGEDY, C. Explaining and harnessing adversarial examples. In *Proceedings of the 2015 International Conference on Learning Representations* (2015), Computational and Biological Learning Society.

[16] GU, S., AND RIGAZIO, L. Towards deep neural network architectures robust to adversarial examples. In *Proceedings of the 2015 International Conference on Learning Representations* (2015), Computational and Biological Learning Society.

[17] HE, K., ZHANG, X., REN, S., AND SUN, J. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *IEEE International Conference on Computer Vision* (2015).

[18] HINTON, G., DENG, L., YU, D., DAHL, G. E., MOHAMED, A.-R., JAITLY, N., SENIOR, A., VANHOUCKE, V., NGUYEN, P., SAINATH, T. N., ET AL. Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups. *IEEE Signal Processing Magazine 29*, 6 (2012), 82–97.

[19] HINTON, G., VINYALS, O., AND DEAN, J. Distilling the knowledge in a neural network. In *Deep Learning and Representation Learning Workshop at NIPS 2014* (2014), arXiv preprint arXiv:1503.02531.

[20] HUANG, L., JOSEPH, A. D., NELSON, B., RUBINSTEIN, B. I., AND TYGAR, J. Adversarial machine learning. In *Proceedings of the 4th ACM workshop on Security and artificial intelligence* (2011), ACM, pp. 43–58.

[21] KANTCHELIAN, A., TYGAR, J., AND JOSEPH, A. D. Evasion and hardening of tree ensemble classifiers. *arXiv preprint arXiv:1509.07892* (2015).

[22] KNORR, E. How paypal beats the bad guys with machine learning, 2015.

[23] KRIZHEVSKY, A., SUTSKEVER, I., AND HINTON, G. E. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems* (2012), pp. 1097–1105.

[24] LECUN, Y., BOTTOU, L., BENGIO, Y., AND HAFFNER, P. Gradient-based learning applied to document recognition. *Proceedings of the IEEE 86*, 11 (1998), 2278–2324.

[25] LECUN, Y., AND CORTES, C. The mnist database of handwritten digits, 1998.

[26] LEHMANN, E. L., AND ROMANO, J. P. *Testing Statistical Hypotheses*. Springer Texts in Statistics, August 2008.

[27] MNIH, V., KAVUKCUOGLU, K., SILVER, D., GRAVES, A., ANTONOGLOU, I., WIERSTRA, D., AND RIEDMILLER, M. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602* (2013).

[28] NAIR, V., AND HINTON, G. E. Rectified linear units improve restricted boltzmann machines. In *Proceedings of the 27th International Conference on Machine Learning (ICML-10)* (2010), pp. 807–814.

[29] NVIDIA. Nvidia tegra drive px: Self-driving car computer, 2015.

[30] PAPERNOT, N., MCDANIEL, P., JHA, S., FREDRIKSON, M., CELIK, Z. B., AND SWAMI, A. The limitations of deep learning in adversarial settings. In *Proceedings of the 1st IEEE European Symposium on Security and Privacy* (2016), arXiv preprint arXiv:1511.07528.

[31] PAPERNOT, N., MCDANIEL, P., WU, X., JHA, S., AND SWAMI, A. Distillation as a defense to adversarial perturbations against deep neural networks. *arXiv preprint arXiv:1511.04508* (2015).

[32] PENNINGTON, J., SOCHER, R., AND MANNING, C. D. Glove: Global vectors for word representation. *Proceedings of the Empirical Methods in Natural Language Processing (EMNLP 2014) 12* (2014), 1532–1543.

[33] PETERSON, C., AND SÖDERBERG, B. A new method for mapping optimization problems onto neural networks. *International Journal of Neural Systems 1*, 01 (1989), 3–22.

[34] RANZATO, M. A., HUANG, F. J., BOUREAU, Y.-L., AND LECUN, Y. Unsupervised learning of invariant feature hierarchies with applications to object recognition. In *IEEE Conference on Computer Vision and Pattern Recognition* (2007), pp. 1–8.

[35] SALAKHUTDINOV, R., AND HINTON, G. E. Deep boltzmann machines. In *International Conference on Artificial Intelligence and Statistics* (2009), pp. 448–455.

[36] SHIN, E. C. R., SONG, D., AND MOAZZEZI, R. Recognizing functions in binaries with neural networks. In *24th USENIX Security Symposium (USENIX Security 15)* (2015), pp. 611–626.

[37] STALLKAMP, J., SCHLIPSING, M., SALMEN, J., AND IGEL, C. Man vs. computer: Benchmarking machine learning algorithms for traffic sign recognition. *Neural Networks*, 0 (2012), –.

[38] SZEGEDY, C., ZAREMBA, W., SUTSKEVER, I., BRUNA, J., ERHAN, D., GOODFELLOW, I., AND FERGUS, R. Intriguing properties of neural networks. In *Proceedings of the 2014 International Conference on Learning Representations* (2014), Computational and Biological Learning Society.

[39] WERBOS, P. J. Generalization of backpropagation with application to a recurrent gas market model. *Neural Networks 1*, 4 (1988), 339–356.

[40] XU, W., QI, Y., AND EVANS, D. Automatically evading classifiers. In *Proceedings of the 2016 Network and Distributed Systems Symposium* (2016).

## Acknowledgments

## Appendix

This appendix contains a table and a figure used in Sections 5, 6, and 7 of the paper:

- Table 3 describes all DNN architectures used in the paper, and provides an identifier for each of them.

- Figure 12 extends Figure 11 for all architectures.

| ID | In | Out | C | MP | C | MP | RL | RL | RL | S | S | SM |
|----|-----|-----|----|----|-----|----|------|-----|-----|-----|-----|-----|
| A | 784 | 10 | 32 | ✓ | 64 | ✓ | 200 | 200 | - | - | - | 10 |
| B | 3,072 | 43 | 64 | ✓ | 128 | ✓ | 256 | 256 | - | - | - | 43 |
| C | 3,072 | 43 | 32 | ✓ | 64 | ✓ | 200 | 200 | - | - | - | 43 |
| D | 3,072 | 43 | 32 | ✓ | 64 | ✓ | 200 | 200 | - | - | - | 43 |
| E | 3,072 | 43 | 64 | ✓ | 64 | ✓ | 200 | 200 | 100 | - | - | 43 |
| F | 784 | 10 | 32 | ✓ | 64 | ✓ | 200 | - | - | - | - | 10 |
| G | 784 | 10 | 32 | ✓ | 64 | ✓ | - | - | - | - | - | 10 |
| H | 784 | 10 | 32 | ✓ | - | - | 200 | 200 | - | - | - | 10 |
| I | 784 | 10 | - | - | - | - | 200 | 200 | 200 | - | - | 10 |
| J | 784 | 10 | - | - | - | - | 1000 | 200 | - | - | - | 10 |
| K | 784 | 10 | - | - | - | - | 1000 | 500 | 200 | - | - | 10 |
| L | 784 | 10 | 32 | ✓ | - | - | 1000 | 200 | - | - | - | 10 |
| M | 784 | 10 | 32 | ✓ | - | - | - | - | - | 200 | 200 | 10 |

Table 3: **DNN architectures used in the paper:** ID: architecture identifier used in the paper, In: input dimension, Out: output dimension, C: convolutional layer with 2x2 kernels [13, 24], MP: max-pooling layer with kernel 2x2 [34], RL: rectified linear layer [28], S: sigmoid layer, SM: softmax layer [33].



Figure 12: **Frequencies of cost gradient sign matrix components equal between substitute DNNs used in our evaluation section and the oracle:** partitioned by source class (columns) and substitute architectures (rows). The same scale than Figure 11 is used here: darker shades correspond to higher probabilities.