# BRITTLE INTERPRETATIONS: THE VULNERABILITY OF TCAV AND OTHER CONCEPT-BASED EXPLAINABILITY TOOLS TO ADVERSARIAL ATTACK

**Davis Brown**[1] **& Henry Kvinge**[1,2]
[1] Pacific Northwest National Laboratory
[2] Department of Mathematics, University of Washington
`{davis.brown, henry.kvinge}@pnnl.gov`

## ABSTRACT

Methods for model explainability have become increasingly critical for testing the fairness and soundness of deep learning. A number of explainability techniques have been developed which use a set of examples to represent a human-interpretable concept in a model's activations. In this work we show that these explainability methods can suffer the same vulnerability to adversarial attacks as the models they are meant to analyze. We demonstrate this phenomenon on two well-known concept-based approaches to the explainability of deep learning models: TCAV and faceted feature visualization. We show that by carefully perturbing the examples of the concept that is being investigated, we can radically change the output of the interpretability method, e.g. showing that stripes are not an important factor in identifying images of a zebra. Our work highlights the fact that in safety-critical applications, there is need for security around not only the machine learning pipeline but also the model interpretation process.

## 1 INTRODUCTION

Deep learning models have achieved superhuman performance in a range of activities from image recognition to complex games (LeCun et al., 2015; Silver et al., 2017). Unfortunately, these gains have come at the expense of model interpretability, with massive, overparametrized models being used to achieve state-of-the-art results. This is a salient problem when deep learning is applied to domains such as healthcare (Miotto et al., 2018), criminal justice (Li et al., 2018), and finance (Huang et al., 2020), where a prediction needs to be explainable to the user, leading to a surge in interest in tools that can illuminate the underlying decision making process of deep learning models.

Besides being inherently black-box in nature, deep learning models have also been shown to be vulnerable to adversarial attacks where small perturbations to model input result in dramatic changes to model output (Szegedy et al., 2013). This phenomenon is concerning when deep learning tools are deployed in safety-critical environments. A range of approaches have been developed to improve a model's robustness to adversarial attack (Silva & Najafirad, 2020), including the use of explainability methods to detect adversarial examples (Zhang et al., 2021; Wang & Gong, 2021). But if explainability methods are an important component in a machine learning system, then the robustness of these methods are nearly as important as the robustness of the model itself. In this paper we explore the vulnerability of concept-based interpretability methods. That is, methods that interrogate a model and its decisions based on a concept.

Concept-based interpretability methods usually rely on a user provided collection of positive examples (tokens) of a concept. While this flexibility makes these methods an attractive approach for understanding deep learning models, it also introduces vulnerabilities since corrupted tokens can result in misleading interpretations of a model's decisions. We describe a threat model for concept-based interpretability methods, outlining adversary goals, knowledge, and capabilities. Then we introduce a family of attacks fitting this threat model which we call *token pushing (TP) attacks*. These learn small perturbations that when added to tokens of a concept result in remarkably different

output for the interpretability method. Specifically, we optimize our perturbations so that when they are added to a token, they significantly change a model's internal representation of the input.

We test TP attacks against two popular concept-based interpretability methods: Testing with Concept Acitvation Vectors (TCAV) (Kim et al., 2018) and Faceted Feature Visualization (FFV) (Goh et al., 2021). While TCAV and FFV are similar in that they are both concept-based, their output is quite different. TCAV quantifies the extent to which a concept is important to a model's prediction for a specific input dataset, while FFV visualizes how individual neurons represent specific aspects of a concept. We show that TP attacks are effective for both TCAV and FFV. For example, a TP attack causes TCAV to give output indicating that stripes are not an important feature to the class 'zebra'. On the other hand, a TP attack can radically change the feature visualizations generated by FFV (Figure 3).

We evaluate TP attacks on pretrained ImageNet models (Deng et al., 2009; Marcel & Rodriguez, 2010) using the Describable Textures Dataset (Cimpoi et al., 2014) for concept tokens. Through our experiments we show that a TP attack does not require the adversary to know what interpretability method is being used. The same perturbations that cause TCAV to fail, also cause FFV to fail. Finally TP attack possesses moderate transferability, meaning that as long as a surrogate model is available, it can be applied even when the defender model architecture is unknown.

In summary, our contributions in this paper include the following.

- Formalization of an adversarial threat model for concept-based interpretability methods.
- Introduction of TP attacks which perturb examples of a concept in such a way that concept-based interpretability methods give misleading output.
- Demonstration of the effectiveness of TP attacks on two concept-based interpretability methods, TCAV and FFV.

## 2 RELATED WORK AND BACKGROUND

**Interpretability methods:** Because of the size and complexity of modern deep learning architectures, skill is required to extract interpretations of how these models make decisions. Established methods range from those that focus on highlighting the importance of individual input features to those that can give clues to the importance of specific neurons to a particular class. Popular examples of interpretability methods that focus on input feature importance include saliency map methods (Selvaraju et al., 2017; Sundararajan et al., 2017; Ribeiro et al., 2016) which rely on gradients to identify those input features (pixels in an image for example) whose change is most likely to change the network's prediction.

Concept-based interpretability methods focus on decomposing the hidden layers of deep neural networks with respect to human-understandable concepts. One of the best known approaches in this direction involves the use of concept activation vectors (CAVs) (Kim et al., 2018) which we describe in detail in the next section. Work that is either related or extends these ideas includes (Zhou et al., 2018; Graziani et al., 2018; 2019).

Feature visualization is a set of interpretability techniques (Szegedy et al., 2014; Mahendran & Vedaldi, 2015; Wei et al., 2015; Nguyen et al., 2016b) concerned with optimizing model input so that it activates some specific node or set of nodes within the network. However, a challenge arises when one tries to analyze 'poly-semantic neurons' (Olah et al., 2018), neurons that activate for several conceptually distinct ideas. For example, a neuron that fires for both a boat and a cat leg is poly-semantic. Interpretability methods have imposed priors to disambiguate neurons by clustering the training images (Wei et al., 2015; Nguyen et al., 2016b) or the hidden layer activations (Carter et al., 2019) and using the average of the cluster as a coarse-grained image prior, parameterizing the feature visualization image with a learned GAN (Nguyen et al., 2016a), or using a diversity term in the feature visualization objective (Wei et al., 2015).

**Robustness of interpretability methods:** This is not the first work that has shown that interpretability methods can be brittle. Saliency methods have been shown to produce output maps that appear to point to semantically meaningful content even when they are extracted from untrained models, indicating that these methods may sometimes simply function as edge detectors (Adebayo et al.,

2018). From a more adversarial perspective, a number of works have shown that saliency methods are vulnerable to small perturbations made to either an input image or to the model itself that cause the model to offer radically different interpretations. (Heo et al., 2019; Ghorbani et al., 2019; Viering et al., 2019; Subramanya et al., 2019). To our knowledge, this is the first work that shows that concept-based interpretability methods are also vulnerable to adversarial attack.

## 2.1 TCAV AND LINEAR INTERPRETABILITY

In this section we describe the method of testing with concept activation vectors (TCAV) (Kim et al., 2018). Let $f : X \to \mathbb{R}^d$ be a neural network which is composed of $n$ layers and designed for the task of classifying whether a given input $x \in X$ belongs to one of $d$ different classes. Write $f_\ell : X \to \mathbb{R}^{d_\ell}$ for the composition of the first $\ell$ layers so that $f_n = f$ and $d_n = d$ and let $h_\ell : \mathbb{R}^{d_\ell} \to \mathbb{R}^d$ be the composition of the last $n - \ell$ layers of the network so that $f = h_\ell \circ f_\ell$ for any $1 \le \ell \le n - 1$. Let $C$ be a concept for which we have a set of positive examples (tokens) $P_C = \{x_i^P\}_i$ and negative examples $N_C = \{x_i^N\}_i$, both belonging to $X$. These are represented in the $\ell$th layer of $f$ as the points $f_\ell(P_C)$ and $f_\ell(N_C)$ respectively. One can apply a binary linear classifier to separate these two sets of points, resulting in a hyperplane in $\mathbb{R}^{d_\ell}$. This hyperplane can be represented by two unit normal vectors. We choose the one, $v_C^\ell \in \mathbb{R}^{d_\ell}$, that points into the region corresponding to the points $f_\ell(P_C)$. $v_C^\ell$ is called the *concept activation vector* in layer $\ell$ associated with concept $C$. One can think of $v_C^\ell$ as the vector that points toward $C$-ness in the $\ell$th layer of the network.

Let $h_{\ell,k}$ denote the $k$th output coordinate of $h_\ell$ corresponding to class $k$. In the classification setting, $h_{\ell,k}$ then represents the model's confidence that input belongs to class $k$. To better understand the extent to which concept $C$ influences the model's confidence of $x \in X$ belonging to class $k$ we compute:

$$S_{C,k,l} = \nabla h_{\ell,k}\left(f_\ell(x)\right) \cdot v_C^l. \tag{1}$$

A positive value of $S_{C,k,l}$ indicates that increasing $C$-ness of $x$ makes the model more confident that $x$ belongs to class $k$. The *magnitude TCAV score* for a dataset $D$ is defined as

$$\mathrm{TCAV}_{Q_{C,k,\ell}} = \frac{1}{|D_k|} \sum_{x \in D_k} S_{C,k,l}(x),$$

where $D_k$ is the subset of $D$ consisting of all instances predicted as belonging to class $k$. We compare the concept magnitude with the TCAV magnitudes for random images in the layer, and use a standard two-sided $t$-test to test for significance. We can also compute the *relative TCAV score*, which replaces the set of negative natural images in $N_C$ with images representing a distinct concept. An example experiment and attack using the relative TCAV score is in the Appendix.

## 2.2 FACETED FEATURE VISUALIZATION

Goh et al. (2021) introduced a new concept-based feature visualization objective for neuron-level interpretability, *Faceted Feature Visualization (FFV)*. The objective disambiguates poly-semantic neurons by imposing a prior towards a linear concept $C$ in the optimization objective. Goh also utilizes a set of positive and negative examples of a concept $C$ ($P_C$ and $N_C$ respectively). Similar to the TCAV method, one trains a binary linear classifier on the image of these two sets under the map $f_\ell$ to obtain $v_C^l$. To visualize output that tends to activate a neuron at layer $\ell$, position $i$, while at the same time steering the visualization toward a specific context, the authors solve the following optimization problem:

$$\arg\max_{x \in X} f_{\ell,i}(x) + v_C^l \cdot (f_\ell(x) \odot \nabla f_{\ell,i}(x)), \tag{2}$$

where $\odot$ is the Hadamard product. Note that the first term helps find $x$ which result in a strong activation of $f_{\ell,i}$, while the second term finds $x$ that tends to point in the direction of $v_C^\ell$.

## 3 ADVERSARIAL ATTACKS ON INTERPRETABILITY

In this section we describe a family of adversarial attacks on concept-based model interpretability methods. An adversarial attack (Goodfellow et al., 2015) on a model $f$ is a small perturbation $\delta$ that, when applied to a specific input $x$, results in large changes to model prediction $f(x)$. The meaning of

'small' is usually specified by a metric such as an $\ell_p$-norm and can either be a hard or soft constraint. In this work we use projected gradient descent (PGD) (Madry et al., 2018) to construct our attacks, though other optimization approaches could doubtless be used.

We frame the notion of a concept-based interpretability method abstractly in order to better understand its attack surface. We view such a method as a map that takes (1) a model, (2) positive tokens of the concept that we would like to steer our interpretation, (3) negative tokens of the concept and (4) an *interpretation target* which will be the focus of the interpretation. We call the output of an interpretability method an *interpretation object*. An interpretation object might be a single scalar value (as in the case of TCAV), or it may be an image (as in the case of FFV). In all cases, an interpretation object is designed to help the user better understand a model's decision making process. Thus, we can understand an interpretability method as a function $I : \mathcal{M} \times \mathcal{P} \times \mathcal{N} \times \mathcal{T} \to \mathcal{O}$, where $\mathcal{M}$ is the collection of models that can be interpreted, $\mathcal{P}$ is the space of all possible positive token sets, $\mathcal{N}$ is the space of all possible negative token sets, $\mathcal{T}$ is the space of interpretation targets, and $\mathcal{O}$ is the space of interpretation objects that the method produces. We note that in the case of TCAV, the interpretation target is a dataset $D_k$ of examples of some class $k$, while the interpretation target of FFV is a specific node position $(i, j, k)$ in the model.

## 3.1 A THREAT MODEL FOR ATTACKS ON CONCEPT-BASED INTERPRETABILITY METHODS

Following a suggestion given in (Carlini et al., 2019), we state the threat model that we will consider in this paper. Since we will only be considering images as input in our experiments, we specify to that setting here. Otherwise, we use the formalism that we developed above. Specifically, we assume there exists an interpretability method $I$, a model $f \in \mathcal{M}$, set of positive image tokens $P_C = \{x_i^P\} \in \mathcal{P}$, set of negative image tokens $N_C \in \mathcal{N}$, and interpretation target $T \in \mathcal{T}$. We also assume a function $F : \mathcal{O} \times \mathcal{O} \to \mathbb{R}$ that quantitatively captures meaningful difference between interpretation objects.

**Adversary's goal:** The adversary's goal is to find perturbations $\{\delta_i\}_i$ such that $\hat{P}_C = \{x_i^P + \delta_i\}$ maximizes the value of $F(I(f, P_C, N_C, T), I(f, \hat{P}_C, N_C, T))$. That is, the change from $P_C$ to $\hat{P}_C$ maximizes the difference in interpretation as measured by $F$. In order to avoid detection, $\hat{P}_C$ is subject to the constraint: $\max_i ||\delta_i||_\infty \le \epsilon$, for some fixed $\epsilon > 0$.

**Adversary knowledge and capabilities:**

1. In this paper we assume that the adversary has read and write access to the tokens $P_C$ either before or after they have been collected.

2. We do not assume that the adversary has access to either $T$ (the dataset of examples predicted as belonging to class $k$ in the case of TCAV or the specific neuron position that is being targeted in the case of FFV). We do assume that the adversary knows the hidden layer that is being targeted for both TCAV and FFV.

3. We assume that the adversary has read access to at least a surrogate model trained on the same dataset as $f$. We do not assume that this surrogate model needs to have the same architecture as $f$.

4. Finally, we do not assume that the adversary knows the interpretability method that will be used.

The adversary's goal is framed in terms of a function $F$ that depends on the specific interpretability method. This might seem to be in conflict with assumption 4 that says that the adversary does not have knowledge of the interpretability method being used. Actually, we show that TP attacks, which we propose below, work for $F$ specific to both TCAV and FFV simultaneously by optimizing for an objective function that disrupts the fundamental mechanism by which TCAV, FFV, and other concept-based interpretability methods work.

## 3.2 ATTACKING TOKENS OF A CONCEPT

In this section we introduce the *token pushing (TP) attack*. The basic idea is simple; we find perturbations $D_C = \{\delta_i\}_i$ that significantly alter a model's internal representation of the concept tokens $P_C = \{x_i^P\}_i$. Using the notation developed in 3.1, we assume that the adversary has access
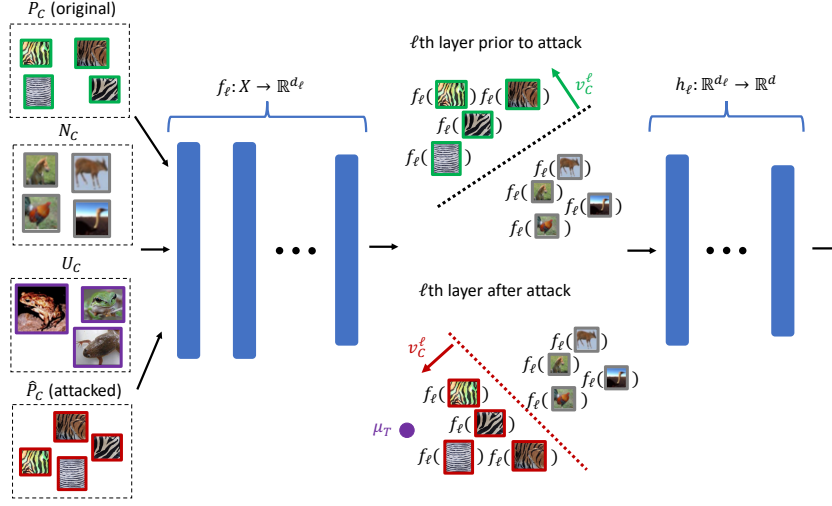
Figure 1: A schematic of the TP attack. $P_C$ is the original set of positive examples of concept $C$, $N_C$ is the set of negative examples of concept $C$, $U_C$ are the unrelated examples that are used to calculate $\mu_T$, and $\hat{P}_C$ is the set of positive examples after the attack. Note that positive examples get pulled toward unrelated concept example centroid $\mu_T$, changing the direction of $v_C^\ell$.

to a copy of the defender's model (or a surrogate model) $f : X \to \mathbb{R}^d$, the hidden layer that the interpretation method will use, and write access to the set of tokens $P_C$ that represent a concept $C$.

The perturbations added to each element in $P_C$ shifts its hidden representation in layer $\ell$ so that it no longer correlates with concept $C$. In order to find a point that can guide this shift, the first step is for the adversary to choose some collection of images that are unrelated to $C$, $U_C := \{x_i^T\}_i$. The adversary calculates the centroid of $f_\ell(U_C)$, which we denote by $\mu_T$, which will serve as a representative of "unrelatedness" to $C$. Then for each $x_i^P \in P_C$, the adversary uses PGD to compute

$$\delta_i := \underset{\|\delta\|_\infty \leq \epsilon}{\arg\min} \|f_l(x_i^P + \delta) - \mu_T\|. \tag{3}$$

This is related to the hidden layer attack that is described in (Wang et al., 2018). A schematic of TP attack can be found in Figure 1. An example of the perturbation can be found in Figure 7 in the Appendix.

In Section 4, we show that in spite of the fact that Equation 3 is neither the interpretability objective of TCAV nor FFV, it is still effective when applied to either method. In fact, objective function 3 makes the TP attack more flexible since it acts against the underlying mechanism common to both these and other interpretability methods: the spatial proximity of hidden representations of input that are semantically related. This means that the adversary does not need to know the specific interpretability method that the defender is using. This also means that the defender does not need access to the interpretation target, as they would if they were to optimize against the interpretability objective directly.

## 4 EXPERIMENTS

To better understand the effectiveness of the methods proposed in Section 3.2, we apply our attacks to TCAV and FFV in the setting where these are used to interrogate an InceptionV1 model (Szegedy et al., 2015) that has been trained on ImageNet-1k (Deng et al., 2009). We choose InceptionV1 because it is a model commonly used in the interpretability literature (Kim et al., 2018; Olah et al., 2020) and choose ImageNet-1k since it is easy to obtain high-quality weights for this model/dataset combination. In our case, we used the pretrained weights from Torchvision (Marcel & Rodriguez, 2010). The token sets that we used to capture concepts come from the Describable Textures Dataset (DTD) (Cimpoi et al., 2014). We perform all PGD attacks with $\epsilon = 4/255$ and 20 steps. We use

| | InceptionV1 Layer | | | |
| Attacks | mixed3a | mixed3b | mixed4a | mixed4b |
|---|---|---|---|---|
| Baseline TCAV (no attack) | $0.63 \pm 0.1$ | $0.80 \pm 0.1$ | $0.45 \pm 0.15$ | $0.52 \pm 0.23$ |
| Gaussian noise | $0.01 \pm 0.14$ | $0.27 \pm 0.17$ | $0.00 \pm 0.16$ | $0.02 \pm 0.20$ |
| *PGD attack on* | | | | |
| Logit | $0.01 \pm 0.14$ | $0.27 \pm 0.17$ | $0.00 \pm 0.18$ | $0.02 \pm 0.2$ |
| mixed3a centroid | $\mathbf{0.07 \pm 0.15}$ | $0.32 \pm 0.16$ | $0.03 \pm 0.12$ | $0.01 \pm 0.24$ |
| mixed3b centroid | $-0.02 \pm 0.13$ | $\mathbf{0.40 \pm 0.13}$ | $-0.07 \pm 0.14$ | $0.01 \pm 0.21$ |
| mixed4a centroid | $0.10 \pm 0.14$ | $0.39 \pm 0.15$ | $\mathbf{-0.07 \pm 0.15}$ | $0.13 \pm 0.21$ |
| mixed4b centroid | $0.12 \pm 0.16$ | $0.38 \pm 0.17$ | $0.07 \pm 0.11$ | $\mathbf{0.29 \pm 0.10}$ |
| mixed4c centroid | $0.10 \pm 0.13$ | $0.28 \pm 0.18$ | $0.19 \pm 0.15$ | $0.14 \pm 0.22$ |
| mixed4d centroid | $0.08 \pm 0.15$ | $0.37 \pm 0.13$ | $0.05 \pm 0.15$ | $0.04 \pm 0.22$ |

Table 1: The difference of the magnitude of the TCAV score for the 'striped' concept before and after the TP attacks on InceptionV1. The Baseline TCAV row gives the TCAV magnitude, all other rows are subtracted from the Baseline TCAV score. The rows below 'PGD attack on' indicate the layer that is being targeted by the TP attack. The columns are the InceptionV1 layer that TCAV is being applied to. We bold those values where the layer targeted by the TP attack and the layer TCAV is applied to are the same.

the Captum (Kokhlikyan et al., 2020) implementation of TCAV with a linear classifier trained via stochastic gradient descent and $\ell_2$-regularization. See the Appendix for more experiment details.

To test the TP attack against TCAV, we choose two concept/class pairs with straightforward associations: 'stripes'/'zebra' and 'honeycombed'/'honeycomb'. We select 20 sets of 35 randomly chosen images from ImageNet $\{N_C^i\}$ which do not intersect. The same $\{N_C^i\}$ will be used for both concept/class pairs. We also fix a set of unrelated images $U_C$ of size 1000 that are also randomly sampled from ImageNet. Finally, we choose random sets of 35 images from the classes 'stripes' $P_{\text{striped}}$ and 'honeycombed' $P_{\text{honey}}$ from DTD. $D_k$ is a collection of images which the InceptionV1 model predicts as 'zebra' or 'honeycombed' respectively.

For each layer of InceptionV1 we run the TP attack against $P_{\text{striped}}$ and $P_{\text{honey}}$. For each of the resulting pairs $(P_{\text{striped}}, \hat{P}_{\text{striped}})$ (respectively $(P_{\text{honey}}, \hat{P}_{\text{honey}})$) and each layer of InceptionV1, we then apply TCAV 20 times (once for each $N_C^i$), calculating the difference in magnitude TCAV score between $P_{\text{striped}}$ and $\hat{P}_{\text{striped}}$ (respectively $(\hat{P}_{\text{honey}}$ and $\hat{P}_{\text{honey}})$). Numerical results for 'stripes'/'zebra' can be found in Table 1. The larger the value, the greater the change in magnitude TCAV score before and after the attack. Plots of the raw TCAV magnitude scores for both the clean positive tokens and the attacked positive tokens (where each attack targeted a different layer of InceptionV1) are found in Figure 2. Sample stripe images before and after the attack can be found in the appendix.

We include 95% confidence intervals for each layer based on the 20 different $N_C^i$ sets. The point of this is to verify that the result does not depend on having the "right" negative examples. To test that the attack perturbations work for reasons than other than the fact that they are perturbations, we also apply TCAV to positive token sets to which we have added random Gaussian noise of the per-channel mean and standard deviation of the PGD logit attack. Finally, note that we also include the results showing what happens when a TP attack targets a different hidden layer than TCAV is being applied to (these are in the off-diagonal of Table 1).

We evaluate the token perturbation attack on FFV by performing feature visualizations for InceptionV1 on every channel neuron for the layers mixed3a, mixed3b, mixed4a, and mixed4b. We use the feature visualization objective equation 2, and compare feature visualizations with clean concept images $P_C$, concept images with Gaussian noise, and a concept set with perturbations created by PGD on the respective hidden layer with equation 3. We give an example of FFV and the attack in the layer mixed4d in Figure 3.

We quantitatively test the effectiveness of the attack on FFV by using a variant of the Fréchet Inception Distance (FID) (Heusel et al., 2017) as a measure of the distance between feature visualizations. Namely, we compare feature visualizations created with a channel objective (i.e., using only the first
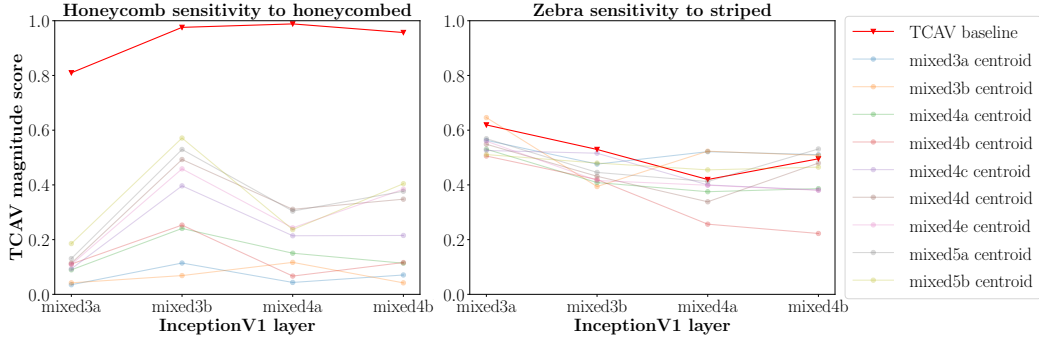
Figure 2: Adversarial attacks on the 'honeycombed' concept set (left) and the striped concept set (right). Each curve represents a TP attack targeting a different layer of InceptionV1. The $x$-axis records different layers of the InceptionV1 network and the $y$-axis is the TCAV magnitude score when TCAV is applied to that layer.
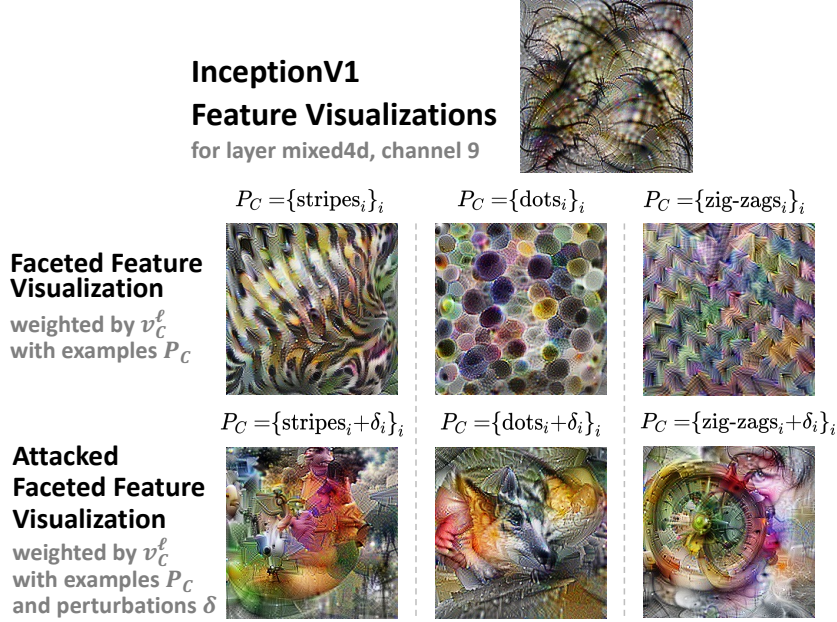


Figure 3: A feature visualization without any concept prior (first row), faceted feature visualization on the same neuron for 'stripes', 'dots', 'zig-zags' (second row), and the faceted feature visualization after it has been attacked (third row) of channel 9 on InceptionV1 layer mixed4d. Note that while visualizations in the second row reflect the concept prior, the visualizations in the third row do not.

term in equation 2), faceted feature visualizations created with two sets of random images with clean stripe concepts $P_C$, faceted feature visualizations with a set of stripe concepts $P_C$ with a perturbation created via targeting the layer mixed3b with equation 3, and faceted feature visualizations where we add Gaussian noise to $P_C$. The FID score is calculated across layers for every channel neuron in InceptionV1 layers: mixed3a (256 channels), mixed3b (480 channels), mixed4a (512 channels), and mixed4b (512 channels), shown in Figure 4. We use a PyTorch implementation of FID (Seitzer, 2020) and use the second block of InceptionV3 as the visual similarity encoder (due to the smaller dataset size).
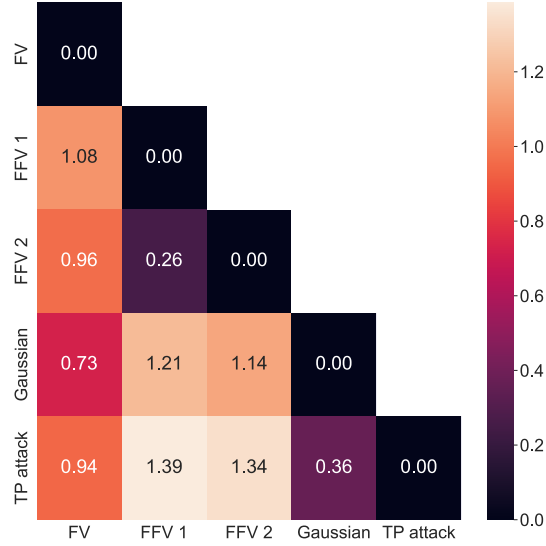
Figure 4: Fréchet Inception Distance scores for feature visualizations on the layers mixed3a, mixed3b, mixed4a, and mixed4b. Five different feature visualizations are performed: channel feature visualization (FV), two separate runs of Faceted Feature Visualization with different sets of positive and negative concept images (FFV 1 and FFV2), with Gaussian noise added to the positive concept images (Gaussian), and with the token pushing attack (TP attack).

## 5  RESULTS

Our results show that TP attacks effectively changes the output of both TCAV and FFV from the baseline interpretability results. For TCAV, we can consistently lower the TCAV magnitude score that indicates the relative importance of a concept to an output class. In Table 1, we measure the TCAV magnitude score on four early layers of InceptionV1. For each run and layer, we take the average difference between the TCAV magnitude score for the striped concept set and a random concept set over 20 sets of random images. We note that, unsurprisingly, attack success tends to increase when the layer that an attack was developed for and the layer TCAV is being applied to are the same. However, we also find that the attack is often remains effective even when these are not the same. For example, in Table 1, the attack targeting the layer 'mixed4b' is successful across all layers examined. We also observe this in Figure 2, where an attack targeting layer mixed4b for the honeycombed concept set effectively modifies the TCAV magnitude of honeycomb to honeycombed for the four layers examined.

For FFV, we can observe the TP attack effectiveness from the visual differences between 1) a channel feature visualization (i.e., a feature visualization that optimizes the first term in equation 2), 2) the faceted feature visualization with a clean concept set $P_C$, and 3) the faceted feature visualization with a perturbed concept set $\hat{P_C}$. We give three such examples separately using the striped, dotted, and zig-zagged concept sets in Figure 3. We use FID as a measure of visual difference, and test the effectiveness of the TP attack on FFV for the 1,760 channel neurons in the InceptionV1 layers 'mixed3a', 'mixed3b', 'mixed4a', and 'mixed4b'. We use the striped concept set and perform two separate FFV visualizations for each neuron with different sets of negative concept set images. Figure 4 shows that the FID scores between the separate clean FFV runs is $0.26$, while the FID score between the TP attack and the clean FFV runs are $1.39$ and $1.34$. The significantly larger FID scores suggest that the TP attack modifies the FFV output more than the variation between runs. This, along with visualizations such as 3, suggest that a TP attack can drastically change the semantic meaning associated with the feature visualizations produced by FFV.

Finally, we find that both the TCAV magnitudes (Table 1) and the FFV FID scores (Figure 4) are susceptible to Gaussian noise added to the concept set. This suggests that, even independent of adversarial attacks, concept-based interpretability methods are brittle. This brittleness likely means

that these methods are also vulnerable to natural distribution shifts in data, e.g., between the concept set and training images. We see a need for research into robust interpretability methods.
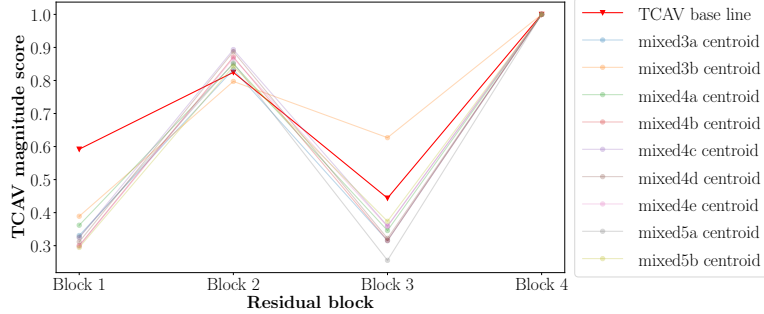


Figure 5: TCAV sensitivity scores for the zebra class with the stripe images for a ResNet-18 (He et al., 2016) trained on ImageNet-1K. The attacks uses perturbations made on the stripe concept images for InceptionV1 using centroids for different hidden layers.

## 5.1 TRANSFERRABILITY

As noted in 3.1, the knowledge required for an adversary to implement an attack is decreased significantly if they do not need to know the specific model being used by the defender. We therefore test the transferrability of TP attack by applying TCAV to an ImageNet-trained ResNet-18, before and after it has been attacked using TP perturbations developed for an InceptionV1 model. We again consider the concept/class pair stripes/zebra with the same set of $N_C$, $U_C$, and $P_{\text{stripes}}$ that were used to generate Table 1. We compute the TCAV magnitude score for stripes/zebra for each of the four residual blocks in ResNet-18. Figure 5 compares a baseline score with scores for TP attacks applied to different layers.

We find that strikingly, TP attacks targeting any of the 7 layers of InceptionV1 result in significant decreases in TCAV magnitude score when applied to the first block of ResNet18. We also see less significant decreases in TCAV magnitude score when Block 3 is targeted (with the exception of the 'mixed3b' TP attack which actually increases the TCAV score). The transfer TP attack does not seem to be effective against Block 2 and Block 4. We believe effectiveness against the first block may be the result of similar base features being extracted in both networks. This of course does not explain why we thereafter see a decrease in magnitude score in Block 3. These results point toward TP attack being moderately transferable, especially when TCAV is being applied to earlier layers of the defender's model.

## 5.2 LIMITATIONS

In this work we chose two concept-based interpretability methods to test TP attacks on. While TCAV and FFV capture some of the diversity of such methods, they do not capture the full breadth. In particular, it would be useful to understand how TP attacks behave when they are applied to other types of feature visualization methods, namely those that average over a large number of images or activations (Nguyen et al., 2016b; Carter et al., 2019) to build a visualization. Further, while we only consider image classification models, TCAV is agnostic to modality. Evaluating interpretability method brittleness in other critical modalities such as NLP would give a more complete picture of these method's vulnerabilities. Finally, the attack model that we focus on only considers perturbations to concept tokens. To fully understand the attack surfaces of concept-based interpretability methods it would make sense to look at attacks on the other inputs to the methods: the model itself, negative examples, and the interpretation targets. As a limited example, an adversarial attack may be designed to be 'triggered' for only certain token and dataset interpretation target combinations.

## 6 CONCLUSION

In this work we show that concept-based interpretability methods, like much of the deep learning modeling pipeline, are vulnerable to adversarial attacks. By subtly changing the examples of a concept that a user wishes to use to interrogate a model, an adversary can induce radically different interpretations. The attacks we describe are general enough that they work for multiple interpretability methods without modification (FFV and TCAV). We hope that these result of this paper will promote better security practices, not only around the model pipeline itself, but also around the method that is being used to interpret the model.

## 7 REPRODUCIBILITY STATEMENT

In the interest of making our results reproducible and able to be easily expanded upon, we make our codebase available to the public, including our implementations of the centroid PGD attack and faceted feature visualization we used. We also include attack and evaluation scripts with sensible defaults and examples. Finally, we provide the data used throughout this paper, including our feature visualizations. This entire repository will be available on a public GitHub repository once the anonymous review period has completed.

## REFERENCES

Julius Adebayo, Justin Gilmer, Ian Goodfellow, and Been Kim. Local explanation methods for deep neural networks lack sensitivity to parameter values. *arXiv preprint arXiv:1810.03307*, 2018.

Nicholas Carlini, Anish Athalye, Nicolas Papernot, Wieland Brendel, Jonas Rauber, Dimitris Tsipras, Ian Goodfellow, Aleksander Madry, and Alexey Kurakin. On evaluating adversarial robustness. *arXiv preprint arXiv:1902.06705*, 2019.

Shan Carter, Zan Armstrong, Ludwig Schubert, Ian Johnson, and Chris Olah. Activation atlas. *Distill*, 2019. doi: 10.23915/distill.00015. https://distill.pub/2019/activation-atlas.

M. Cimpoi, S. Maji, I. Kokkinos, S. Mohamed, , and A. Vedaldi. Describing textures in the wild. In *Proceedings of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2014.

Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.

Amirata Ghorbani, Abubakar Abid, and James Zou. Interpretation of neural networks is fragile. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pp. 3681–3688, 2019.

Gabriel Goh, Nick Cammarata †, Chelsea Voss †, Shan Carter, Michael Petrov, Ludwig Schubert, Alec Radford, and Chris Olah. Multimodal neurons in artificial neural networks. *Distill*, 2021. doi: 10.23915/distill.00030. https://distill.pub/2021/multimodal-neurons.

Ian Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *International Conference on Learning Representations*, 2015. URL `http://arxiv.org/abs/1412.6572`.

Mara Graziani, Vincent Andrearczyk, and Henning Müller. Regression concept vectors for bidirectional explanations in histopathology. In *Understanding and Interpreting Machine Learning in Medical Image Computing Applications*, pp. 124–132. Springer, 2018.

Mara Graziani, James M Brown, Vincent Andrearczyk, Veysi Yildiz, J Peter Campbell, Deniz Erdogmus, Stratis Ioannidis, Michael F Chiang, Jayashree Kalpathy-Cramer, and Henning Müller. Improved interpretability for computer-aided severity assessment of retinopathy of prematurity. In *Medical Imaging 2019: Computer-Aided Diagnosis*, volume 10950, pp. 109501R. International Society for Optics and Photonics, 2019.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.

Juyeon Heo, Sunghwan Joo, and Taesup Moon. Fooling neural network interpretations via adversarial model manipulation. *Advances in Neural Information Processing Systems*, 32:2925–2936, 2019.

Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017.

Jian Huang, Junyi Chai, and Stella Cho. Deep learning in finance and banking: A literature review and classification. *Frontiers of Business Research in China*, 14:1–24, 2020.

Been Kim, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler, Fernanda Viegas, et al. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav). In *International conference on machine learning*, pp. 2668–2677. PMLR, 2018.

Narine Kokhlikyan, Vivek Miglani, Miguel Martin, Edward Wang, Bilal Alsallakh, Jonathan Reynolds, Alexander Melnikov, Natalia Kliushkina, Carlos Araya, Siqi Yan, and Orion Reblitz-Richardson. Captum: A unified and generic model interpretability library for pytorch, 2020.

Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436–444, 2015.

Oscar Li, Hao Liu, Chaofan Chen, and Cynthia Rudin. Deep learning for case-based reasoning through prototypes: A neural network that explains its predictions. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.

Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations*, 2018.

Aravindh Mahendran and Andrea Vedaldi. Understanding deep image representations by inverting them. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 5188–5196, 2015.

Sébastien Marcel and Yann Rodriguez. Torchvision the machine-vision package of torch. In *Proceedings of the 18th ACM international conference on Multimedia*, pp. 1485–1488, 2010.

Riccardo Miotto, Fei Wang, Shuang Wang, Xiaoqian Jiang, and Joel T Dudley. Deep learning for healthcare: review, opportunities and challenges. *Briefings in bioinformatics*, 19(6):1236–1246, 2018.

Anh Nguyen, Alexey Dosovitskiy, Jason Yosinski, Thomas Brox, and Jeff Clune. Synthesizing the preferred inputs for neurons in neural networks via deep generator networks. *Advances in neural information processing systems*, 29:3387–3395, 2016a.

Anh Nguyen, Jason Yosinski, and Jeff Clune. Multifaceted feature visualization: Uncovering the different types of features learned by each neuron in deep neural networks. *arXiv preprint arXiv:1602.03616*, 2016b.

Chris Olah, Arvind Satyanarayan, Ian Johnson, Shan Carter, Ludwig Schubert, Katherine Ye, and Alexander Mordvintsev. The building blocks of interpretability. *Distill*, 2018. doi: 10.23915/distill.00010. https://distill.pub/2018/building-blocks.

Chris Olah, Nick Cammarata, Ludwig Schubert, Gabriel Goh, Michael Petrov, and Shan Carter. An overview of early vision in inceptionv1. *Distill*, 2020. doi: 10.23915/distill.00024.002. https://distill.pub/2020/circuits/early-vision.

Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. Model-agnostic interpretability of machine learning. *arXiv preprint arXiv:1606.05386*, 2016.

Maximilian Seitzer. pytorch-fid: FID Score for PyTorch. https://github.com/mseitzer/pytorch-fid, August 2020. Version 0.1.1.

Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*, pp. 618–626, 2017.

Samuel Henrique Silva and Peyman Najafirad. Opportunities and challenges in deep learning adversarial robustness: A survey. *arXiv preprint arXiv:2007.00753*, 2020.

David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of go without human knowledge. *nature*, 550(7676):354–359, 2017.

Akshayvarun Subramanya, Vipin Pillai, and Hamed Pirsiavash. Fooling network interpretation in image classification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 2020–2029, 2019.

Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *International Conference on Machine Learning*, pp. 3319–3328. PMLR, 2017.

Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.

Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In *2nd International Conference on Learning Representations, ICLR 2014*, 2014.

Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1–9, 2015.

Tom Viering, Ziqi Wang, Marco Loog, and Elmar Eisemann. How to manipulate cnns to make them lie: the gradcam case. *arXiv preprint arXiv:1907.10901*, 2019.

Bolun Wang, Yuanshun Yao, Bimal Viswanath, Haitao Zheng, and Ben Y Zhao. With great training comes great vulnerability: Practical attacks against transfer learning. In *27th {USENIX} Security Symposium ({USENIX} Security 18)*, pp. 1281–1297, 2018.

Shen Wang and Yuxin Gong. Adversarial example detection based on saliency map features. *Applied Intelligence*, pp. 1–14, 2021.

Donglai Wei, Bolei Zhou, Antonio Torrabla, and William Freeman. Understanding intra-class knowledge inside cnn. *arXiv preprint arXiv:1507.02379*, 2015.

Zhun Zhang, Qihe Liu, and Shijie Zhou. Ggcad: A novel method of adversarial detection by guided grad-cam. In *International Conference on Wireless Algorithms, Systems, and Applications*, pp. 172–182. Springer, 2021.

Bolei Zhou, Yiyou Sun, David Bau, and Antonio Torralba. Interpretable basis decomposition for visual explanation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 119–134, 2018.

# A APPENDIX

## A.1 EXPERIMENT DETAILS

To run TCAV, FFV, and our attacks, we use PyTorch with an NVIDIA Tesla T4 GPU provided with Google Colab Pro as well as a single NVIDIA Tesla P100 GPU.

## A.2 TP ATTACK ON RELATIVE TCAV

Here, we give an example experiment showing that TP attacks are also effective for a variant of TCAV, using relative TCAV scores. The results in Figure 6 uses concept sets for stripes, ziz-zags, and polka-dots of 35 images each. Perturbations are made on the striped concept set using the final logit layer, towards an unrelated class (the $999^{th}$ 'toilet tissue' ImageNet-1k class).
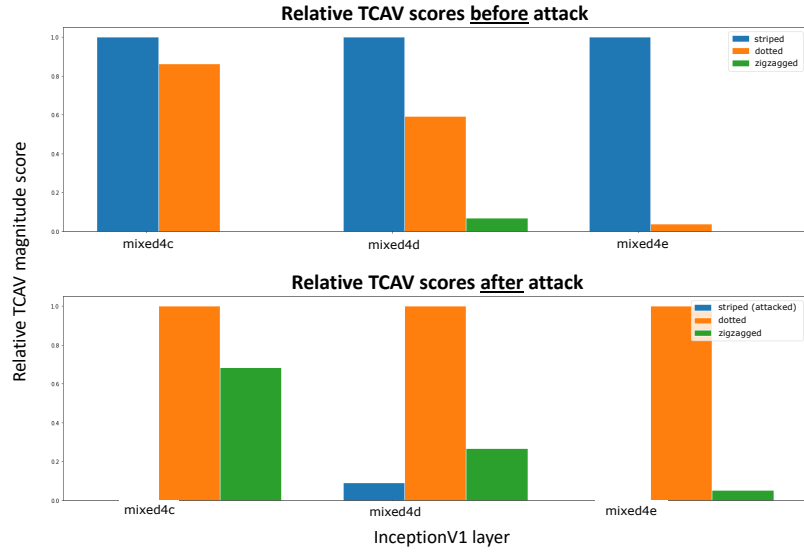
Figure 6: Relative TCAV magnitude scores before (top) and after (bottom) the PGD logit attack on the striped concept images for the striped, zig-zagged, and dotted concept sets.
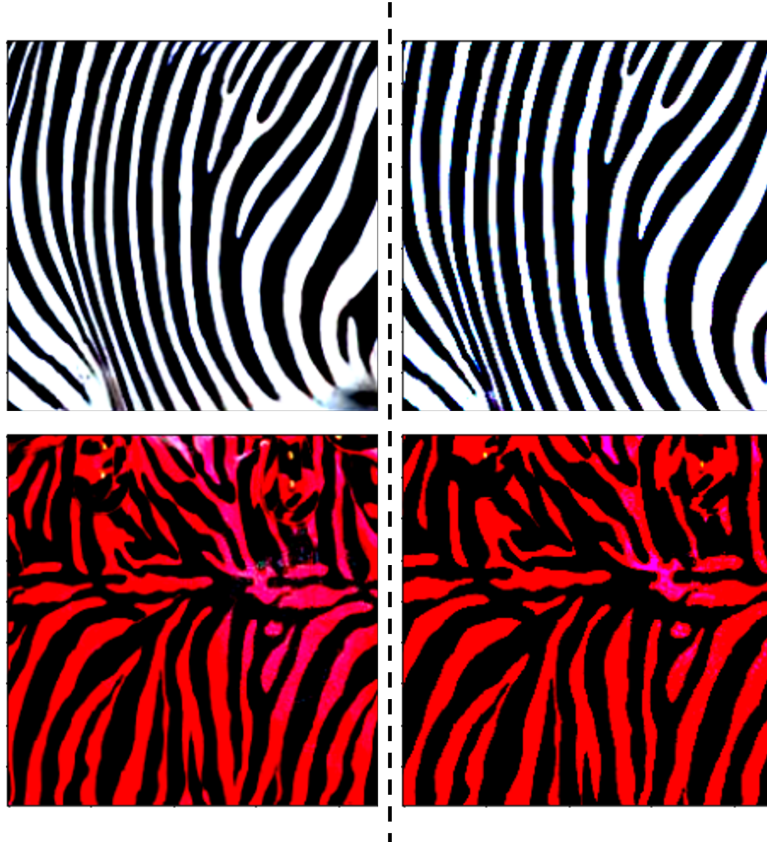


Figure 7: Example of stripe concept images before (left) and after (right) a TP attack. We use $\epsilon = 4/255$ and 20 iterations for all PGD experiments. The perturbation shown targets InceptionV1 layer mixed3a.