

Evaluating the Adversarial Robustness of CNNs Layer by Layer

Anonymous authors

Paper under double-blind review

Abstract

In order to measure the adversarial robustness of a feature extractor, Bhagoji et al. introduced a distance on example spaces measuring the minimum perturbation of a pair of examples to achieve identical feature extractor outputs. They related these distances to the best possible robust accuracy of any classifier using the feature extractor. By viewing initial layers of a neural network as a feature extractor, this provides a method of attributing adversarial vulnerability of the classifier as a whole to individual layers. However, this framework views any injective feature extractor as perfectly robust: any bad choices of feature representation can be undone by later layers. Thus the framework attributes all adversarial vulnerabilities to the layers that perform dimensionality reduction. Feature spaces at intermediate layers of convolutional neural networks are generally much larger than input spaces, so this methodology provides no information about the contributions of individual layers to the overall robustness of the network. We extend the framework to evaluate feature extractors with high-dimensional output spaces by composing them with a random linear projection to a lower dimensional space. This results in non-trivial information about the quality of the feature space representations for building an adversarial robust classifier.

1 Introduction

Over the last several years, a line of work has characterized the optimal adversarial loss achievable for a particular data distribution and adversarial constraint specification Bhagoji et al. (2019; 2021); Pydi & Jog (2020); Trillos et al. (2023); Dai et al. (2023). These works characterize the problem of minimizing adversarial loss over all possible classifier functions to some kind of optimal transport problem on the class conditional distributions. In order to set up this problem, one needs the distance structure of the data distribution with regard to the distance appearing in the adversarial constraint. These characterizations of optimal loss provide a lower bound for loss achievable using a particular architecture and a benchmark performance level that a trained classifier can be compared to. At a higher level, they provide an objective method to determine whether a given adversarial constraint is reasonable to use with a particular dataset.

More recently, this method has been extended to analyze the robustness of a feature extractor, or equivalently the initial layers of a feed-forward network Bhagoji et al. (2023). Just as the minimum loss on a data distribution is a minimization over all possible classifiers, the robustness of a feature extractor is a minimization over all classifiers using that feature extractor. The robustness of a feature extractor can be computed via a characterization that uses the distance structure of the data distribution using a specially defined distance that takes into account the behavior of the feature extractor. This allows the adversarial loss of the model to be decomposed layer-by-layer: including a layer of the model in the feature extractor decreases the size of the class of models using that feature extractor and thus increases the minimum loss. This increase can be considered the contribution of that layer to the overall adversarial loss of the model.

Unfortunately, for wide models, the decomposition tends to be uninformative. Layers are assigned none of the adversarial loss until dimensionality reduction to below the input space dimension occurs. Since the development of CNNs, state of the art architectures have used huge feature spaces. Dimensionality reduction has to occur eventually, but this methodology often results in the attribution of all the loss to the final layer.

In this paper, we propose a variation on the feature extractor distance tailored to wide networks, or equivalently feature extractors with high dimension feature spaces. We explore basic properties of this feature extractor distance and propose and implement an algorithm to approximately compute it for ReLU CNNs feature extractors. We use the algorithm to perform some initial exploration into the effects of architectural and training method decisions on layer-by-layer robustness.

As pre-trained models and the transfer learning paradigm are more widely used, analysis of these models as feature extractors rather than fully specified classifiers become more and more relevant. We illustrate the potential feasibility of the layer-by-layer approach to this task.

At this point, adversarial examples have been known to be a serious problem in neural network for over a decade Szegedy et al. (2013); Goodfellow et al. (2015). Adversarial training has succeeded at mitigating this vulnerability to an extent, but is computationally expensive and does not achieve information-theoretic limits on adversarial robustness Bhagoji et al. (2021). Certified methods give stronger guarantees, but are usually even more computationally expensive and limited in the scale of the adversarial budget.

We would like to provide new analysis tools to enable continued progress on training adversarially robust models. With these tools, we can better identify different behaviors in models as training method is varied. By localizing adversarial vulnerabilities within a model, we allow computational resources to be spent efficiently.

Contributions. We introduce a novel collision search algorithm [that leverages gradient information to evaluate differentiable feature extractors](#). We provide a theoretical framework for understanding the robustness of feature extractors [in high-dimensional feature spaces](#). We utilize random projections to transform feature extractors with high-dimensional feature spaces to representative low-dimensional versions whose robustness can be efficiently analyzed. From our experiments [on convolutional networks](#), we find that different training methods produce significant variation in layer-by-layer robustness.

2 Related Work

Robustness of Feature Extractors. Our work primarily extends the theoretical framework established by Bhagoji et al. (2023), which links adversarial robustness to optimal transport distances. However, a significant limitation of prior work is its reliance on combinatorial search methods that scale poorly with depth. Also, it applies only to restricted activation types. We advance this framework in two key ways: First, we introduce a gradient-based collision search algorithm that evaluates robustness for *any* differentiable feature extractor, enabling efficient scaling to deep architectures supported by various differentiable activation functions. Second, we integrate *random projections* to handle high-dimensional feature spaces, overcoming the dimensionality constraints that limited previous theoretical evaluations.

Layer-wise Analysis and Model Utilization. Understanding how robustness and information flow evolve through network layers is an active area of research. Song et al. (2021) focused on optimizing layer-wise quantization under Lipschitz bounds, while Craig et al. (2023) derived bounds for information flow through subnetworks. More recently, Gavrikov et al. (2024) investigated how training schemes (including adversarial training) alter layer utilization, identifying which layers become critical for performance.

While Gavrikov et al. (2024) focuses on functional criticality, our work provides a *geometric* lower bound on untargeted robustness via feature collisions. Unlike prior approaches that bound Lipschitz constants or analyze parameter utilization, our method identifies concrete collision points. This dual functionality allows us to not only measure vulnerability but also generate the specific adversarial examples that exploit it, offering a direct probe into the geometry of middle-layer representations.

Random Projections and Optimal Transport. Our use of Gaussian random projections is philosophically rooted in the Sliced Wasserstein (SW) distance. Optimal adversarially robust classification is intimately connected to optimal transport distances between class distributions Bhagoji et al. (2019); Pydi & Jog (2020); Trillos et al. (2023). However, estimating these distances suffers from the curse of dimensionality. SW distance bypasses this by projecting high-dimensional measures onto lower-dimensional subspaces Kolouri et al. (2018); Deshpande et al. (2019); Bonneel et al. (2015); Nadjahi et al. (2021). Whereas SW methods typically average

1-D distances over many slices, we employ a higher-dimensional (e.g., 10-D) random projection to approximate the optimal loss, balancing computational efficiency with the preservation of geometric structure.

3 Framework

In this section, we build toward the definition of our robustness measure by discussing key technical aspects of prior work. Section 3.1 explains the characterized the optimal adversarial cross-entropy loss in terms of a convex optimization problem with constraints described by a bipartite conflict graph. Section 3.2 explains an earlier definition of a robustness measure for feature extractors based on the optimal adversarial loss of a classifier built on that feature extractor.

We consider the following adversarial classification scenario. (For a comprehensive summary of notation, please refer to Table ?? in Appendix A.) There is an example space $\mathcal{X} = \mathbb{R}^{n_0}$ equipped with a distance metric d and a label space $\mathcal{Y} = \{0, 1, \dots, k-1\}$. A labeled example (x, y) is sampled from a distribution P over $\mathcal{X} \times \mathcal{Y}$. An adversary selects a point $\tilde{x} \in \mathcal{X}$ such that $d(x, \tilde{x}) \leq \varepsilon$. The classifier receives \tilde{x} , uses a classifier h to produce a probability vector $h(\tilde{x}) \in \mathbb{P}(\mathcal{Y})$, and incurs loss $\ell(h(\tilde{x}), y)$. The expected adversarial loss of h is $\tilde{L}(\varepsilon, P, h) = E_{(x,y) \sim P}[\sup_{\tilde{x}: d(x, \tilde{x}) \leq \varepsilon} \ell(h(\tilde{x}), y)]$.

3.1 Characterizing Optimal Adversarial Loss

We briefly summarize the characterization of optimal adversarial cross entropy loss for a two-class discrete data distribution Bhagoji et al. (2021).

The distance structure of the data can be summarized in a bipartite *conflict graph* $(\mathcal{V}, \mathcal{E}_\varepsilon)$ in which, the vertices \mathcal{V} represent labeled examples, and the edges \mathcal{E}_ε indicate the presence of intersections under a given attack budget ε . When such an intersection exists, it signifies a fundamental limitation: no classifier can correctly classify both points in the adversarial setting. In the graph the left vertices are examples from class zero, the right vertices are examples from class one, and an edge (x, x') is present when $d(x, x') \leq 2\varepsilon$. The polytope $\mathcal{P}_\varepsilon \subset \mathbb{R}^V$ is defined by the following inequalities:

$$0 \leq q_v \leq 1, v \in \mathcal{V}; \quad q_v + q_{v'} \leq 1, (v, v') \in \mathcal{E}_\varepsilon.$$

For a labeled example (x, y) and soft classifier $h : \mathcal{X} \rightarrow \mathbb{P}(\mathcal{Y})$, $h(x)_y$ is the probability that h assigns to the correct class. Bhagoji et al. show that there is an h such that $\inf_{\tilde{x}: d(x, \tilde{x}) \leq \varepsilon} h(\tilde{x})_y \geq q_{(x,y)}$ for each labeled example (x, y) in the support of the data distribution if and only if $q \in \mathcal{P}_\varepsilon$. Thus \mathcal{P}_ε is the space of achievable correct-classification-probability vectors.

Using this characterization of the correct classification vectors, we obtain a convex optimization problem for the minimum adversarial loss:

$$\tilde{L}^*(\varepsilon, P) = \min_h \tilde{L}(\varepsilon, P, h) = \min_{q \in \mathcal{P}_\varepsilon} \sum_{v \in \mathcal{V}} p_v \log \frac{1}{q_v}.$$

The optimal loss represents the lower bound on the loss that any possible classifier can achieve when attempting to correctly classify an adversarial example under the current perturbation.

3.2 A Previous Approach to Feature Extractor Robustness

This analysis has been extended to multi-layer classifiers. Consider a classifier $h : \mathbb{R}^{n_0} \rightarrow \mathbb{P}(\mathcal{Y})$ that is the composition of ℓ layers: $h = \text{softmax} \circ h_\ell \circ h_{\ell-1} \circ \dots \circ h_1$, where $h_i : \mathbb{R}^{n_{i-1}} \rightarrow \mathbb{R}^{n_i}$ and $n_\ell = k = |\mathcal{Y}|$. We can consider its initial j layers $f_j : \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_j}$, $f_j = h_j \circ \dots \circ h_1$ as a *feature extractor*. The key idea of Bhagoji et al. is to evaluate any feature extractor f via the set of all classifiers that use f as their initial layers: $\mathcal{H}_f = \{g \circ f \text{ such that } g : \mathbb{R}^{n_j} \rightarrow \mathbb{P}(\mathcal{Y})\}$ Bhagoji et al. (2023). To measure the adversarial robustness of f , they consider the best achievable adversarial loss over this family: $\tilde{L}^*(\varepsilon, P, f) = \inf_{h \in \mathcal{H}_f} \tilde{L}(\varepsilon, P, h)$.

This definition has some nice properties. As more layers are considered to be part of the feature extractor, the robustness can only decrease because $\mathcal{H}_{g \circ f} \subseteq \mathcal{H}_f$. The set of all classifiers is \mathcal{H}_{id} , so the identity feature extractor has perfect robustness, i.e. no reduction from input space loss lower bound.

However, there is a key weakness to the definition: *any* injective feature extractor achieves perfect adversarial robustness. A feature extractor $f : \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_j}$ will typically be injective or close enough to only introduce a negligible amount of loss. By this definition, only feature extractors that perform significant dimensionality reduction typical cause loss of robustness.

Colliding adversarial examples and feature layer distance Let f be a feature extractor and let (x, y) and (x', y') be labeled example from different classes ($y \neq y'$). A pair adversarial examples collide when the perturbations cause the examples to be mapped to the same point in the feature space: \tilde{x} and \tilde{x}' satisfy $d(x, \tilde{x}) \leq \epsilon$, $d(x', \tilde{x}') \leq \epsilon$, and $f(\tilde{x}) = f(\tilde{x}')$. This guarantees a classification error for at least one example for any backend g : $g(f(\tilde{x})) = g(f(\tilde{x}'))$. The minimum adversarial budget required to find inputs that collide in the feature space defines a distance between examples:

$$d_f(x, x') = \inf_{\tilde{x}, \tilde{x}' : f(\tilde{x}) = f(\tilde{x}')} \max(d(x, \tilde{x}), d(x', \tilde{x}')).$$

Computing the Robustness Measure The distance d_f can be used in place of the input distance d to build a conflict graph and characterize optimal adversarial loss over \mathcal{H}_f . Define $\mathcal{E}_{f, \epsilon}$, the edge set of the conflict graph for d_f , analogously to \mathcal{P}_ϵ : $\mathcal{E}_{f, \epsilon} = \{(x, x') : d_f(x, x') \leq \epsilon\}$. Define $\mathcal{P}_{f, \epsilon} \subset \mathbb{R}^V$ is defined analogously to \mathcal{P}_ϵ , but with constraints coming from $\mathcal{E}_{f, \epsilon}$.

Bhagoji et al. show that for a discrete distributions P , $\tilde{L}^*(\epsilon, P, f)$ can be characterized similarly to $\tilde{L}^*(\epsilon, P)$:

$$\tilde{L}^*(\epsilon, P, f) = \inf_{h \in \mathcal{H}_f} \tilde{L}(\epsilon, P, h) = \min_{q \in \mathcal{P}_{f, \epsilon}} \sum_{v \in V} p_v \log \frac{1}{q_v}.$$

If $\mathcal{E}_{f, \epsilon}$ is known, then $\tilde{L}^*(\epsilon, P, f)$ can be efficiently computed by solving a convex optimization problem.

4 Measuring Robustness of High-dimensional Feature Extractor

To evaluate a feature extractor $f : \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_j}$, we consider $\tilde{L}^*(\epsilon, P, a \circ f)$, where $a : \mathbb{R}^{n_j} \rightarrow \mathbb{R}^m$ is a random linear function: $a(x) = Ax$, $A \in \mathbb{R}^{m \times n_j}$ is a matrix with i.i.d. standard Gaussian entries. Note that because a is a random function, $d_{a \circ f}(x, x')$ and $\tilde{L}^*(\epsilon, P, a \circ f)$ are random variables.

Discussion An injective feature extractor does not commit to any aspect of the behavior of the classifier using it: subsequent layers of the classifier could invert it. However, a feature extractor causes some classifiers extending it to have simpler representations than others. We would like to measure the robustness of a “typical” classifier extending f and we use the composition of a random dimension-reducing linear function with f as a proxy for this that is computationally tractable to analyze. The random linear function is comparable to a randomly initialized but untrained classifier network using the features.

4.1 Properties of the Robustness Measure

Now we discuss key properties of the feature extractor distance $d_{a \circ f}$ and the robustness measure $\tilde{L}^*(\epsilon, P, a \circ f)$.

Invariance of d_f to scaling of f or a For any $c \neq 0$, $\tilde{L}^*(\epsilon, P, cf) = \tilde{L}^*(\epsilon, P, f)$. A collision $Af(\tilde{x}) = Af(\tilde{x}')$ for random linear layer A is still a collision after scaling by a constant. This also means that the scale of the distribution of A does not affect our definitions.

Rotation invariance The distribution of random linear maps is rotation invariant, so any rotation of the feature space leaves all distances unchanged. When $m \leq n_j$, the random matrix A has rank m with probability one. It has singular value decomposition $A = U\Sigma V^T$, where $V \in \mathbb{R}^{n_j \times m}$ and $U \in \mathbb{R}^{m \times m}$ are orthogonal and $\Sigma \in \mathbb{R}^{m \times m}$ is diagonal and invertible. Then $U\Sigma V^T f(\tilde{x}) = U\Sigma V^T f(\tilde{x}')$ if and only if $V^T f(\tilde{x}) = V^T f(\tilde{x}')$. Thus using i.i.d. Gaussian matrices results in the same distance as a uniform distribution over orthogonal projections to a dimension- d subspace.

Invertible function non-invariance On the other hand, applying a generic invertible linear transformation b to the feature space \mathbb{R}^{n_j} *does* affect the distribution of the distance: it is not true that $\tilde{L}^*(\varepsilon, P, b \circ f) = \tilde{L}^*(\varepsilon, P, f)$. Applying the linear transformation to $f(x)$ is equivalent to applying the inverse to the random matrix A , which changes its distribution. All linear maps still appear, but with skewed densities.

Adding, removing, rescaling features While d_f does not depend on the absolute scaling of features, it does depend on the relative scaling. It is easier for the adversary to satisfy $Af(\tilde{x}) = Af(\tilde{x}')$ by modifying the examples in ways that change large magnitude features.

If we modify a feature extractor f to add an additional feature, giving f' , this could increase the distance between x and x' or decrease it. Adding a constant feature leaves the distance unchanged.

Figure 1 illustrates how a “good” feature extractor (where level sets are separated) can become a “bad” one (where level sets are cluttered), when enough weak features are added.

Specializing to ℓ_2 In the remainder of this section and in our experiments, we use ℓ_2 distance as our input space distance.

4.2 Linear Feature Extractors

Suppose that f is affine linear, so $f(x) = Hx + k$. Suppose that $\text{rank}(H) = r$, the singular value decomposition of H is $H = U\Sigma V^T$ where $\Sigma \in \mathbb{R}^{r \times r}$. Then Bhagoji et al. observe that the optimal adversarial examples are $\tilde{x} = x + \delta$ and $\tilde{x}' = x' - \delta$ where $\delta = \frac{1}{2}VV^T(x - x')$. The perturbation is the projection of half the difference of examples onto the orthogonal complement of the nullspace of H : the adversary should not allocate any of their budget to moving within the level sets of f . Thus $d_f(x, x')$ is $\frac{1}{2}\|V^T(x - x')\|$. The left part of Figure 2 illustrates this.

Proposition 1 (Informal version). *Let $f(x) = ALx$, where $L \in \mathbb{R}^{n_1 \times n_0}$ is a fixed linear feature extractor and the entries of $A \in \mathbb{R}^{m \times n_1}$ are i.i.d. standard Gaussian. Let $r_4 = \frac{\|L\|_4^4}{\|L\|_F^4}$, where $\|L\|_4$ is a Schatten p -norm and $\|L\|$ is the ℓ_2 operator norm. Then r_4 is a variation of stable rank for L .*

If $1 \ll m \ll \sqrt{r_4}$, then $d_f(x, x')$ concentrates tightly around $\frac{\sqrt{m}\|L(x - x')\|}{2\|L\|_F}$.

Precise statement and proof are in Appendix ?? . Note that $\|L(x - x')\|$ is just the distance between the feature representations of x and x' . If we were only interested in linear layers, computing $d_{a \circ h}$ would not be justified: $\|L(x - x')\|$ and $\|L\|_F$ are much more straightforward to compute, and not random! However, these norms do not obviously generalize to a nonlinear h , which would have different linear approximations at

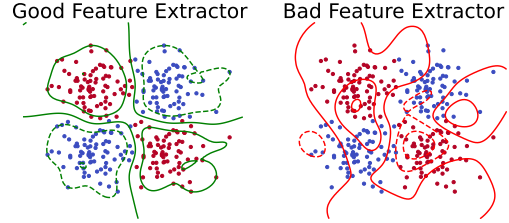


Figure 1: Random non-linear level sets for compositions of a random linear function with good and bad feature extractors. The lines represent the level sets of a random linear function composed with the feature extractor. The bad feature extractor contains all of the features in the good extractor plus many more. Most of these are low quality so a random linear combination has low robustness.

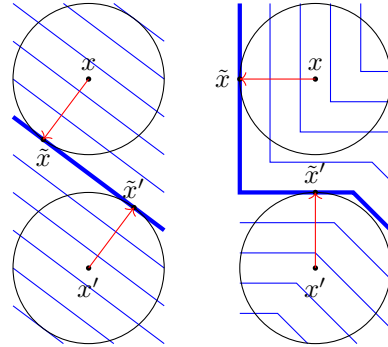


Figure 2: Visualization of $d_f(x, x')$ for two feature extractors $f: \mathbb{R}^2 \rightarrow \mathbb{R}^1$, one linear and one piecewise-linear. Blue lines are level sets of f . Each of the balls has radius $d_f(x, x')$. The level set of the colliding adversarial examples is highlighted. These are not classifiers and the highlighted line is not a decision boundary. When the feature space is more than one dimension smaller than the input space, level sets do not disconnect the input space.

different points in the input space, while $d_{a \circ h}$ does still make sense. When evaluating nonlinear h , we would still like to pick m so that $d_{a \circ h}$ concentrates so that we can efficiently and accurately estimate it by sampling.

Input space calculation The special case of $d_a(x, x')$ for a random linear function a is more intuitive. This corresponds to $L = I \in \mathbb{R}^{n_0 \times n_0}$ and $\|L\|_F = n_0$. The orthogonal complement of the nullspace of $A \in \mathbb{R}^{m \times n_0}$ is a uniformly distributed dimension- m subspace. So $d_a(x, x')^2$ has mean $\frac{m}{4n_0}\|x - x'\|_2^2$ and $d_a(x, x')$ concentrates around $\frac{\sqrt{m}}{2\sqrt{n_0}}\|x - x'\|_2$ when m is large enough.

4.3 Relationship between distances at different layers

Now we compare $d_{a \circ f_j}(x, x')$ and $d_{a \circ f_{j+1}}(x, x')$. As we noted before, $d_{f_j}(x, x') \geq d_{f_{j+1}}(x, x')$: adding an additional layer can only decrease distances by introducing new collisions between inputs. In contrast, $d_{a \circ f_{j+1}}$ could be either larger or smaller than $d_{a \circ f_j}$. Including an additional layer h_{j+1} to the feature extractor introduces new possibilities for collisions, which reduces robustness, but it also reweights the existing feature information in a way that may reduce the effect of weak features. Figure 3 presents paired adversarial examples identified for the 7-layer model. We observe that as the feature extractor becomes deeper, the average adversarial distance $d_{a \circ f_i}(x, x')$ increases as well.

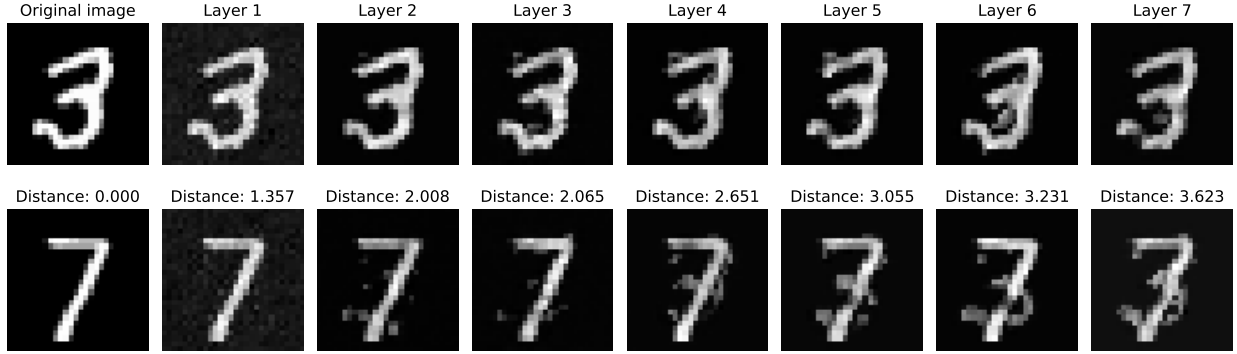


Figure 3: Colliding adversarial examples and their ℓ_2 perturbation distance for one pair across different layers.

4.4 Varying the output dimension of the random projection

As we vary the output dimension m of the random linear function, the distance $d_{a \circ f}$ reveal different information about the structure of the feature space representation. At one extreme, if $m \geq n_j$, then the linear map a is injective with probability one and $d_f(x, x') = d_{a \circ f}(x, x')$. If $m = 1$, we roughly learn how much of the feature representation is usable by a linear classifier. By using a single fixed m to analyze multiple layers with different feature space dimensions n_i , we eliminate some the effects of this dimension variation on measured robustness. For most of our experiments, we use $m = 10$ to facilitate comparison with the logits.

5 Algorithm for approximating feature-layer distances

When we use ℓ_2 distance to constrain our adversary, for $f : \mathbb{R}^{n_0} \rightarrow \mathbb{R}^d$, $d_f(x, x')$ becomes

$$\min_{\tilde{x}, \tilde{x}'} \|\tilde{x} - x\|_2^2 + \|\tilde{x}' - x'\|_2^2 \quad \text{subject to} \quad f(\tilde{x}) = f(\tilde{x}')$$

This is a slight relaxation: to get a smoother problem we have replaced the original objective $\max(\|\tilde{x} - x\|_2, \|\tilde{x}' - x'\|_2)$ with the sum of squares. Now we present the algorithm that we use to approximately solve this optimization problem.

Our algorithm is related to a method used to compute distance from an example to the decision boundary of a classifier. The linear approximated decision hyperplane is used in DeepFool Moosavi-Dezfooli et al. (2016),

and the FAB attack paper Croce & Hein (2020) employs convex combinations to make the solution closer to the original point.

Algorithm 1 Collision Search Algorithm

Require: Original points x, x' ; initial step $\alpha \in (0, 1)$; decay rate $\gamma \in (0, 1)$; maximum number of iteration N_{\max} ; convergence threshold ϵ .

- 1: Initialize $\tilde{x}^{(0)} = (x + x')/2$, $\tilde{x}'^{(0)} = (x + x')/2$, $i = 0$.
- 2: **for** $i = 0 \dots N_{\max} - 1$ **do**
- 3: Compute Jacobians, build collision subspace:

$$\mathcal{H}_i = \{ (\tilde{x}, \tilde{x}') \in \mathbb{R}^{2n_0} : f(\tilde{x}'^{(i)}) - f(\tilde{x}^{(i)}) + Df(\tilde{x}'^{(i)})(\tilde{x}' - \tilde{x}'^{(i)}) - Df(\tilde{x}^{(i)})(\tilde{x} - \tilde{x}^{(i)}) = 0 \}$$

- 4: Convex combination and projection:

$$(\tilde{x}^{(i+1)}, \tilde{x}'^{(i+1)}) = \text{proj} \left((1 - \alpha^{(i)})(\tilde{x}^{(i)}, \tilde{x}'^{(i)}) + \alpha^{(i)}(x, x'), \mathcal{H}_i \right)$$

- 5: Break if converged:

$$\left\| f(\tilde{x}'^{(i+1)}) - f(\tilde{x}^{(i+1)}) \right\| < \epsilon \quad \text{and} \quad (\|\tilde{x}^{(i)} - x\|_2^2 + \|\tilde{x}'^{(i)} - x'\|_2^2) - (\|\tilde{x}^{(i+1)} - x\|_2^2 + \|\tilde{x}'^{(i+1)} - x'\|_2^2) < \epsilon$$

- 6: $\alpha^{(i+1)} = \alpha^{(i)} \times \gamma$.

- 7: **end for**

- 8: **Return** colliding adversarial examples $\tilde{x}^{(i+1)}, \tilde{x}'^{(i+1)}$.
-

The algorithm can be viewed as an approximate version of projected gradient descent. Because f is a nonconvex piecewise-linear function, it is intractable to project onto the true constraint set $f(\tilde{x}) = f(\tilde{x}')$. Instead, we locally approximate the constraint set as an affine subspace. The derivative of $(\tilde{x}, \tilde{x}') \mapsto f(\tilde{x}) - f(\tilde{x}')$ at (\tilde{x}, \tilde{x}') is the matrix $(Df(\tilde{x}) \quad -Df(\tilde{x}')) \in \mathbb{R}^{m \times 2n_0}$, where D represents the Jacobian matrix of the gradient. The affine subspace on which the linear approximation is zero has dimension at least $2n_0 - m$. In contrast to related algorithms to compute minimum distance to the classification boundary, the subspace usually has codimension greater than one. i.e. is not a hyperplane. Thus projection onto the subspace requires us to solve a linear system and it is preferable to choose m to be not too large. Details are shown in Appendix A.

We know that $\tilde{x} = \tilde{x}'$ implies $f(\tilde{x}) = f(\tilde{x}')$, so we initialize \tilde{x} and \tilde{x}' by minimizing the objective subject to this constraint. Thus we set $\tilde{x}^{(0)} = \tilde{x}'^{(0)} = (x + x')/2$, the midpoint of the original examples. We found that starting with points satisfying the collision condition worked much better than any algorithm initialized at the original examples ($\tilde{x}^{(0)} = x$ and $\tilde{x}'^{(0)} = x'$).

5.1 Convex Combination Projection

We address the NP-hard collision-search problem Katz et al. (2017) via a simple iterative method that alternates between convex combination and projection onto a locally linearized subspace. At iteration i , this subspace is \mathcal{H}_i . If the subspace were the true collision manifold, then the minimizer would be $\text{proj}((x, x'), \mathcal{H}_i)$. However, \mathcal{H}_i is only an accurate approximation locally, so we improve the objective by mixing the projection of the current estimate $(\tilde{x}^{(i)}, \tilde{x}'^{(i)})$ with the projection of original point (x, x') :

$$\begin{aligned} (\tilde{x}^{(i+1)}, \tilde{x}'^{(i+1)}) &= (1 - \alpha^{(i)}) \text{proj}((\tilde{x}^{(i)}, \tilde{x}'^{(i)}), \mathcal{H}_i) + \alpha^{(i)} \text{proj}((x, x'), \mathcal{H}_i) \\ &= \text{proj}((1 - \alpha^{(i)})(\tilde{x}^{(i)}, \tilde{x}'^{(i)}) + \alpha^{(i)}(x, x'), \mathcal{H}_i) \end{aligned}$$

The second expression for the iterate holds because the projection is affine. Using this saves computation by performing the convex combination first.

Since the midpoint $\frac{1}{2}(x + x')$ is itself a perfect collision in function space, we initialize our search there and compute its Jacobian to define the first affine subspace. In the very first projection step, we set $\alpha^{(0)} = 1$, so that $\tilde{x}^{(1)} = \text{proj}(x, \mathcal{H}_0)$, using both original points to project onto \mathcal{H}_0 . Thereafter, we gradually reduce $\alpha^{(i)}$

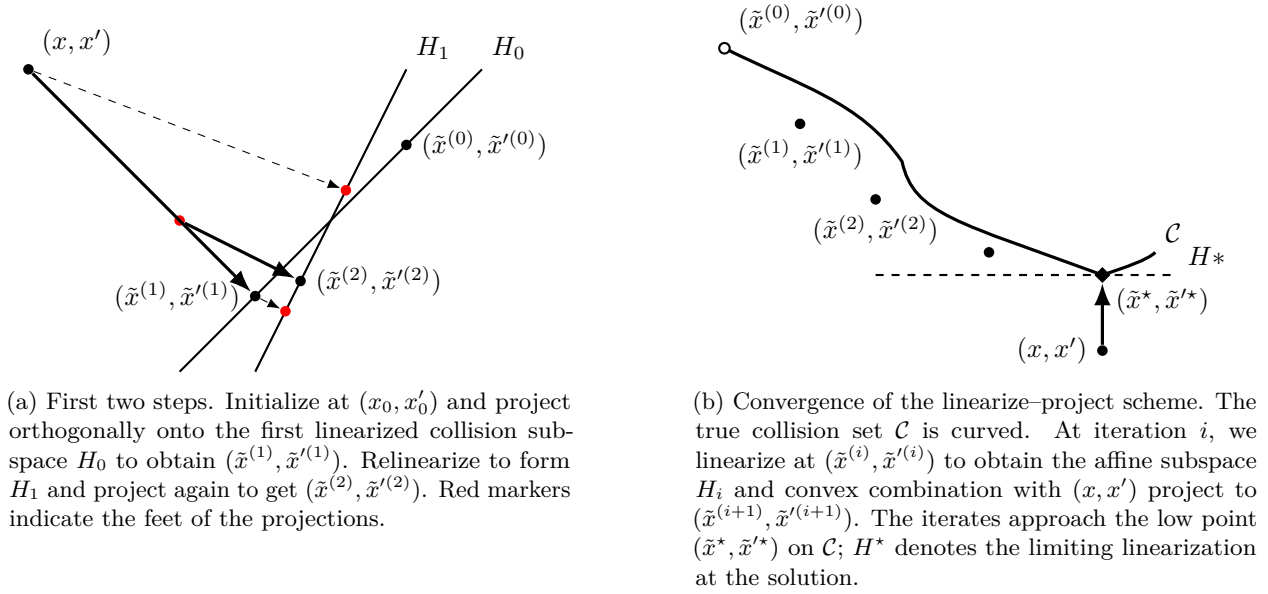


Figure 4: Geometry of our linearize-project algorithm on the collision manifold: (a) initialization and the first two projections; (b) iterative convergence toward the optimal point on the true collision set.

to introduce bias toward the original point x . This diminishing bias keeps every iterate within a controlled neighborhood of a feasible input, preventing large jumps that might violate constraints or break the validity of the local linear approximation. Figure 4 shows the convex-combination and projection operation in first two step as well as the convergent case in which the initial point x exactly projects to the current iteration point $\tilde{x}^{(i)}$.

The algorithm proposed in Bhagoji et al. (2023) is tailored to ReLU networks: it greedily searches ReLU activation patterns to formulate a sequence of linear programs. However, as network depth increases or when non-ReLU activations are used, this combinatorial approach becomes intractable. In contrast, our Jacobian-subspace method sidesteps these limitations by relying solely on gradient information. By constructing an affine subspace from the model’s local Jacobians, we can handle [differentiable](#) activation functions and deep architectures efficiently...

5.2 Computation and Complexity Analysis

Our algorithm alternates between (i) computing local Jacobians to linearize the collision constraint and (ii) projecting onto the resulting affine subspace.

Time Complexity. The computational cost to generate a collision pair scales as:

$$O(N \cdot T \cdot (m \cdot \mathcal{C}_{\text{bp}} + m^2 n_0))$$

where N is the number of random projections, T is the number of iterations, n_0 is the input dimension, and m is the projection dimension (default $m = 10$). Here, \mathcal{C}_{bp} represents the cost of a standard backpropagation pass through the feature extractor. The term $m \cdot \mathcal{C}_{\text{bp}}$ reflects the computation of the Jacobian (via m vector-Jacobian products), while $m^2 n_0$ accounts for solving the linear system for the projection. In practice, with small m , the runtime is dominated by the backpropagation steps.

Memory Complexity. Peak memory usage is determined by the storage of intermediate activations required for gradient computation. This scales linearly with $O(\text{Batch} \times \text{Depth} \times \text{Resolution})$. For our 7-layer CNN on CIFAR-10, we can process batches of 100 pairs on a 40GB GPU; however, for deeper architectures like WideResNet or higher-resolution inputs, the batch size must be reduced (e.g., to ~ 15) to remain within VRAM limits.

6 Experiments

In this section, we test several theories and evaluate the robustness of convolutional models layer by layer on the MNISTLeCun et al. (1998) and CIFAR-10Krizhevsky & Hinton (2009) datasets. Since MNIST is a relatively simple dataset, it allows us to easily explore robustness properties. And the MNIST dataset highlights a larger robustness gap between layers and provides a range of ϵ values for testing. For our baseline, we use an adversarially trained model on MNIST with projected gradient descent (PGD)Madry et al. (2017) and set perturbation budget $\epsilon = 2$.

To design the experiment, we selected 100 examples from two class in the training dataset for evaluation. For the robustness analysis, we focus on fixed binary class pairs to keep the setup simple and reproducible. On MNIST and Cifar-10 we use digits 3 vs. 7. For each dataset we randomly sample $S = 100$ training examples per class and evaluate our metric on *all* cross-class pairs, resulting in $S^2 = 10000$ pairs. Since our collision-search algorithm is run on every pair and for multiple random projections, the overall cost scales quadratically in S , and larger subsets quickly become computationally prohibitive. In Appendix Section C, we show that increasing the subset to $S = 200$ and changing the class pair or the sampled subset yields qualitatively similar layer-wise robustness curves, indicating that our conclusions are robust to these choices.

Algorithm 2 Optimal Loss Calculation via Conflict Graph

- 1: Choose feature extractor $f_j : \mathbb{R}^n \rightarrow \mathbb{R}^{n_j}$
 - 2: Specify the number of random projections N , and projection dimension m
 - 3: **for** $i = 1$ to N **do**
 - 4: Generate random projection matrix $A_i \in \mathbb{R}^{m \times n_j}$, with elements drawn i.i.d. from $\mathcal{N}(0, 1)$
 - 5: Build new feature extractor $(a_i \circ f_j)(x) = A_i f_j(x')$
 - 6: **for** example pairs (x, x') **do**
 - 7: Approximate $d_{a_i \circ f_j}(x, x')$ using Algorithm 1,
 - 8: **end for**
 - 9: **for** $\epsilon = \epsilon_{\min} \dots \epsilon_{\max}$ **do**
 - 10: Build conflict graph (V, E_ϵ) , where:
 - 11: V is the set of data, $x, x' \in V$
 - 12: $E_\epsilon = \{(x, x') : d_{a_i \circ f_j}(x, x') \leq \epsilon\}$
 - 13: Find $L_{\text{opt}}(\epsilon)$ using Algorithm B
 - 14: **end for**
 - 15: **end for**
 - 16: **Return** List of averaged optimal losses L_{opt}
-

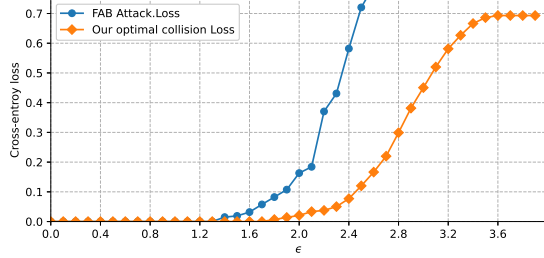
For the selection of CNN models, we trained three different architectures, following the work in Gowal et al. (2019), where all layers utilize ReLU activation functions. We employed two robust training methods: PGD adversarial training, and TRADES adversarial training Zhang et al. (2019). Each of these methods was applied with robustness constraints based on the l_2 norm. For all models, we used a cosine learning rate schedule Loshchilov & Hutter (2016) with an initial learning rate set to 0.01 to ensure consistency across training approaches.

Optimal Loss We use optimal loss $\tilde{L}^*(\epsilon, P, a \circ f)$ defined in Section 3.2 to represent the robustness of the feature extractor. Recall that higher loss represent lower robustness performance. Algorithm 2 computes the optimal loss using a conflict graph-based procedure. The process begins by selecting a feature extractor and generating a sequence of random matrices with i.i.d. Gaussian entries. Each random matrix A_i is combined with the feature extractor f_j to create a lower-dimensional feature mapping $x \mapsto A_i f_j(x)$, making the representation suitable for robustness analysis. The algorithm then evaluates all example pairs to estimate pairwise distances and identifies collisions according to a perturbation-based threshold ϵ . These collisions define the edges of a conflict graph, where nodes represent examples and edges connect pairs within the specified perturbation distance. For each adversarial budget the algorithm invokes Algorithm B (corresponding to Algorithm 1 from Bhagoji et al. (2021)) to compute the optimal loss over the conflict graph. Finally, the results are averaged over multiple random matrices to produce the final optimal loss estimates.

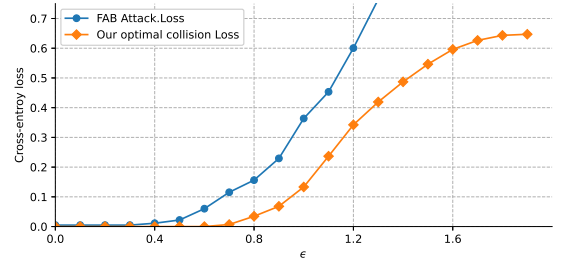
6.1 Validating the robustness metric

To evaluate the true robustness loss on both the training and test datasets, we additionally compare our collision-based loss to the FAB attack, which provides a standard decision-boundary-based robustness metric, rather than relying only on the training-time attack loss. FAB computes, for a single input, the minimum perturbation required to change its predicted label, whereas our collision-based procedure at the logit layer jointly moves a pair (x, x') from different classes to a point where their logits coincide.

In our setup, the worst-case loss is $\log(2)$, indicating that even when collisions occur, we still have a 50% chance of correctly classifying a randomly chosen element of the pair. As shown in the figure 5a, when the value of ϵ increases, the adversarial loss measured by FAB rises, while our loss saturates at its maximum value $\log(2)$. Thus, the FAB loss remains consistently higher than our collision-based loss, which is expected because FAB directly optimizes misclassification of a single input, whereas our metric is designed to capture how easily feature representations of different classes can be brought to coincide.



(a) MNIST Medium CNN (PGD-2).



(b) CIFAR-10 Medium CNN (PGD-0.5).

Figure 5: Comparison of FAB Attack Loss vs. Our Optimal Collision Loss.

6.2 Varying the dimension of the projected features

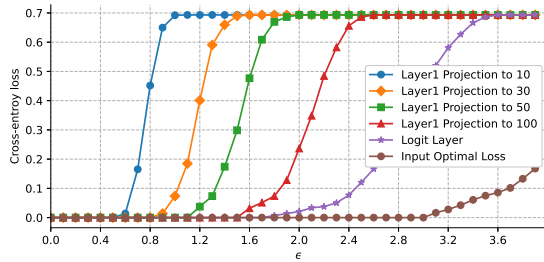


Figure 6: MNIST Medium CNN (PGD-2): First conv layer with varying random projection dimensions ($m \rightarrow 784$).

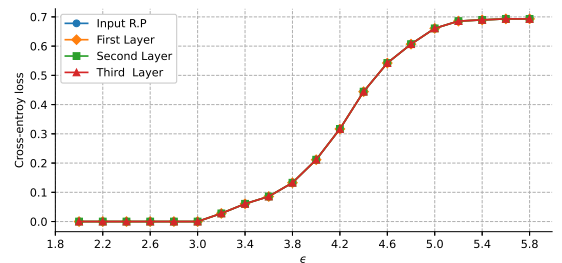


Figure 7: Medium CNN optimal loss using the method in Bhagoji et al. (2023).

Since the output dimension of the convolution is higher than the input, when no random projection is applied, the only collisions occur when the pairs are identical. Therefore, the distance between pairs reflects the exact input distance, as shown by the brown line in the figure 6. When projecting to a higher-dimensional space, we introduce less noise while preserving more useful feature information. As a result, the dimension increase, and the loss decreases.

The distances scale roughly with \sqrt{m} , but not exactly, which reflects the somewhat different information about feature structure accessed by different choices of m .

6.3 Baseline layer-by-layer robustness results

When applying the method of Bhagoji et al. (2023), the search algorithm cannot find any nontrivial collisions because none of the initial linear layers have an output dimension smaller than the example space dimension. The robustness measure for these initial layers remains essentially the same as at the input level. This is illustrated in 7. Layer dimensions for all models investigate are listed in Table ?? in Appendix ??.

In Figures 8 and 9, we have baseline layer-by-layer results for MNIST and CIFAR10 respectively. As we move through the layers, robustness increases, suggesting that low-quality feature are being pruned. This effect is far outweighing any reductions in robustness due to the introduction of new opportunities for collisions as the feature dimension spaces decrease. Clearly, the dimensions of the feature spaces have a significant impact on the robustness measure, but they do not explain the curves by themselves. For both datasets, there is a relatively large improvement from layer one to layer two and from layer four to layer five. These layers perform dimensionality reduction by a factors of four and six, which are larger than reductions in adjacent layers. However, from layer six to the logits, we do not observe much robustness improvement despite dimensionality reduction from 512 to 10.

As we move through the layers, we observe that the optimal loss decreases, indicating that the model gradually learns more about the data’s robustness, and the features become increasingly separable. Notably, there is a distinct gap between the convolutional layers and the fully connected layers. Additionally, the random projection appears to perform similarly to the learned matrix. The last few linear layers, however, seem to contribute little to the overall robustness.

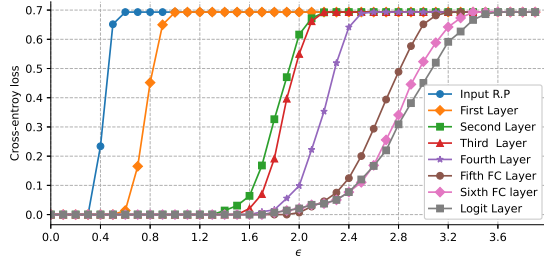


Figure 8: MNIST Medium CNN Model trained with PGD-2: All layers robustness.

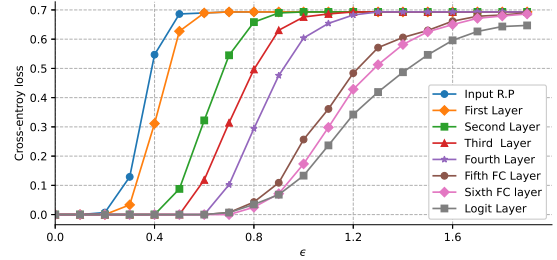


Figure 9: CIFAR-10 Medium CNN Model trained with PGD-0.5: All layer robustness.

6.4 Effects of architecture on layer-by-layer robustness

When comparing the layer-wise robustness of the small CNN in figure 10 and the Large CNN in figure 11 on MNIST, we observe that robustness does not change uniformly throughout the models. The same observation shows in robustness on CIFAR in figure 12. Although the small CNN exhibits lower overall robustness, its first convolutional layer appears to perform better, and its second layer achieves a similar level of robustness to that of the medium CNN.

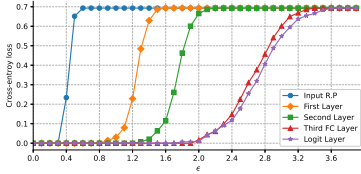


Figure 10: MNIST Small CNN Model trained with PGD-2: All layer robustness.

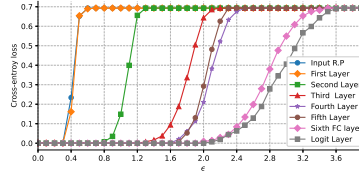


Figure 11: MNIST Large CNN Model trained with PGD-2: All layer robustness.

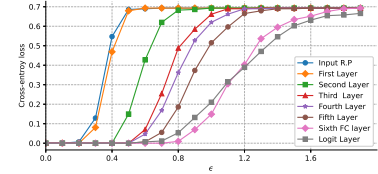


Figure 12: CIFAR-10 Large CNN Model trained with PGD-0.5: All layer robustness.

6.5 Effects of training method on layer-by-layer robustness

In Figures 13,14,15, we examine the medium CNN models trained with methods other than plain PGD. These can be compared to Figure 8. TRADES have layer-one robustness cliffs at lower adversarial budgets than PGD. Besides, TRADES achieves more separation between fully-connected layer curves and more gradual growth of those curves.

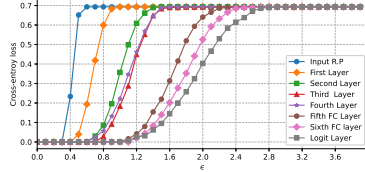


Figure 13: MNIST Medium CNN Model with benign training: All layer robustness.

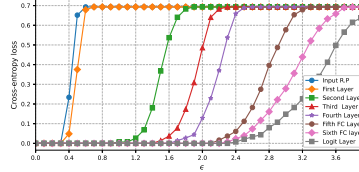


Figure 14: MNIST Medium CNN Model trained with TRADES-2: All layer robustness.

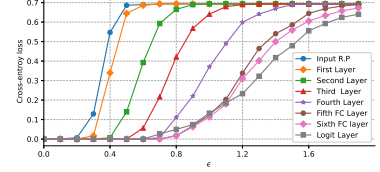


Figure 15: CIFAR-10 Medium CNN Model trained with TRADES-0.5: All layer robustness.

Relationship between intermediate layer robustness and test robust accuracy In Figures ?? and ??, we see the test robust accuracy of the models. As expected, benign training achieves very limited robustness, TRADES significantly outperforms plain PGD. It is unclear to us whether the layer-one behavior of TRADES should be viewed as an undesirable side-effect that suggests room for improvement, or an indication that layer-one robustness curves are not appropriate objectives.

6.6 Conclusion

An interesting finding is that, for models exhibiting greater robustness under the norm-based attack, the optimal loss at the logit layer remains relatively low, which indicates high overall robustness, yet the first convolutional layer does not consistently follow this trend. In fact, this layer often learns less robust features even when the model as a whole is quite robust. This suggests that prioritizing the robustness of the first convolutional layer early in the training process is a potential avenue for improving a model’s robustness.

At the same time, our experiments show that robustness does not grow uniformly across layers or blocks or depend in on feature space dimension a simple manner. Some adjacent layers have similar robustness, and robustness levels often stabilize before the final linear layer. Of course, this behavior varies with different datasets, training methods, and architectures.

Limitations and Future Work. Our method requires Jacobian computations for input pairs, imposing a memory overhead that scales linearly with input resolution and network depth. This trade-off necessitates smaller batch sizes for very deep or high-resolution models compared to standard training. Additionally, while our collision search is algorithmically compatible with various constraints, our current theoretical concentration bounds are derived for the rotationally invariant l_2 geometry. Developing rigorous theoretical guarantees for the l_∞ norm involves addressing non-trivial challenges regarding the geometry of intersection sets in high dimensions, which constitutes a key direction for our future work.

Broader Impact Statement Analysis of feature extractor robustness can be used to detect weaknesses in classifiers and localize these to particular layers. This is intended to serve the goal of building more robust and trustworthy models and is likely to be most useful for these purposes. For example, weak layers can be retrained or rearchitected. However, there is some potential for any method of assessing robustness to also be useful for the design and implementation of attacks: locating weaknesses in classifiers so they can be exploited rather than fixed. The potential for misuse is mitigated by the fact that our collision finding algorithm has been designed for use in the white-box setting.

References

- Arjun Nitin Bhagoji, Daniel Cullina, and Prateek Mittal. Lower bounds on adversarial robustness from optimal transport. In *Advances in Neural Information Processing Systems*, pp. 7496–7508, 2019.
- Arjun Nitin Bhagoji, Daniel Cullina, Vikash Sehwal, and Prateek Mittal. Lower bounds on cross-entropy loss in the presence of test-time adversaries. In *International Conference on Machine Learning*, pp. 863–873. PMLR, July 2021.
- Arjun Nitin Bhagoji, Daniel Cullina, and Ben Y Zhao. A theoretical perspective on the robustness of feature extractors. In *The Second Workshop on New Frontiers in Adversarial Machine Learning*, 2023.
- Nicolas Bonneel, Julien Rabin, Gabriel Peyré, and Hanspeter Pfister. Sliced and radon wasserstein barycenters of measures. *Journal of Mathematical Imaging and Vision*, 51(1):22–45, 2015.
- Jared Craig, Jayson Andle, Thomas S. Nowak, and Sahand Y. Sekeh. A theoretical perspective on subnetwork contributions to adversarial robustness. *arXiv preprint arXiv:2307.03803*, 2023.
- Francesco Croce and Matthias Hein. Minimally distorted adversarial examples with a fast adaptive boundary attack. In *International Conference on Machine Learning*, pp. 2196–2205. PMLR, November 2020.
- Sihui Dai, Wenxin Ding, Arjun Nitin Bhagoji, Daniel Cullina, Heather Zheng, Ben Zhao, and Prateek Mittal. Characterizing the optimal 0-1 loss for multi-class classification with a test-time attacker. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 49467–49489. Curran Associates, Inc., 2023.
- Ishan Deshpande, Ziyu Hu, Liyuan Sun, Nicholas Kolkin, Justin Solomon, Changil Kim, Sanja Fidler, Deva Ramanan, and Krishna Singh. Max-sliced wasserstein distance and its use for gans. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10648–10656, 2019.
- Paul Gavrikov, Shashank Agnihotri, Margret Keuper, and Janis Keuper. How do training methods influence the utilization of vision models? In *NeurIPS 2024 Workshop on Interpretable AI: Past, Present and Future*, 2024. URL <https://openreview.net/forum?id=zJFvjdw9JS>.
- Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *International Conference on Learning Representations*, 2015.
- Sven Gowal, Krishnamurthy (Dj) Dvijotham, Robert Stanforth, Jonathan Uesato, Relja Arandjelovic, Brendan O’Donoghue, and Pushmeet Kohli. Scalable verified training for provably robust image classification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 4842–4851, 2019.
- Guy Katz, Clark Barrett, David L. Dill, Kyle Julian, and Mykel J. Kochenderfer. Reluplex: An efficient smt solver for verifying deep neural networks. In *Computer Aided Verification (CAV)*, volume 10426 of *Lecture Notes in Computer Science*, pp. 97–117. Springer, 2017. doi: 10.1007/978-3-319-63387-9_5.
- Soheil Kolouri, Phillip E. Pope, Charles E. Martin, and Gustavo K. Rohde. Sliced wasserstein auto-encoders. In *International Conference on Learning Representations*, 2018.
- Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. Technical report, University of Toronto, 2009.
- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016.
- Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.

- Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, and Pascal Frossard. Deepfool: a simple and accurate method to fool deep neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2574–2582, 2016.
- Kimia Nadjahi, Ievgen Redko, Guillaume Bouchard, Marc Sebban, and Nicolas Courty. Fast approximation of sliced-wasserstein distances using concentration of measure. In *Advances in Neural Information Processing Systems*, volume 34, pp. 22238–22250, 2021.
- Muni Sreenivas Pydi and Varun Jog. Adversarial risk via optimal transport and optimal couplings. In *Proceedings of the 37th International Conference on Machine Learning*, pp. 7814–7823, 2020.
- Chulin Song, Rajeev Ranjan, and Haoliang Li. A layer-wise adversarial-aware quantization optimization for improving robustness. *arXiv preprint arXiv:2110.12308*, 2021.
- 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.
- Nicolás García Trillos, Matt Jacobs, and Jakwang Kim. The multimarginal optimal transport formulation of adversarial multiclass classification. *Journal of Machine Learning Research*, 24(45):1–56, 2023.
- Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric Xing, Laurent El Ghaoui, and Michael Jordan. Theoretically principled trade-off between robustness and accuracy. In *International Conference on Machine Learning*, pp. 7472–7482. PMLR, 2019.