

A Singular Value Perspective on Model Robustness

Malhar Jere
UC San Diego
mjjere@ucsd.edu

Maghav Kumar
UIUC
mkumar10@illinois.edu

Farinaz Koushanfar
UC San Diego
fkoushanfar@ucsd.edu

Abstract

Convolutional Neural Networks (CNNs) have made significant progress on several computer vision benchmarks, but are fraught with numerous non-human biases such as vulnerability to adversarial samples. Their lack of explainability makes identification and rectification of these biases difficult, and understanding their generalization behavior remains an open problem. In this work we explore the relationship between the generalization behavior of CNNs and the Singular Value Decomposition (SVD) of images. We show that naturally trained and adversarially robust CNNs exploit highly different features for the same dataset. We demonstrate that these features can be disentangled by SVD for ImageNet and CIFAR-10 trained networks. Finally, we propose **Rank Integrated Gradients (RIG)**, the first rank-based feature attribution method to understand the dependence of CNNs on image rank.

1. Introduction

Deep Neural Networks have made significant progress on several challenging tasks in computer vision such as image classification [22], object detection [34] and semantic segmentation [12, 29]. However, these networks have been shown to possess numerous non-human biases, such as high facial recognition misclassification error rates against certain races and genders [3], vulnerability to numerous classes of adversarial samples [44, 11, 15, 14, 19, 32], and vulnerability to training-time backdoor attacks [23].

A line of recent efforts focused on explaining the generalization behavior of neural networks through adversarial robustness has shown significant promise. Such methods involve characterizing network inputs based on robust and non-robust features [17, 8, 47], understanding their effects on feature maps [51], interpreting their frequency components [53, 25] and interpreting their principal component properties [18, 1]. Surprisingly, prior work has shown that neural nets often generalize to test sets based on superficial correlations in the training set [17, 47, 10, 46, 28].

In this work and inspired by previous works [20, 17, 47], we investigate the hypothesis that naturally trained CNNs leverage such superficial correlations in the dataset. How-

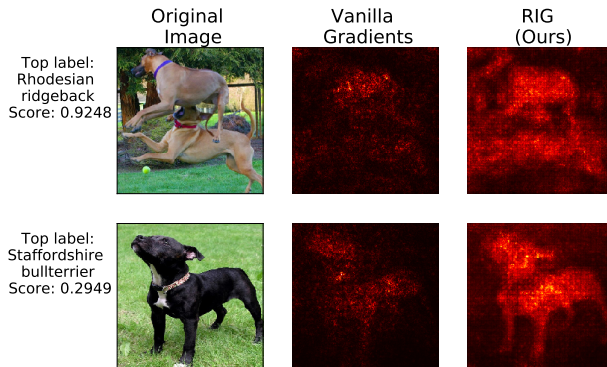


Figure 1: Left-to-right: The original image, Vanilla Gradients (as obtained by backpropagation with respect to the top label), and **Rank Integrated Gradients (RIG)**, our pixel importance method for the top label that averages saliency map information across low-rank representations of the same image. Notice that the visualizations obtained from RIG are better at identifying distinctive features of the image.

ever, different from prior works, we argue that these superficial correlations can be distilled from an image via low-rank image approximation, a claim that was previously refuted [47]. We further argue that naturally trained neural networks and adversarially robust neural networks exploit highly different features from the same image, and that these features can be separated by singular value decomposition (SVD). Our contributions are as follows:

- We identify for the first time that image rank (obtained from SVD) yields several novel insights about CNN robustness and interpretability (for example Figure 4). We provide arguments in favor of using image rank as a potential human-aligned image robustness metric.
- We show empirically that naturally trained CNNs place a large importance on human-imperceptible higher-rank components, and that adversarial retraining increases reliance on human-aligned lower-ranked components. Furthermore, we demonstrate that neural networks trained on imperceptible, higher-rank fea-

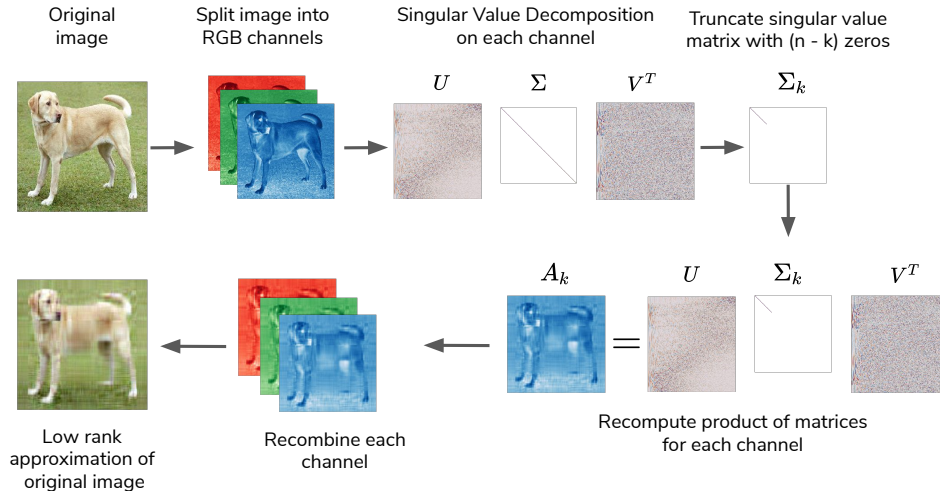


Figure 2: Generating a rank- k image via truncated SVD. Given an $n \times n$ RGB image, we decompose the image into its individual color channels, zero out the last $(n - k)$ singular values obtained via SVD, and then reconstruct the image with its k nonzero singular values. Low rank images are often more blurry than their full-rank counterparts.

tures generalize to the test set.

- We provide experimental evidence that neural networks trained on low-rank images are more adversarially robust than their naturally trained counterparts for the same dataset, and capture the accuracy-robustness tradeoff in CNNs in this new lens.
- We propose **Rank-Integrated Gradients (RIG)**, the first rank-based feature attribution method. Saliency maps generated by RIG highlight features more in line with human vision and offer a new way to interpret the decisions of CNNs (Figure 1).

Our work provides a new methodology to capture model robustness, and allows us to distinguish between naturally trained and robust models outside of the traditional L_p -norm robustness framework. It suggests that data approximation strategies such as low-rank approximation can be leveraged to improve out-of-distribution CNN performance, such as against adversarial samples. Finally, we show that saliency maps that incorporate rank information highlight more visually meaningful features. We hope our work will encourage researchers to include image approximation techniques when studying CNN generalization.

2. Background and Related Work

2.1. Notation

We consider a neural network $f(\cdot)$ used for classification where $f(x)_i$ represents the softmax probability that image

x corresponds to class i . Images are represented as $x \in [0, 1]^{w \times h \times c}$, where w, h, c are the width, height and number of channels of the image. We denote the classification of the network as $r(x) = \arg \max_i f(x)_i$, with $r^*(x)$ representing the ground truth of the image. Given an image x and an L_p norm bound ϵ , an adversarial sample $x' = x + \delta$ has the following properties:

- For a perturbation $\delta \in [0, 1]^{w \times h \times c}$ added to an image x such that $x' = x + \delta$, $L_p(\delta) = (\sum_{i=1}^h \sum_{j=1}^w |\delta_{i,j}|^p)^{1/p} \leq \epsilon$ where $p = (1, 2, \infty)$.
 - $p = 1$ is the Manhattan norm, defined as the sum of the absolute values of δ .
 - $p = 2$ is the Euclidean norm of δ .
 - $p = \infty$ is the infinity norm or max-norm of δ , defined as the largest absolute value in δ .
- $r(x') \neq r^*(x) = r(x)$. This means that the prediction on the adversarial sample is incorrect while the original prediction is correct.

2.2. Adversarial Samples

In this work we consider adversaries with white-box access to the neural network. In the white box threat model all information about the neural network is accessible. Using this information, adversaries can compute gradients with respect to inputs by backpropagation. White box attacks can be either targeted or untargeted. In targeted attacks, adversaries seek to generate an adversarial sample x' from an image x to force the neural network $f(x)$ to predict a pre-specified target t that is different from the true class

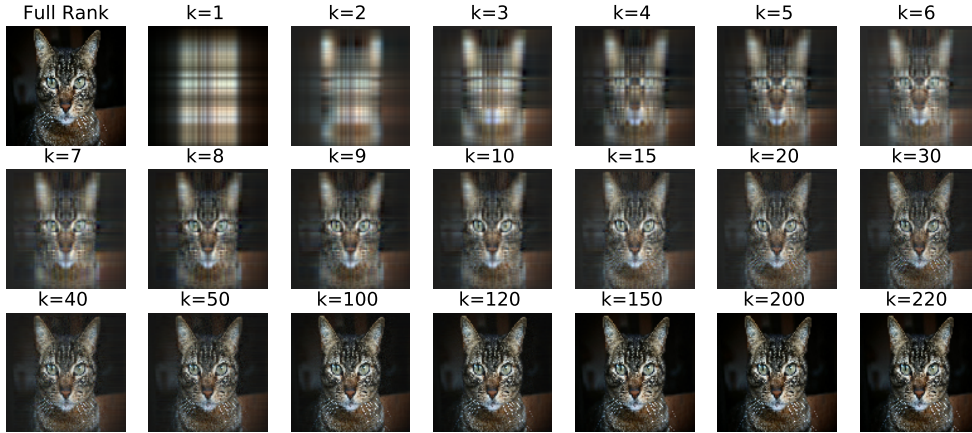


Figure 3: Low-rank approximations of the same image. Transitioning from low-rank approximations of images to higher-rank approximations yields better image quality.

$r^*(x)$, while in untargeted attacks adversaries seek to find an adversarial sample x' whose prediction is simply different from that of the true class.

Numerous methods to generate adversarial samples have been proposed [31, 11, 4, 30, 27]. In this work, we focus on the PGD attack with random starts, which has been shown to be an effective universal first-order adversary against neural networks [27]. For a neural network f , PGD is an iterative adversarial attack method that seeks to generate a targeted adversarial sample x' from an original image x with maximum perturbation limit ϵ . At each iteration, it performs a gradient descent step in the loss function w.r.t the image pixel values and the target class t and projects the perturbed image onto the feasible space, which is either a maximum per-pixel perturbation of ϵ (for L_∞ perturbations) or a maximum Euclidean perturbation distance from x of ϵ (for L_2 perturbations).

Adversarial training. Adversarial training defends against adversarial samples by training networks on adversarial perturbations that are generated on-the-fly. Adversarial training with L_∞ PGD samples has been shown to be among the most effective methods in mitigating these attacks [51, 27].

2.3. Explaining Adversarial Samples

Recent work has begun to understand the origin of adversarial samples. Ilyas et al. demonstrate that models trained on adversarial samples can generalize to test sets [17], and posit that adversarial samples are generalizable features that neural networks learn which are invisible to humans. Yin et al. and Wang et al. [53, 47] propose that adversarial samples for non-robust neural networks are in the high-frequency domain. Jere et al. [18] hypothesize that adversarial samples require significantly more principal components of an image to reach the same prediction compared to natural images. Our work is most similar to that of Yin et al. [53],

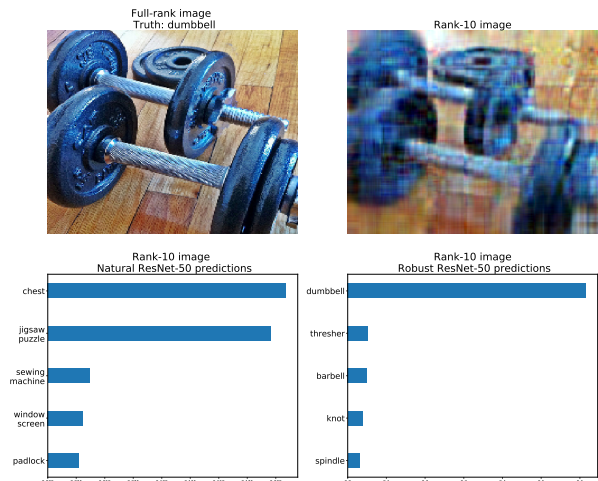


Figure 4: Accuracy for full-rank and rank-10 images. Top left is the full-rank image, top-right is its rank-10 approximation, bottom left are top-5 predictions for a naturally trained ResNet-50 on the rank-10 image, bottom right are top-5 predictions for an $L_\infty = 4/255$ adversarially robust ResNet-50 on the rank-10 image. Note that the adversarially robust version makes better predictions on rank-based distorted versions of the same image, even though they do not fall under an L_p -distortion framework.

in that we observe naturally trained CNNs are sensitive to higher-rank features, and that adversarial training makes them more biased to low-rank features. We explore the relationship between Fourier and low-rank features in the appendix.

2.4. Feature Attribution Methods

The problem of feature attribution seeks to *attribute the prediction of deep neural networks to its input features.*

Most methods of feature attribution involve variants of visualizing the gradients $\frac{\partial f}{\partial x}$ of the network with respect to the top predicted class i [37, 39, 2]. A significant challenge in designing attribution techniques is that attributions without respect to a fixed baseline are hard to evaluate, a problem which was successfully addressed by Integrated Gradients [43]. In Integrated Gradients, the baseline for an image is established with respect to a completely black image \tilde{x} , and a straight line is defined from \tilde{x} to x with increasing brightness values. Gradients are weighed and computed at each of these steps to result in the final saliency map, which is mathematically equivalent to the path integral along a straightline path from the baseline to the input. In our work, we perform a similar evaluation with our baseline set by the minimum *rank* of an image, and with gradients computed along increasing image rank. Further details can be found in Section 4.3, with numerous examples of RIG saliency maps in Figures 1, 8 and in the appendix.

2.5. Low-rank approximations

Low-rank representations of matrices (obtained via the Singular Value Decomposition) can capture a significant amount of information while simultaneously eliminating spurious correlations, and have recently been used in several Deep Learning applications, such as compression of CNN filters [7] and compression of internal representations of attention matrices in transformers [5, 48].

3. Methodology

In this section we first introduce the theoretical basis behind singular value decomposition, followed by the algorithm used in the rest of the paper to perform low-rank decomposition of RGB images.

3.1. Eigendecomposition of images

Eigendecomposition is commonly used to factor matrices into a canonical form, where the matrix is represented in terms of its eigenvalues and eigenvectors. In this work, we focus on utilizing the Singular Value Decomposition (SVD) to obtain low-rank approximations of an input to a neural network. In particular, let an image $x \in [0, 1]^{w \times h \times c}$, where w, h, c are the width, height and number of channels of the image respectively, and $w \leq h$. For each channel $m \in 1, 2, \dots, c$ the singular value decomposition on the matrix $A \in [0, 1]^{w \times h}$ yields:

$$A = U \Sigma V^T \quad (1)$$

U and V are orthonormal matrices, and Σ is a $w \times h$ ($w \leq h$) diagonal matrix with entries $(\sigma_1, \sigma_2, \dots, \sigma_w)$ denoting the singular values of A such that $\sigma_1 \geq \sigma_2, \dots \geq \sigma_w \geq 0$. According to the Eckart-Young-Mirsky theorem, the best

k rank approximation to the matrix A in the spectral norm $\|\cdot\|_2$ is given by:

$$A_k = \sum_{j=1}^k \sigma_j u_j v_j^T \quad (2)$$

where u_j, v_j denote the j th column of U and V respectively. This process is also termed as *truncated SVD*. For an n -rank matrix A , its k -rank approximation can also be expressed as $A_k = U \Sigma_k V^T$, where Σ_k is constructed from Σ by setting the smallest $n - k$ diagonal entries set to zero.

3.2. Algorithm

We define the top class $r(x) = \arg \max_i f(x)_i$ as the predicted class on the full-rank image x . For an image $x \in [0, 1]^{w \times h \times c}$, we perform singular value decomposition for each color channel c , reconstruct a low-rank approximation using k singular values, and perform inference on the rank- k image. Algorithm 1 highlights the steps that occur, and Figure 2 illustrates the steps involved in generating a rank- k image by truncating each of the RGB channels and reconstructing the image. Experimental results regarding latency and runtime can be found in the appendix.

Algorithm 1: Finding the accuracies of a model $f(\cdot)$ for a batch of images $x_{1:N}$, where $x_i \in [0, 1]^{w \times h \times c}$ as a function of the input rank.

Result: Accuracies for each rank $k = 0 : w$

```

full_rank_preds ← f(x1:N)
rank_k_acc ← zeros(w + 1)
for k = 0 : w do
    rank_k_x = zeros_like(x1:N)
    for i = 1 : N do
        for channel = 1 : c do
            u, σ, v = SVD(x[i][channel])
            σ[k : w] = 0
            rank_k_x[i][channel] = u diag(σ) v
        end
    end
    rank_k_acc[k] = (f(rank_k_x) ==
                    full_rank_preds)
end
return rank_k_acc

```

4. Towards robustness metrics beyond L_p distortions

In this section we highlight several limitations of L_p distortions, followed by experimental results for naturally trained and adversarially robust CNNs. Finally, we introduce Rank-Integrated Gradients.

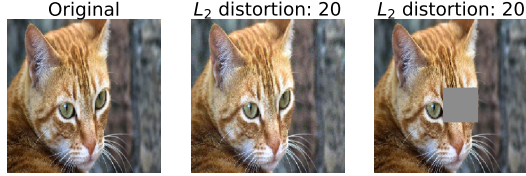


Figure 5: Limitations of L_p distortions. The middle image is generated by adding uniform noise to the first one. Both distorted images have identical L_2 distortions, but are perceptually dissimilar. The patched image (far right) represents one of a near infinite number of possible images with identical L_2 distortions as the middle image.

Dataset	Naturally trained Models	Robust models
ImageNet	ResNet-50, VGG-19, DenseNet-201	ResNet-50 ($L_\infty = 4/255, 8/255$), ($L_2 = 3.0$)
CIFAR-10	ResNet-50	ResNet-50 ($L_\infty = 8/255$), ($L_2 = 0.25, 0.5, 1.0$)

Table 1: Experimental setup to investigate behavior of benign v/s adversarially robust neural networks.

4.1. Limitations of L_p distortions

Extensive experiments have been conducted to secure neural networks against L_p -norm bounded perturbations, such as adversarial training [27, 21, 51, 38] and certified defenses against adversarial samples [33, 50, 54, 6]. Unfortunately, L_p -distortions represent a small fraction of potential image modifications. An infinite number of modified images exist that possess identical norm-bounded perturbations with respect to a base image. Furthermore, identical L_p norm-bounded distortions may be extremely different perceptually, pointing to L_p -norm robustness potentially being misaligned with human perception (Figure 5). Based on these observations and limitations, we argue that image rank and rank-based robustness metrics might be better suited to capture image modifications than L_p -distortions for the following reasons.

- Matrix rank for an image $x \in [0, 1]^{w \times h \times c}$ ($w \leq h$) is restricted to the set of integers $1, 2, \dots, w$. The set of images generated by rank truncation is a bijective mapping from rank k to generated images x_k . This is in contrast to images generated with L_p distortions whose mapping contains infinite possibilities.
- Low-rank image approximation effectively captures a much larger range of image modifications (Figure 3) that is more perceptually aligned with human vision

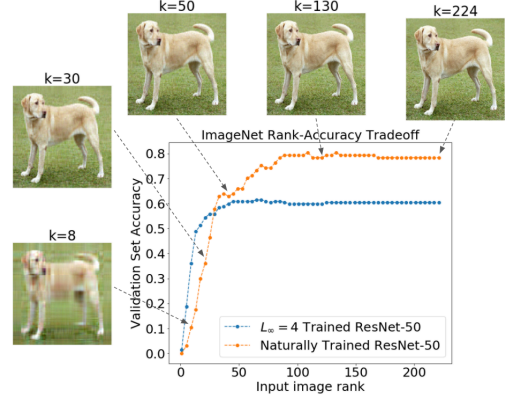


Figure 6: Dependence of test accuracy of naturally trained and $L_\infty = 4/255$ robust ResNet-50 models on input rank of natural images (subset of ImageNet validation images).

that might not be captured with L_p distortions. Transitioning from low-rank approximations of images (unrecognizable to humans) to higher-rank approximations (recognizable to both DNNs and humans) effectively allows us to better understand the gap between human and computer vision.

4.2. Rank Dependence of CNNs

We seek to understand the dependence of image rank and classifier accuracy for ResNet-50. We observe that naturally trained and robust CNNs use features that are highly different in their rank properties. Particularly, naturally trained CNNs use high-rank features that are often invisible to humans, while robust CNNs do not respond to these features but rather rely on highly visible low-rank features. For both models we randomly sampled 1000 images from the ImageNet validation set and cropped and reshaped each image to have a shape of $(224 \times 224 \times 3)$. For each image, we performed low-rank approximation for every possible rank prior to inference according to Algorithm 1.

4.2.1 Behavior of Naturally Trained CNNs

We investigate the behavior of the ImageNet-trained ResNet-50 [13], VGG-19 [41] and DenseNet-201 [16] (Table 1) CNN architectures trained on the ImageNet dataset [22]. Experimental results for the VGG-19 and DenseNet models can be found in the Appendix. In Figure 6 we observe the top-1 accuracy for ResNet-50 (orange) on these truncated images for both naturally trained and robust models. We make the following notable observations for this as well as for Figures 7a and 7b:

- Classifier accuracy sharply increases for lower-ranked images (rank-50 to rank-100) followed by saturation around rank-100 for ImageNet trained models. We observe similar behavior at rank-15 for CIFAR-10 trained models (Figure 7a).

- We observe that the features corresponding to this increase in accuracy, namely rank-50 to rank-100, contribute no meaningful semantic content to the image (Figure 3), indicating that naturally trained CNNs exploit features that are often invisible to humans. (7b)

4.2.2 Behavior of Adversarially Robust CNNs

We investigate the behavior of pretrained adversarially robust CNNs from the robustness library [9] as highlighted in column 2 of Table 1. Our results for an $L_\infty = 4/255$ robust ResNet-50 model can be seen in Figure 6 (blue). In contrast to naturally trained models we observe very different behaviors for robust models. Notably, we observe that:

- Robust CNN accuracy for full-rank images is lower than that of naturally trained CNNs for both ImageNet and CIFAR-10 trained models (which has been observed in previous work [45]).
- Robust CNN accuracy for lower-rank images is higher than that of naturally trained CNNs. Surprisingly, it is significantly superior for lower-ranked images, with a $> 20\%$ validation set accuracy improvement in both datasets (Figure 7c and 7d).
- Robust CNN accuracy increases much more quickly than naturally trained CNNs for lower-rank images, and does not exhibit the same dependence on features between rank-50 and rank-100 for the ImageNet dataset.

The rank-accuracy tradeoff as well as superior performance of robust models for lower-rank images has not been observed before, and to the best of our knowledge ours is the first work to identify such phenomena. We further observe in Figures 7a and 7b that this rank behavior persists across CNNs trained with different L_∞ bounds, different L_p -norm metrics and different datasets. While there exist minor differences between the rank behavior of L_∞ and L_2 robust CNNs, their behaviors are largely distinct from those of naturally trained CNNs.

4.3. Rank Integrated Gradients (RIG)

Based on these observations, we seek to visualize the rank-dependency of CNNs. Generating visual explanations for CNN image classifiers typically involves computing saliency maps that take the gradient of the output corresponding to the correct class with respect to a given input vector such as GradCAM and guided-backprop [40, 37, 42]. However, such methods often only capture local explanations for a given image, and are not robust to perturbations to the original image. Other methods involve training simpler, more interpretable surrogate models [24, 35] to understand model predictions in a local neighborhood around a

given input, but these cannot sufficiently capture rank-based image modifications nor scale to models such as ResNets.

Feature Attribution methods that seek to capture network predictions while being invariant to perturbations and implementations of the method [43] have been successful in capturing such attributes among image pixels. In particular, Blur Integrated Gradients (BIG) [52] has been effective in capturing feature attributes through Gaussian-blurred versions of the original image. In this regard our work is most similar to BIG, but we differentiate ourselves in calculating gradients through low-rank representations of an image rather than gaussian-blurred representations.

Intuitively, our technique weighs low-rank representations with their contributions to the gradients of an input for the top class predicted on the full-rank input. Formally, let us denote a classifier f and an input signal $x \in [0, 1]^{w \times h \times c}$ where $w \leq h$, with x_k as a rank- k image obtained by Algorithm 1. Let us denote an image classifier f , and the top predicted class i for a full rank image x_w . Let $(\frac{\partial f(x_k)}{\partial x_k})_i$ denote the maximum gradient across all color channels c . Then, our method computes RIG as:

$$RIG(x, f, i) = \sum_{k=1}^w \frac{w-k}{w} \times (\frac{\partial f(x_k)}{\partial x_k})_i \quad (3)$$

RIG requires no modification to the model and is extremely easy to implement, requiring less than 10 lines of PyTorch code and using a few calls to the gradient operation, thereby allowing even novice practitioners to easily apply the technique. Examples can be found in Figure 1 and 8.

5. Transferability of Rank-Based Features

Motivated by the disparities in the behavior between naturally trained and adversarially robust CNNs in Section 4, we proceeded to test the following hypotheses:

- Do CNNs trained solely on high-rank representations generalize to a full-rank test set?
- Do CNNs trained solely on low-rank representations generalize to a full-rank test set?
- Do CNNs trained solely on low-rank representations improve robustness to L_p -norm bounded attacks?

5.1. Training on solely high-rank representations

We conducted experiments to test the hypothesis: *Do CNNs trained on solely high-rank representations generalize to the test set?* We modified Algorithm 1 to zero out the k -largest singular values (instead of the k -smallest singular values), thereby creating images that consist solely of higher-rank features that are largely imperceptible and difficult to interpret even when visualized (Figure 9). We trained the ResNet-50 architecture on a modified version of the

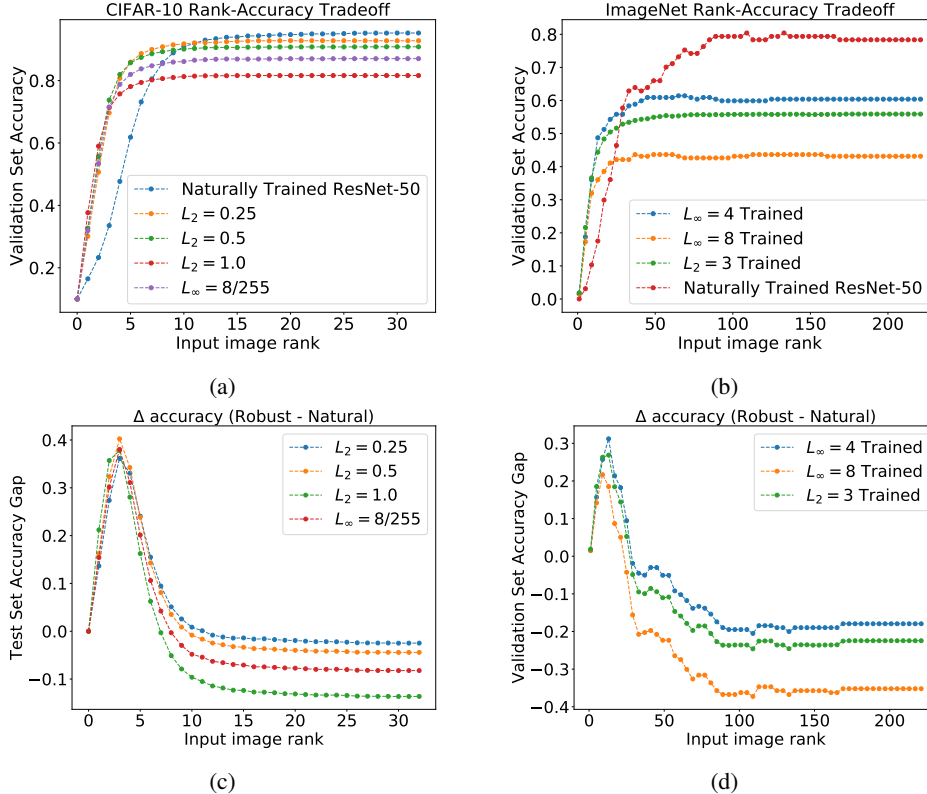


Figure 7: (a) CIFAR-10 rank spectrum for naturally trained and robust ResNet-50 models (b) ImageNet rank spectrum for naturally trained and adversarially robust ResNet-50 models. (c) Gap in test accuracy between robust and naturally trained CIFAR-10 ResNet-50 models. (d) Gap in test accuracy between robust and naturally trained ImageNet ResNet-50 models.

	CIFAR-10			ImageNet			
Trained on	Attack success rate	Recovery rate	Top-1 accuracy	Trained on	Attack success rate	Recovery rate	Top-1 accuracy
Full rank	99.93%	0.03%	95.21%	Full rank	95.87%	0.01%	78.35%
20	99.70%	0.15%	95.41%	100	81.81%	7.07%	73.99%
10	99.43%	0.19%	94.90%	50	73.99%	8.57%	70.16%
5	99.27%	0.34%	91.54%	30	73.19%	5.15%	69.07%

Table 2: Robustness of low-rank CIFAR-10 and ImageNet-trained ResNet-50 models. Attack success rate and recovery rate measured for targeted 20-step PGD attacks with L_∞ -bounds of 4/255. Top-1 accuracy is measured for full-rank test sets.

CIFAR-10 dataset consisting solely of these higher-ranked representations, and evaluated it on the full-rank CIFAR-10 test sets. Each network was trained for 350 epochs with an SGD optimizer, with learning rate of 0.1, momentum 0.9 and weight decay of 0.0005. We decreased the learning rate by 10 after the 150th and 250th epochs.

We observe that for training images where the first 3 singular values are truncated (Figure 9), we get a full-rank test accuracy of 28.02%. Such non-random accuracy shows that high-rank representations of an image contain meaningful features for generalization. However, this accuracy quickly

decreases as we increase the number of truncated singular values. Further details can be found in the appendix.

5.2. Training on low-rank representations

We conducted experiments to address the hypothesis: *Do CNNs trained solely on low-rank representations generalize to the test set?* We trained ResNet-50 models for CIFAR-10 and ImageNet solely on low-rank representations and observed their accuracy on the held out full-rank test sets.

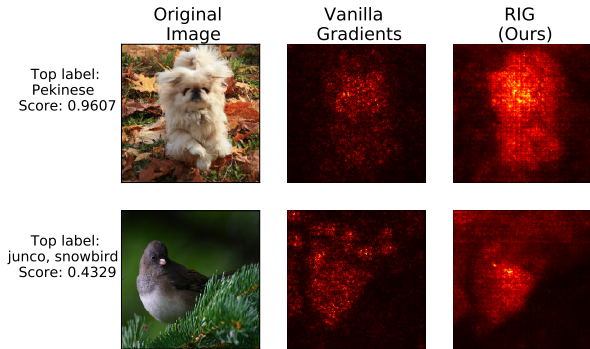


Figure 8: Rank Integrated Gradient Images. In the second row, vanilla gradients highlight features in the background that contribute to the top class but have no meaningful semantic content. RIG highlights these features as well as object specific features, such as the bird’s beak.



Figure 9: Reverse rank truncated CIFAR-10 image. Middle column corresponds to images where the 3 largest singular values were set to 0, with the difference between the original image and truncated image on the right column.

5.2.1 CIFAR-10

We trained ResNet-50 models on low-rank representations with similar hyperparameters as Section 5.1. We observe that low-rank representations are sufficient to achieve test accuracy of more than 90% for full-rank CIFAR-10 test sets. Surprisingly, training CIFAR-10 on rank-5 images yields test accuracy of 91.54%, which is only 3.67% lower than when trained on a full-rank dataset (Table 2). Increasing the rank of training-set images quickly closes this gap, and rank-20 images have almost identical test accuracy to full-rank images, indicating that the high-rank information in images is largely irrelevant to make predictions for CIFAR-10. This corroborates the results we obtained from Figure 7a, where we observed that test accuracy does not rely on features past rank-15 for CIFAR-10. Further experimental details can be found in the Appendix.

5.2.2 ImageNet

We trained ResNet-50 using the same hyperparameters as described in the original ResNet-50 paper [13] on low-rank versions of the ImageNet dataset. Specifically, we trained on rank-30, 50 and 100 representations (Table 2). Despite rank-100 and full-rank images being visually identical, we observe that the full-rank validation set accuracy for rank-

100 trained ResNet-50 is 4.3% lower than that of ResNet-50 trained on full-rank images. This indicates that the ImageNet data consists of a large number of imperceptible, high-rank features that do not contain semantically meaningful content but contribute to test accuracy.

5.3. Robustness as an emergent property of low-rank representations

In this section we conducted experiments to test the hypothesis: *Do CNNs trained solely on low-rank representations improve adversarial robustness to L_p -norm bounded attacks?* To tackle this, we performed 20-step, $L_\infty = 4/255$ PGD adversarial attacks on the low-rank CIFAR-10 and ImageNet-trained ResNet-50 models from Section 5.2. Our experimental results can be found in Table 2. Notably, we observe that adversarial robustness to L_∞ attacks improves with training on lower-ranked image representations for ImageNet. However, this does not hold true for CIFAR-10 trained models. Furthermore, strategies such as adversarial training [26] or feature denoising [51] offer superior performance than training on low-rank representations.

6. Discussion

Prior work on interpreting adversarial samples [17, 49] hypothesized that images consist of *robust* and *non-robust* features, where robust features are largely visible to humans while non-robust features are not. Further work argues that robustness leads to improved feature representations [36]. Our findings appear to support these claims. Specifically, we observe that due to their large contribution to image quality and predictive performance for robust networks, low-rank features are analogous to *robust* features and can be generated through low-rank truncation. Conversely, higher-ranked features which do not contribute to robust network predictions are analogous to *non-robust* features. Furthermore, we observe that quantifying network robustness through L_p perturbations does little to capture the massive range of possible image modifications, and often runs into the issue of multiple perceptually different images having identical L_p distortions. Rank-based image modifications simultaneously capture a much larger range of image modifications while offering a 1–1 mapping from modification parameter to perceptual representation.

With respect to feature attribution, we observe that saliency maps that leverage rank information in images are much more aligned with human vision than conventional vanilla gradients, and offer a new lens into understanding the inner workings of these image classifiers. We hypothesize that there exist several other similar forms of matrix decomposition that allow for visualizations that are more perceptually meaningful as well.

7. Conclusion

Closing the gap between computer and human vision is a challenging and an open problem. Human vision remains robust under a variety of image transformations, while neural network based computer vision is still fragile to small L_p -norm limited perturbations, which furthermore do not capture the full range of image modifications. We demonstrate the need for robustness metrics beyond these perturbations, and make several arguments in favor of using image-rank (as obtained by SVD) as a potential alternative. We demonstrate several behavioral differences between naturally trained and adversarially robust CNN classifiers in terms of their generalization that could not be captured in an L_p -bound framework. Finally, we propose a simple rank-based feature attribution technique that produces gradient visualizations that are much more perceptually informative than saliency maps.

8. Acknowledgements

This work was supported by the Semiconductor Research Corporation (AUTO TASK 2899.001) and a Defense Advanced Research Projects Agency (DARPA) Techniques for Machine Vision Disruption (TMVD) grant.

References

- [1] Arjun Nitin Bhagoji, Daniel Cullina, Chawin Sitawarin, and Prateek Mittal. Enhancing robustness of machine learning systems via data transformations. In *2018 52nd Annual Conference on Information Sciences and Systems (CISS)*, pages 1–5. IEEE, 2018. 1
- [2] Alexander Binder, Grégoire Montavon, Sebastian Lapuschkin, Klaus-Robert Müller, and Wojciech Samek. Layer-wise relevance propagation for neural networks with local renormalization layers. In *International Conference on Artificial Neural Networks*, pages 63–71. Springer, 2016. 4
- [3] Joy Buolamwini and Timnit Gebru. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Conference on fairness, accountability and transparency*, pages 77–91, 2018. 1
- [4] Nicholas Carlini and David Wagner. Adversarial examples are not easily detected: Bypassing ten detection methods. In *Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security*, pages 3–14. ACM, 2017. 3
- [5] Krzysztof Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane, Tamas Sarlos, Peter Hawkins, Jared Davis, Afroz Mohiuddin, Lukasz Kaiser, et al. Rethinking attention with performers. *arXiv preprint arXiv:2009.14794*, 2020. 4
- [6] Jeremy M Cohen, Elan Rosenfeld, and J Zico Kolter. Certified adversarial robustness via randomized smoothing. *arXiv preprint arXiv:1902.02918*, 2019. 5
- [7] Emily L Denton, Wojciech Zaremba, Joan Bruna, Yann LeCun, and Rob Fergus. Exploiting linear structure within convolutional networks for efficient evaluation. *Advances in neural information processing systems*, 27:1269–1277, 2014. 4
- [8] Logan Engstrom, Andrew Ilyas, Aleksander Madry, Shibani Santurkar, Brandon Tran, and Dimitris Tsipras. A discussion of ‘adversarial examples are not bugs, they are features’: Discussion and author responses. *Distill*, 4(8):e00019–7, 2019. 1
- [9] Logan Engstrom, Andrew Ilyas, Hadi Salman, Shibani Santurkar, and Dimitris Tsipras. Robustness (python library), 2019. 6, 15
- [10] Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A Wichmann, and Wieland Brendel. Imagenet-trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness. *arXiv preprint arXiv:1811.12231*, 2018. 1
- [11] Ian Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *International Conference on Learning Representations*, 2015. 1, 3
- [12] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017. 1
- [13] 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*, pages 770–778, 2016. 5, 8
- [14] Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adversarial examples. *arXiv preprint arXiv:1907.07174*, 2019. 1
- [15] Hossein Hosseini and Radha Poovendran. Semantic adversarial examples. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 1614–1619, 2018. 1
- [16] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4700–4708, 2017. 5, 12
- [17] Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Logan Engstrom, Brandon Tran, and Aleksander Madry. Adversarial examples are not bugs, they are features. In *Advances in Neural Information Processing Systems*, pages 125–136, 2019. 1, 3, 8
- [18] Malhar Jere, Sandro Herbig, Christine Lind, and Farinaz Koushanfar. Principal component properties of adversarial samples. *arXiv preprint arXiv:1912.03406*, 2019. 1, 3
- [19] Malhar Jere, Briland Hitaj, Gabriela Ciocarlie, and Farinaz Koushanfar. Scratch that! an evolution-based adversarial attack against neural networks. *arXiv preprint arXiv:1912.02316*, 2019. 1
- [20] Jason Jo and Yoshua Bengio. Measuring the tendency of cnns to learn surface statistical regularities. *arXiv preprint arXiv:1711.11561*, 2017. 1
- [21] Harini Kannan, Alexey Kurakin, and Ian Goodfellow. Adversarial logit pairing. *arXiv preprint arXiv:1803.06373*, 2018. 5

- [22] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 1097–1105. Curran Associates, Inc., 2012. 1, 5
- [23] Yingqi Liu, Shiqing Ma, Yousra Aafer, Wen-Chuan Lee, Juan Zhai, Weihang Wang, and Xiangyu Zhang. Trojancing attack on neural networks. In *25nd Annual Network and Distributed System Security Symposium, NDSS 2018, San Diego, California, USA, February 18-221, 2018*. The Internet Society, 2018. 1
- [24] Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In *Advances in neural information processing systems*, pages 4765–4774, 2017. 6
- [25] Tao Luo, Zheng Ma, Zhi-Qin John Xu, and Yaoyu Zhang. Theory of the frequency principle for general deep neural networks. *CoRR*, abs/1906.09235, 2019. 1
- [26] 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. 8
- [27] 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. 3, 5
- [28] Taro Makino, Stanisław Jastrzębski, Witold Oleszkiewicz, Celin Chacko, Robin Ehrenpreis, Naziya Samreen, Chloe Chhor, Eric Kim, Jiyon Lee, Kristine Pysarenko, Beatrice Reig, Hildegard Toth, Divya Awal, Linda Du, Alice Kim, James Park, Daniel K. Sodickson, Laura Heacock, Linda Moy, Kyunghyun Cho, and Krzysztof J. Geras. Differences between human and machine perception in medical diagnosis. *arXiv:2011.14036*, 2020. 1
- [29] Safa Messaoud, Maghav Kumar, and Alexander G Schwing. Can we learn heuristics for graphical model inference using reinforcement learning? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 766–767, 2020. 1
- [30] S. Moosavi-Dezfooli, A. Fawzi, O. Fawzi, and P. Frossard. Universal adversarial perturbations. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 86–94, July 2017. 3
- [31] 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*, pages 2574–2582, 2016. 3
- [32] Paarth Neekhara, Shehzeen Hussain, Malhar Jere, Farinaz Koushanfar, and Julian McAuley. Adversarial deepfakes: Evaluating vulnerability of deepfake detectors to adversarial examples. *arXiv preprint arXiv:2002.12749*, 2020. 1
- [33] Aditi Raghunathan, Jacob Steinhardt, and Percy Liang. Certified defenses against adversarial examples. *arXiv preprint arXiv:1801.09344*, 2018. 5
- [34] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 779–788, 2016. 1
- [35] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "why should i trust you?" explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1135–1144, 2016. 6
- [36] Hadi Salman, Andrew Ilyas, Logan Engstrom, Ashish Kapoor, and Aleksander Madry. Do adversarially robust imagenet models transfer better? *Advances in Neural Information Processing Systems*, 33, 2020. 8
- [37] 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*, pages 618–626, 2017. 4, 6
- [38] Ali Shafahi, Mahyar Najibi, Mohammad Amin Ghiasi, Zheng Xu, John Dickerson, Christoph Studer, Larry S Davis, Gavin Taylor, and Tom Goldstein. Adversarial training for free! In *Advances in Neural Information Processing Systems*, pages 3353–3364, 2019. 5
- [39] Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. Learning important features through propagating activation differences. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 3145–3153. JMLR. org, 2017. 4
- [40] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034*, 2013. 6
- [41] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations*, 2015. 5, 12
- [42] Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin Riedmiller. Striving for simplicity: The all convolutional net. *arXiv preprint arXiv:1412.6806*, 2014. 6
- [43] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 3319–3328. JMLR. org, 2017. 4, 6
- [44] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In *International Conference on Learning Representations*, 2014. 1
- [45] Dimitris Tsipras, Shibani Santurkar, Logan Engstrom, Alexander Turner, and Aleksander Madry. Robustness may be at odds with accuracy. *arXiv preprint arXiv:1805.12152*, 2018. 6, 12
- [46] Haohan Wang, Zexue He, Zachary C. Lipton, and Eric P. Xing. Learning robust representations by projecting superficial statistics out. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019. 1

- [47] Haohan Wang, Xindi Wu, Zeyi Huang, and Eric P Xing. High-frequency component helps explain the generalization of convolutional neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8684–8694, 2020. 1, 3
- [48] Sinong Wang, Belinda Li, Madian Khabsa, Han Fang, and Hao Ma. Linformer: Self-attention with linear complexity. *arXiv preprint arXiv:2006.04768*, 2020. 4
- [49] Zifan Wang, Yilin Yang, Ankit Shrivastava, Varun Rawal, and Zihao Ding. Towards frequency-based explanation for robust cnn. *arXiv preprint arXiv:2005.03141*, 2020. 8
- [50] Eric Wong and J Zico Kolter. Provable defenses against adversarial examples via the convex outer adversarial polytope. *arXiv preprint arXiv:1711.00851*, 2017. 5
- [51] Cihang Xie, Yuxin Wu, Laurens van der Maaten, Alan L Yuille, and Kaiming He. Feature denoising for improving adversarial robustness. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 501–509, 2019. 1, 3, 5, 8
- [52] Shawn Xu, Subhashini Venugopalan, and Mukund Sundarajan. Attribution in scale and space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9680–9689, 2020. 6
- [53] Dong Yin, Raphael Gontijo Lopes, Jon Shlens, Ekin Dogus Cubuk, and Justin Gilmer. A fourier perspective on model robustness in computer vision. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 13255–13265. Curran Associates, Inc., 2019. 1, 3
- [54] Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric P Xing, Laurent El Ghaoui, and Michael I Jordan. Theoretically principled trade-off between robustness and accuracy. *arXiv preprint arXiv:1901.08573*, 2019. 5

9. Appendix

9.1. Relationship between Fourier Features and Rank-based features

Statement: We hypothesize that the rank of a matrix obtained from a $(k \times k)$ low-pass filter in the frequency domain is upper bounded by k .

Proof. Let a matrix $X \in [0, 1]^{w \times h} = U\Sigma V^T$ with rank $r(X) = \min(w, h)$ exist in the spatial domain, and let the Fourier transform operation be represented as $F(\cdot)$. Due to its linearity, the Fourier transform of X can be expressed as $F(X) = WX$. Let us denote the low-pass filtering operation with a window of size k as L . By definition, the window will have a max rank of k . Then, we can express the k -window low-pass filtered version of X as $\tilde{Y} = LWX$ in the frequency domain, and its spatial domain representation as $\tilde{X} = W^{-1}LWX$.

The rank of \tilde{X} can be expressed as $r(\tilde{X}) = r(W^{-1}LWX)$. Due to the rank property of matrix multiplication, $r(W^{-1}LWX) \leq \min(r(W^{-1}), r(LWX))$. Therefore,

$$\begin{aligned} r(W^{-1}LWX) &\leq \min(r(W^{-1}), r(LWX)) \\ &\implies r(LWX) \leq \min(r(X), r(LW)) \\ &\implies r(LW) \leq \min(r(L), r(W)) = r(L) = k \\ &\implies r(W^{-1}LWX) \leq r(LWX) \leq r(LW) \leq r(L) \leq k \end{aligned} \tag{4}$$

Thus, $r(\tilde{X}) \leq k$.

9.2. Rank Integrated Gradients

We provide several more examples of RIG saliency maps for robust and non-robust ResNet-50 models here (Figures 10, 11, 17, 18). RIG highlights rank-based features which are more perceptually-aligned than vanilla gradients for naturally trained as well as adversarially robust networks.

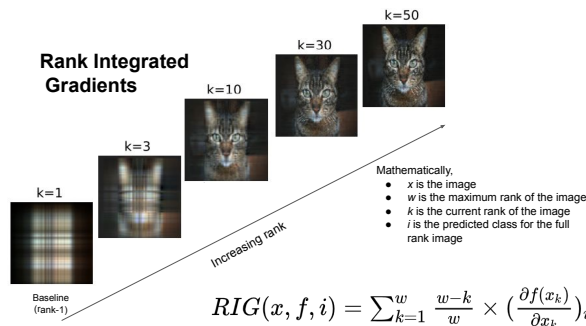


Figure 10: RIG generation.

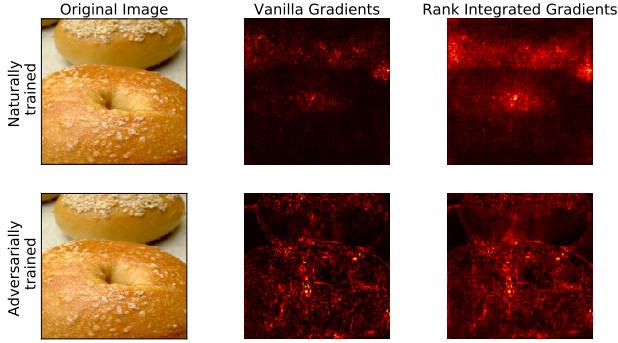


Figure 11: Comparison of RIG for naturally trained and adversarially robust neural networks. Adversarially robust neural networks have image representations that are much more aligned with human perception, which has been previously observed in [45].

9.3. Runtime measurements for low-rank approximation

There is minimal overhead to generating low-rank representations for images, with a distribution over 10 images across all possible ranks shown in Figure 12. The time required to generate rank- k approximations is independent of k , and generating arbitrary rank- k representations of a $(224 \times 224 \times 3)$ RGB image for ImageNet inference takes less than 1 second on an NVIDIA TITAN Xp GPU.

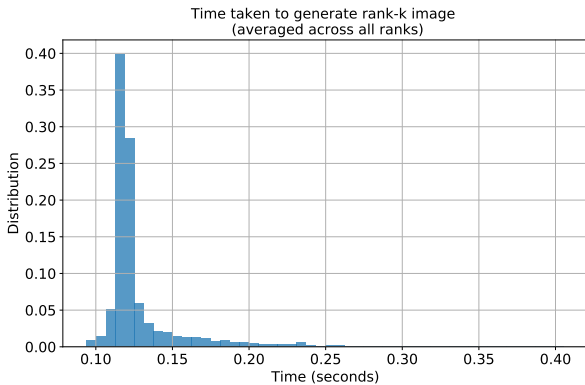


Figure 12: Time required to generate low-rank approximations of $(224 \times 224 \times 3)$ RGB images for ImageNet. Averaged across all possible ranks for 10 images.

9.4. Experimental results for VGG-19 and DenseNet

We observe that other state of the art models such as DenseNet [16] and VGG-19 [41] have similar rank-behavior to ResNet-50 models. In particular, we observe in Figure 13 that VGG-19 is more biased towards higher ranked representations, indicating a potentially larger vul-

# of Largest Singular values Truncated?	Test accuracy
2	31.71%
3	28.02%
5	13.54%
10	11.19%

Table 3: Test accuracy on full-rank CIFAR-10 test set for ResNet-50 trained on images with largest singular values removed from image.

nerability to adversarial examples.

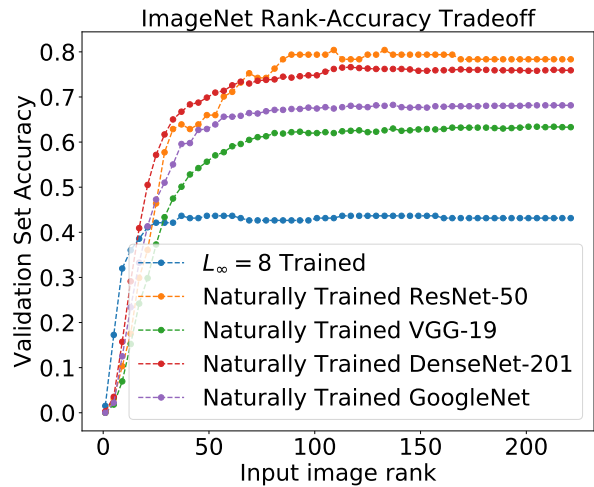


Figure 13: Rank spectrum for naturally trained ResNet-50, DenseNet-201, VGG-19 and GoogleNet architectures compared to a $L_\infty = 8/255$ robust ResNet-50.

9.5. Training on solely high-rank representations

Table 3 has full-rank test accuracies for ResNet-50 trained on modified versions of the CIFAR-10 dataset, where the largest k singular values for each image are deleted, leaving only higher-ranked features. Test accuracy as a function of training epoch can be found in Figure 14.

9.6. Training on solely low-rank representations

9.6.1 ImageNet

Figure 15 highlights full-rank test accuracy for ResNet-50 on modified low-rank versions of the ImageNet dataset. As expected, the performance of the models increases with increasing ranks of the images. We show this on ranks 30, 40, 50 and 100. For efficient training of ResNet-50 on the various ranks, we pre-process and store the low-rank copies of ImageNet. We trained each network for 24 hours on 4 NVIDIA V-100 GPUs.

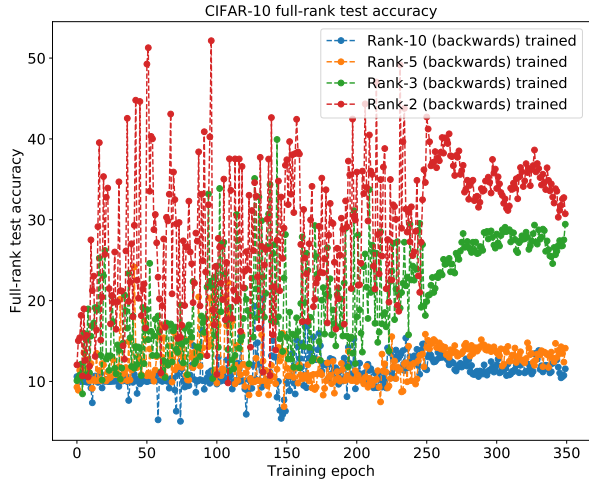


Figure 14: Full-rank test accuracy for ResNet-50 models trained where the largest 10, 5, 3, 2 singular values from each training image of CIFAR-10 are deleted.

9.6.2 CIFAR-10

Figure 16 highlights full-rank test accuracy for ResNet-50 trained on modified versions of the CIFAR-10 dataset. Models trained on rank-10, 20 and full rank have identical test accuracies, indicating that a large component of higher-ranked features do not contribute as much to prediction as those from ImageNet.

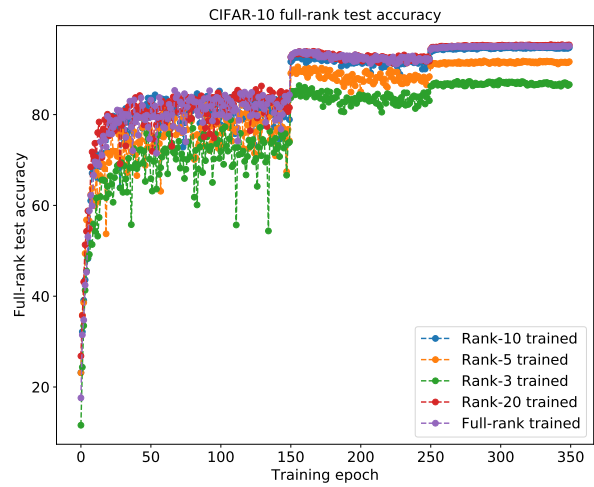


Figure 16: Full-rank test accuracy for ResNet-50 models trained on rank-3, 5, 10, 20 CIFAR-10 datasets.

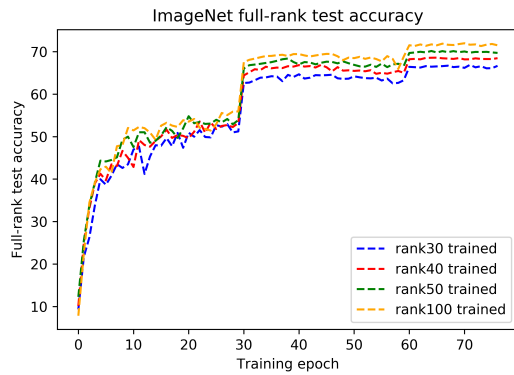


Figure 15: Full-rank test accuracy for ResNet-50 models trained on rank-30, 40, 50, 100 ImageNet datasets.

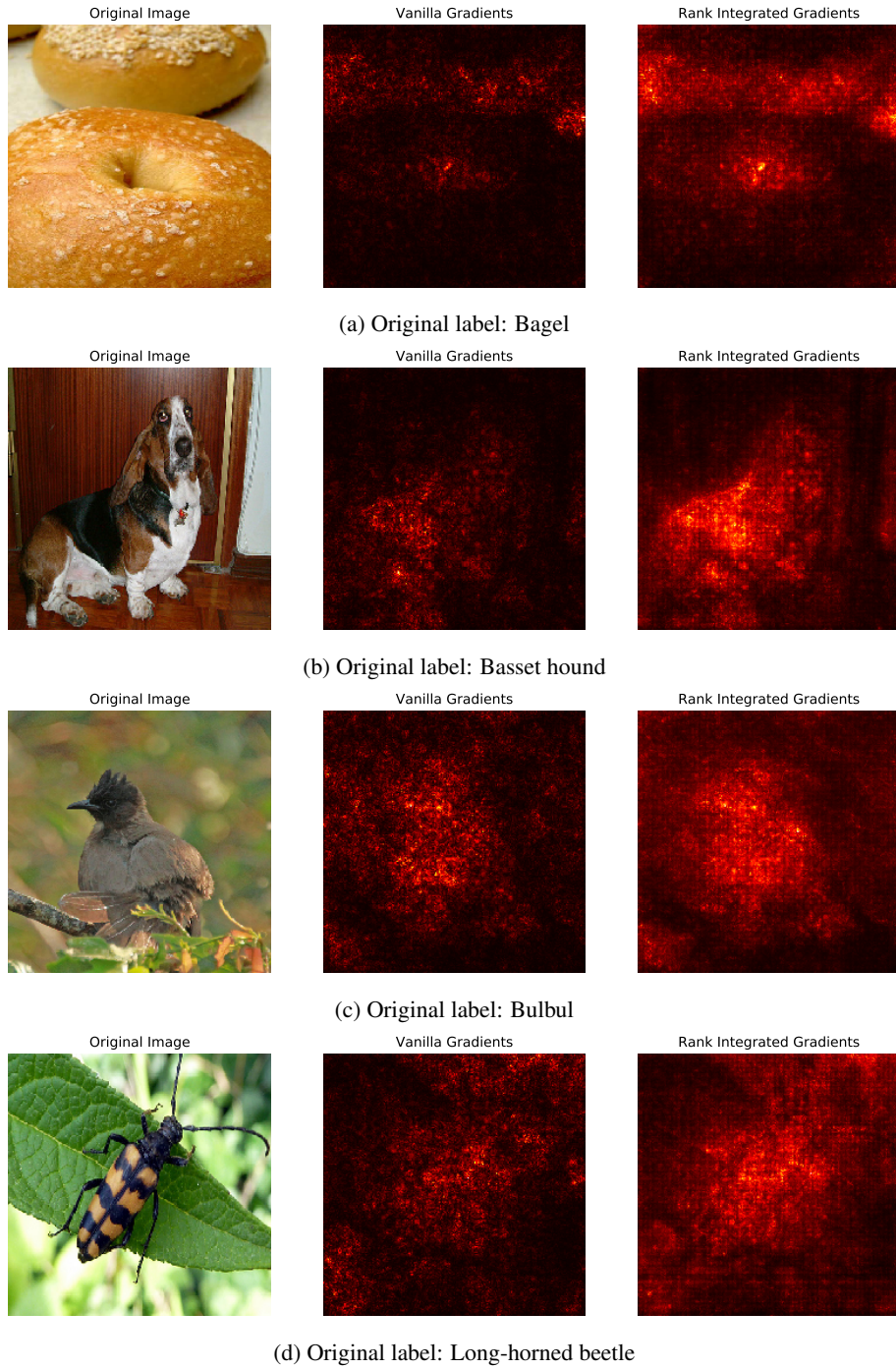
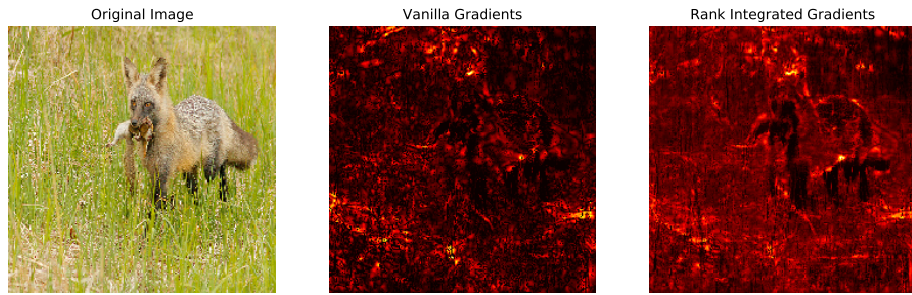
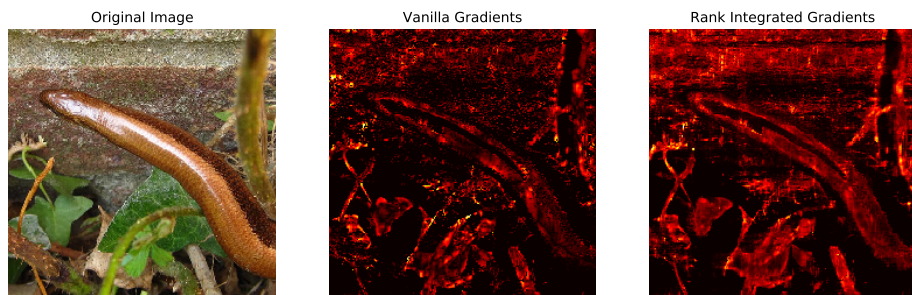


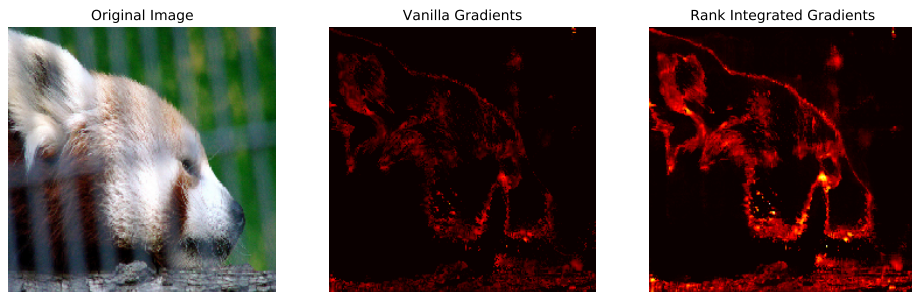
Figure 17: RIG plots for naturally trained ResNet-50 in PyTorch for images randomly chosen from the ImageNet validation set.



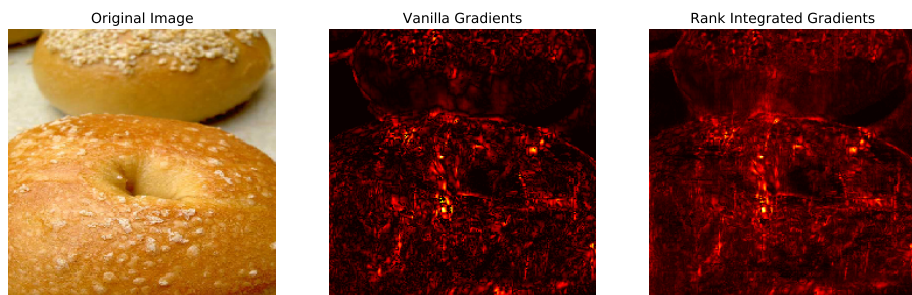
(a) Original label: grey fox



(b) Original label: thunder snake



(c) Original label: lesser panda, red panda



(d) Original label: bagel

Figure 18: RIG plots for $L_2 = 3.0$ robust ResNet-50 [9] for images randomly chosen from the ImageNet validation set.