

On the Adversarial Robustness of Visual Transformers

Rulin Shao
Xi'an Jiaotong University
shaorulin@stu.xjtu.edu.cn

Zhouxing Shi
University of California, Los Angeles
zhouxingshichn@gmail.com

Jinfeng Yi
JD AI Research
yijinfeng@jd.com

Pin-Yu Chen
IBM Research
pin-yu.chen@ibm.com

Cho-Jui Hsieh
University of California, Los Angeles
chohsieh@cs.ucla.edu

Abstract

Following the success in advancing natural language processing and understanding, transformers are expected to bring revolutionary changes to computer vision. This work provides the first and comprehensive study on the robustness of vision transformers (ViTs) against adversarial perturbations. Tested on various white-box and transfer attack settings, we find that ViTs possess better adversarial robustness when compared with convolutional neural networks (CNNs). We summarize the following main observations contributing to the improved robustness of ViTs: 1) Features learned by ViTs contain less low-level information and are more generalizable, which contributes to superior robustness against adversarial perturbations. 2) Introducing convolutional or tokens-to-tokens blocks for learning low-level features in ViTs can improve classification accuracy but at the cost of adversarial robustness. 3) Increasing the proportion of transformers in the model structure (when the model consists of both transformer and CNN blocks) leads to better robustness. But for a pure transformer model, simply increasing the size or adding layers cannot guarantee a similar effect. 4) Pre-training on larger datasets does not significantly improve adversarial robustness though it is critical for training ViTs. 5) Adversarial training is also applicable to ViT for training robust models. Furthermore, feature visualization and frequency analysis are conducted for explanation. The results show that ViTs are less sensitive to high-frequency perturbations than CNNs and there is a high correlation between how well the model learns low-level features and its robustness against different frequency-based perturbations.

1. Introduction

Transformers are originally applied in natural language processing (NLP) tasks as a type of deep neural net-

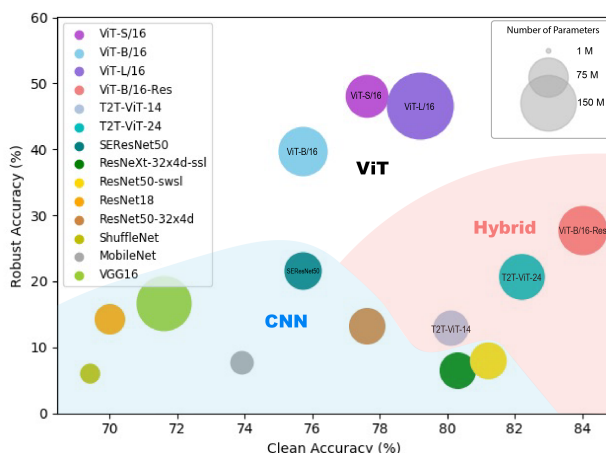


Figure 1. Robust accuracy v.s. clean accuracy. The robust accuracy is tested by the standard evaluation of AutoAttack [8]. The “Hybrid” class includes CNN-ViT and T2T-ViT as introduced in Section 3. Models with attention mechanisms have their names printed at the center of the circles. ViTs have the best robustness against adversarial perturbations. Introducing CNN or T2T modules to ViT can improve the clean accuracy but hurt the adversarial robustness. CNNs are more vulnerable to adversarial attacks.

work (DNN) mainly based on the self-attention mechanism [35, 10, 3], and transformers with large-scale pre-training have achieved state-of-the-art results on many NLP tasks [10, 23, 41, 33]. Recently, [11] applied a pure transformer directly to sequences of image patches (i.e., a vision transformer, ViT) and showed that the Transformer itself can be competitive with convolutional neural networks (CNN) on image classification tasks. Since then transformers have been extended to various vision tasks and show competitive or even better performance compared to CNNs and recurrent neural networks (RNNs) [4, 5, 48]. While ViT and its variants hold promise toward a unified machine learning paradigm and architecture applicable to different

data modalities, it remains unclear on the robustness of ViT against adversarial perturbations, which is critical for safe and reliable deployment of many real-world applications.

In this work, we conduct the first study on examining the adversarial robustness of ViTs on image classification tasks and make comparisons with CNN baselines. As highlighted in Figure 1, our experimental results illustrate the superior robustness of ViTs than CNNs in both white-box and black-box attack settings, based on which we make the following important observations:

- Features learned by ViTs contain less low-level information and have benefits including superior adversarial robustness and better transferability in transfer attack. ViTs achieve a lower attack success rate (ASR) of 51.9% compared with a minimum of 83.3% by CNNs in Figure 1. And they are less sensitive to high frequencies of the adversarial perturbations in our frequency study.
- It will take the cost of adversarial robustness to improve the classification accuracy of ViTs by introducing blocks to help learn low-level features as shown in Figure 1.
- Increasing the proportion of transformer blocks in the model leads to better robustness when the model consists of both transformer and CNN blocks. For example, the ASR decreases from 87.1% to 79.2% when 10 additional transformer blocks are added to T2T-ViT-14. However, increasing the size of a pure transformer model cannot guarantee a similar effect, e.g., the robustness of ViT-S/16 is better than that of ViT-B/16 in Figure 1.
- Pre-training on larger datasets does not improve adversarial robustness though it is critical for training ViT.
- The principle of adversarial training through min-max optimization [24, 46] can be applied to train robust ViTs.

2. Related Work

Transformer [35] architecture has achieved remarkable performance on many important Natural Language Processing (NLP) tasks, so the robustness of transformer has been studied on those NLP tasks. [17, 19, 29, 22, 12, 43] conducted adversarial attacks on transformers including pre-trained models, and in their experiments transformers usually show better robustness compared to models with structures such as LSTM or CNN, with a theoretical explanation provided in [17]. However, due to the discrete nature of NLP models, these studies are focusing on discrete perturbations (e.g., word or character substitutions) which are very different from small and continuous perturbations in computer vision tasks. Besides, [36] improved the robustness of pre-trained transformers from an information-theoretic perspective, and [30, 42, 39] studied the robust-

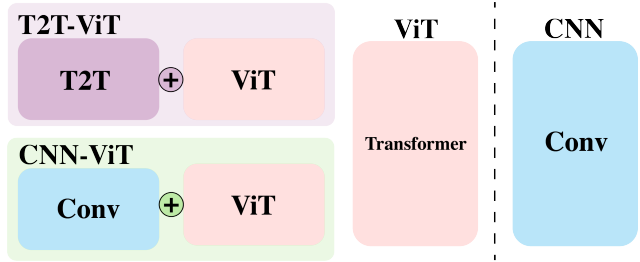


Figure 2. Visual transformers and CNNs investigated in our experiments. T2T-ViT and CNN-ViT take features learned by T2T modules and convolutional modules respectively as the inputs of ViT.

ness certification of transformer-based models. To the best of our understanding, this work is the first study that investigates the adversarial robustness (against small perturbations in the input pixel space) of transformers on computer vision tasks.

In the context of computer vision, the most relevant work is [1] which applies transformer encoder in the object detection task and reports better adversarial robustness. But the model they considered is a mix of CNN and transformer instead of the ViT model considered in this paper. Besides, the attacks they applied were weak, they provided neither study nor explanation on the benefit of adversarial robustness brought by the transformers.

3. Model Architectures

We first review the architecture of models investigated in our experiments including several vision transformers (ViTs) and CNN models. We summarize their information in Table 3. The weights of these models are all publicly available at [27, 37] such that our experiments can be easily reproduced.

3.1. Visual Transformers

We consider three visual transformers including the original Vision Transformer (ViT) [11], the hybrid model of ViT and CNN features (CNN-ViT) [11] and the hybrid model of ViT and T2T features (T2T-ViT) [44] as shown in Figure 2.

Vision Transformer (ViT) ViT [11] mostly follows the original design of Transformer [35, 10] on language tasks. For a 2D image $x \in \mathbb{R}^{H \times W \times C}$ with resolution $H \times W$ and C channels, it is divided into a sequence of $N = \frac{H \cdot W}{P^2}$ flattened 2D patches of size $P \times P$, $x_i \in \mathbb{R}^{N \times (P^2 \cdot C)}$ ($1 \leq i \leq N$). The patches are first encoded into patch embeddings with a simple convolutional layer, where the kernel size and stride of the convolution is exactly $P \times P$. In addition, there are also position embeddings to preserve positional information. Similar to BERT [10], a $[\text{CLS}]$ token is added for

Table 1. Comparison of the target models investigated in our experiments.

Model	ViT backbone			Attention	Params	Pretraining	
	Layers	Hidden size	MLP size			Pretraining dataset	Scale
ViT-S/16 [11]	8	786	2358	Self-attention	49M	ImageNet-21K [9]	14M
ViT-B/16 [11]	12	786	3072	Self-attention	87M	ImageNet-21K	14M
ViT-L/16 [11]	24	1024	4096	Self-attention	304M	ImageNet-21K	14M
ViT-B/16-Res [11]	12	786	3072	Self-attention	87M	ImageNet-21K	14M
T2T-ViT-14 [44]	14	384	1152	Self-attention	22M	-	-
T2T-ViT-24 [44]	24	512	1536	Self-attention	64M	-	-
SEResNet50 [18]	-	-	-	Squeeze-and-Excitation	28M	-	-
ResNeXt-32x4d-ssl [40]	-	-	-	-	25M	YFCC100M [34]	100M
ResNet50-swsl [40]	-	-	-	-	26M	IG-1B-Targeted [25]	940M
ResNet18 [14]	-	-	-	-	12M	-	-
ResNet50-32x4d [14]	-	-	-	-	25M	-	-
ShuffleNet [47]	-	-	-	-	2M	-	-
MobileNet [16]	-	-	-	-	4M	-	-
VGG16 [31]	-	-	-	-	138M	-	-

the classification.

Typically, ViT needs to be pre-trained on large datasets such as JFT-300M [32] and ImageNet-21k [9], and they are then fine-tuned on downstream tasks with smaller datasets such as ImageNet-1k [9] and CIFAR-10 [20]. Transformers do not generalize well when trained on insufficient data due to their lack of some inductive biases inherent to CNNs, such as translation invariance and locality [44]. So for ViT we only consider ViT-S/16, ViT-B/16 and ViT-L/16 pre-trained on ImageNet-21k and fine-tuned on ImageNet-1k. Here ‘‘S’’, ‘‘B’’ and ‘‘L’’ stand for ‘‘small’’, ‘‘base’’, and ‘‘large’’ respectively, and 16 stands for the patch size. ViT-B/16 and ViT-L/16 are provided by [11], and ViT-S/16 is provided by [37]. These ViT models achieve competitive performance compared with CNNs.

Hybrid of CNN and ViT (CNN-ViT) [11] also proposed a hybrid architecture for ViTs by replacing raw image patches with patches extracted from a CNN feature map. This is equivalent to adding learned CNN blocks to the head of ViT as shown in Figure 2. Following [11], we investigate ViT-B/16-Res in our experiments, where the input sequence is obtained by flattening the spatial dimensions of the feature maps from ResNet50.

Hybrid of T2T and ViT (T2T-ViT) [44] proposed to overcome the limitations of the simple tokenization in ViTs. They propose to progressively structurize an image to tokens with a T2T module, which recursively aggregate neighboring tokens into one token such that low-level structures can be better learned. Features learned by T2T modules are then fed into the following ViT backbone as shown in Figure 2. T2T-ViT was shown to perform better than ViT when trained from scratch on a midsize dataset. We inves-

tigate T2T-ViT-14 [44] and T2T-ViT-24 [44] trained from scratch on ImageNet.

3.2. Convolutional Neural Networks

We also study several CNN models for comparison including ResNet18 [14], ResNet50-32x4d [14], ShuffleNet [47], MobileNet [16] and VGG16 [31]. In addition, we also consider a CNN model with a self-attention mechanism, the Squeeze-and-Excitation (SE) block [18], which applies attention to channel dimensions, by fusing both spatial and channel-wise information within local receptive fields at each layer. We take SEResNet50 [18] for experiments.

The aforementioned CNNs are all trained on ImageNet from scratch. For a better comparison with pre-trained transformers, we also consider two CNN models pre-trained on larger datasets and fine-tuned on ImageNet, using semi-weakly supervised methods [40]. We take ResNeXt-32x4d-ssl [40] pre-trained on YFCC100M [34], and ResNet50-swsl that is pre-trained on IG-1B-Targeted [25]. They are both fine-tuned on ImageNet.

4. Adversarial Robustness Evaluation Methods

We consider the commonly used ℓ_∞ adversarial attacks to evaluate the robustness of target models. An ℓ_∞ attack is usually formulated as solving a constrained optimization problem:

$$\arg \max_{\mathbf{x}^{adv}} J(\mathbf{x}^{adv}, y) \quad \text{s.t.} \quad \|\mathbf{x}^{adv} - \mathbf{x}_0\|_\infty \leq \epsilon, \quad (1)$$

where for a clean example \mathbf{x}_0 with label y , we aim to find an adversarial example \mathbf{x}^{adv} within an ℓ_∞ ball with radius ϵ centered at \mathbf{x}_0 , such that the loss of the classifier for input \mathbf{x}^{adv} , denoted as $J(\mathbf{x}^{adv}, y)$, is maximized. We consider

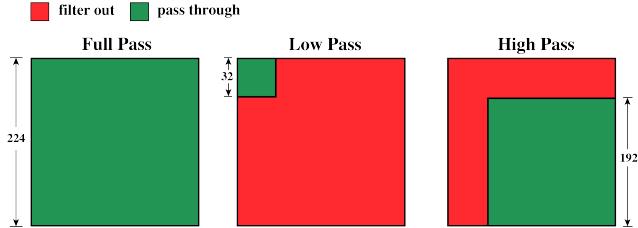


Figure 3. Filters for the frequency-based attack. The frequencies corresponding to the red part are filtered out, and the frequencies corresponding to the green part can pass through. “Full Pass” means all of the frequencies are preserved. “Low Pass” means only low-frequency components are preserved. “High Pass” preserves the high-frequency part.

untargeted attack in this paper, so an attack is successful if the perturbation successfully changes the model’s prediction. The attacks used in this paper are listed below.

Projected Gradient Decent The Projected Gradient Decent (PGD) attack [24] solves Eq. (1) by iteratively taking gradient ascent:

$$\mathbf{x}_{t+1}^{adv} = \text{Clip}_{\mathbf{x}_0, \epsilon}(\mathbf{x}_t^{adv} + \alpha \cdot \text{sgn}(\nabla_{\mathbf{x}} J(\mathbf{x}_t^{adv}, y))), \quad (2)$$

where \mathbf{x}_t^{adv} stands for the solution after t iterations, and $\text{Clip}_{\mathbf{x}_0, \epsilon}(\cdot)$ denotes clipping the values to make each $\mathbf{x}_{t+1, i}^{adv}$ within $[\mathbf{x}_{0, i} - \epsilon, \mathbf{x}_{0, i} + \epsilon]$, according to the ℓ_∞ threat model. As a special case, **Fast Gradient Sign Method (FGSM)** [13] uses a single iteration with $t = 1$.

AutoAttack AutoAttack [8] evaluates adversarial robustness further with a parameter-free ensemble of diverse attacks, including two parameter-free versions of PGD, and two existing complementary attacks, Fast Adaptive Boundary (FAB) attack [7] and Square Attack [6].

Frequency-based attack We also design a simple frequency-based attack to study the effect of different frequencies of the adversarial perturbations on the target models. We generate the initial perturbations following PGD. These perturbations are then transformed to the frequency domain using Discrete Cosine Transform (DCT), multiplied by one of the frequency filters in Figure 3 element-wise, and transformed back to the spatial domain using Inverse Discrete Cosine Transform (IDCT):

$$\mathbf{x}_{freq}^{adv} = \text{IDCT}(\text{DCT}(\mathbf{x}_{pgd}^{adv} - \mathbf{x}_0) \odot \mathbf{M}_f) + \mathbf{x}_0, \quad (3)$$

where \mathbf{x}_{pgd}^{adv} is the adversarial example generated by PGD, and the \mathbf{M}_f stands for the mask metric defined by one of the filter in Figure 3 for constraining the perturbation to have low or high frequency.

Black-box attack We also consider the transfer attack scenario in the black-box setting, where we study whether an adversarial perturbation generated by attacking the *source* model can successfully fool the *target* model. This test not only evaluates the robustness of models under the black-box setting, but also becomes a sanity check for detecting the obfuscated gradient phenomenon [2]. Previous works have demonstrated that single-step attacks like FGSM enjoys better transferability than multi-step attacks [21]. We thus use FGSM for transfer attack in our experiments.

5. Experimental Results

In the experiments, we will study the robustness of ViT and CNNs models under both white-box attacks and transfer attacks. Based on the experiments, we also analyze the reasons of why ViT shows better robustness than CNNs. In addition, we conduct a frequency study and visualize the learned features to verify our hypothesis that learning less high-frequency features makes ViT more robust than other models. Finally, we also conduct preliminary adversarial training experiments for ViT.

We load pre-trained visual transformer models and CNNs from the PyTorch image models library (timm) [37] and torchvision [27] respectively for convenient reproduction. In the evaluation, clean accuracy stands for the accuracy evaluated on standard and clean test examples without perturbation, robust accuracy stands for the accuracy on test examples with adversarial perturbations, and Attack Success Rate (ASR) stands for the rate of test examples such that the model makes a correct prediction on the clean input but predicts wrongly on the corresponding perturbed input. Note that for ASR it is the lower the better in terms of the robustness of a model. For experiments on models without adversarial training, we evaluate the clean accuracy of each model on the whole test set of ImageNet-1k [9], and we sample 1,000 test examples to evaluate ASR or robust accuracy. For adversarial training we use CIFAR-10 [20], detailed in Section 5.6.

5.1. Robustness under White-Box Attacks

Settings We use PGD and AutoAttack respectively to study the robustness under white-box attacks. We consider attack radius ϵ from $\{0.001, 0.003, 0.005, 0.01, 0.1\}$. For PGD attack [24], we fix the attack steps to $n_{iter} = 40$ with other parameters following the default setting of the implementation in Foolbox [28]. AutoAttack does not require any free hyper-parameters.

Results We present the results using PGD attack in Table 2 and the results using AutoAttack in Table 3. PGD attack has approximately 100% ASR on all the models when

$\epsilon = 0.1$. And when ϵ is smaller, PGD attack generally achieves the lower ASR on ViT models compared to CNNs. For example, when $\epsilon = 0.001$, the ASR for ViT-S/16 is only 44.6% while the ASRs for CNNs are at least 70.0%. In terms of AutoAttack in Table 3, results show that AutoAttack is stronger with higher ASR than PGD attack under the same ϵ respectively. And we can observe similar patterns in Table 3 as Table 2 with lower ASRs on ViTs compared to CNNs. These results demonstrate that ViT is more robust than CNNs under these white-box attacks. We also visualize the clean/robust accuracy tradeoff and model size of these 14 models, including 3 ViT-based models, 8 CNN-based models, and 3 hybrid models, in Figure 1.

Table 2. Attack success rate (%) of target models against 40-step PGD attack with different attack radii, and also the clean accuracy (“clean acc”). A model is considered to be more robust if the attack success rates are lower.

Model	Clean Acc	PGD Attack radius				
		0.001	0.003	0.005	0.01	0.1
ViT-S/16	77.6	44.6	75.4	89.8	99.0	100.0
ViT-B/16	75.7	51.1	85.4	94.0	99.1	100.0
ViT-L/16	79.2	44.9	76.6	90.1	98.2	100.0
ViT-B/16-Res	84.0	54.5	91.6	97.7	99.9	100.0
T2T-ViT-14	80.1	62.9	93.0	98.2	100.0	100.0
T2T-ViT-24	82.2	52.3	87.7	96.6	99.8	100.0
SEResNet50	75.7	64.6	95.1	99.2	99.9	100.0
ResNeXt-32x4d-ssl	80.3	77.0	97.1	98.8	99.5	99.7
ResNet50-sws1	81.2	75.3	97.1	98.6	99.6	100.0
ResNet18	70.0	75.1	98.0	99.4	99.9	100.0
ResNet50-32x4d	77.6	71.8	96.8	98.8	99.6	99.9
ShuffleNet	69.4	85.0	99.4	99.8	100.0	100.0
MobileNet	71.9	83.3	99.6	100.0	100.0	100.0
VGG16	71.6	73.7	96.8	98.7	99.4	100.0

5.2. Reasoning of Adversarial Robustness

5.2.1 Convolutional Blocks Make Models Less Robust

One interesting and perhaps surprising finding is that ViTs have worse robustness when modules that help to learn local structures are added ahead of the transformer blocks. For example, T2T-ViT adds several T2T modules to the head of ViT which iteratively aggregates the neighboring tokens into one token in each local perceptive field. ViT-B/16-Res takes the features generated by ResNet as inputs, which has the same effect as incorporating a trained CNN layer in front of the transformer blocks. Both modules help to learn local structures like edges and lines [44].

When the features learned by ResNet are introduced, the ASR of ViT-B/16 rises from 51.1% to 54.5% of ViT-B/16-Res under PGD attack, and from 51.9% to 72.3% under AutoAttack, with attack radius $\epsilon = 0.001$. A similar phenomenon can be observed by comparing the ASR of ViTs and T2T-ViTs. The ASR of T2T-ViT-14 is 18.3%

Table 3. Attack success rate (%) of target models against AutoAttack with different attack radii, in a similar format as Table 2.

Model	Clean Acc	AutoAttack attack radius			
		0.001	0.003	0.005	0.01
ViT-S/16	77.6	51.9	94.0	99.5	100.0
ViT-B/16	75.7	60.2	94.6	99.4	100.0
ViT-L/16	79.2	53.4	91.5	99.0	100.0
ViT-B/16-Res	84.0	72.3	99.1	100.0	100.0
T2T-ViT-14	80.1	87.1	99.9	100.0	100.0
T2T-ViT-24	82.2	79.2	99.7	100.0	100.0
SEResNet50	75.7	78.4	99.4	100.0	100.0
ResNeXt-32x4d-ssl	80.3	93.5	100.0	100.0	100.0
ResNet50-sws1	81.2	91.9	100.0	100.0	100.0
ResNet18	70.0	85.7	99.6	100.0	100.0
ResNet50-32x4d	77.6	86.8	99.8	100.0	100.0
ShuffleNet	69.4	93.9	100.0	100.0	100.0
MobileNet	71.9	92.2	100.0	100.0	100.0
VGG16	71.6	83.3	99.5	100.0	100.0

higher under PGD attack and 35.2% higher under AutoAttack compared with the ASR of ViT-S/16, under attack radius $\epsilon = 0.001$.

One possible explanation is that the introduced modules improve the classification accuracy by remembering the low-level structures that repeatedly appear in the training dataset. These structures such as edges and lines are high-frequent and sensitive to perturbations. Learning such features makes the model more vulnerable to adversarial attacks. Examination of this hypothesis is conducted in Section 5.4 and Section 5.5.

5.2.2 Increasing the Proportion of Transformer Blocks Can Improve Robustness

[15] mentioned that larger model does not necessarily imply better robustness. It can be confirmed by our experiments where ViT-S/16 shows better robustness than larger ViT-B/16 under both PGD attack and AutoAttack. In this case, simply adding transformer blocks to the classifier can not guarantee better robustness.

However, we recognize that for mixed architecture that has both CNN and transformer blocks, it is useful to improve adversarial robustness by increasing the proportion of the transformer blocks in the model. As shown in Table 2 and Table 3, T2T-ViT-24 has lower ASR than that of T2T-ViT-14 under both attacks. Besides the transformer block, we find that other attention mechanism modules such as SE block also improves adversarial robustness – as SEResNet50 has the least proportion of attention, the ASR of SEResNet50 is higher than ViT and T2T-ViT models but lower than other pure CNNs. These two findings are coherent since the attention mechanism is fundamental in Trans-

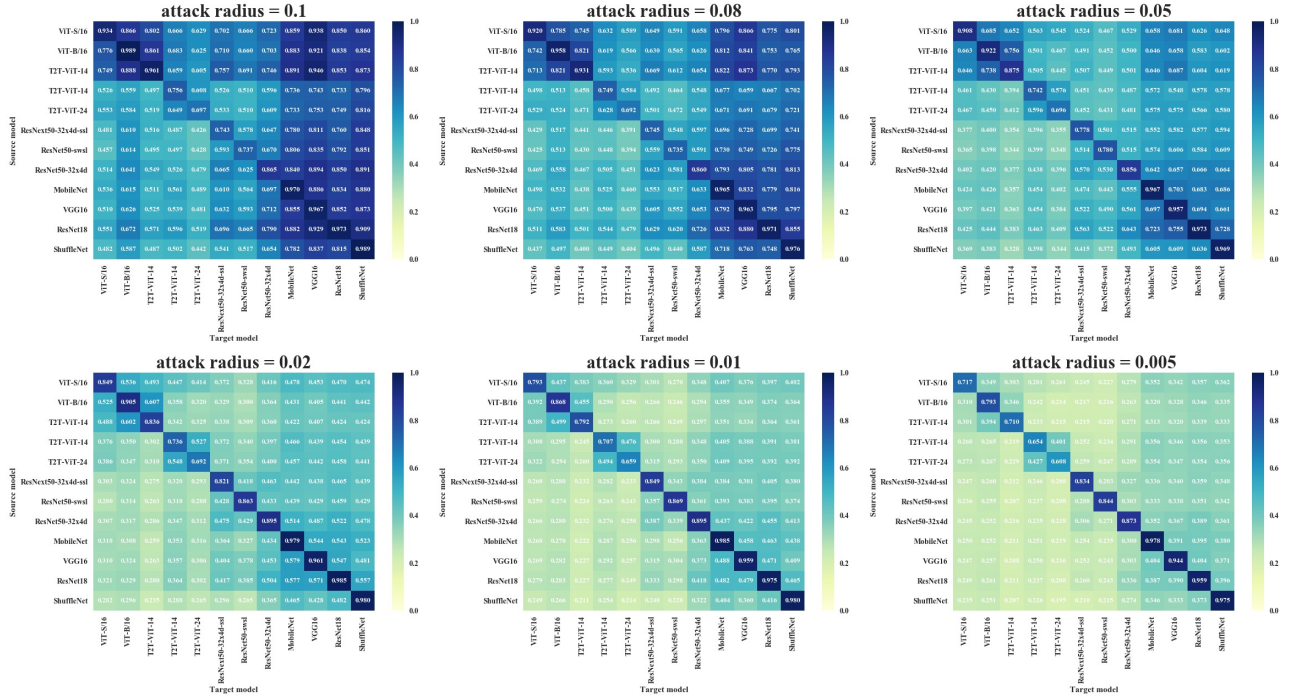


Figure 4. ASR of transfer attack using FGSM with different attack radii. The rows stand for the surrogate models used to generate adversarial examples in the white-box attack approach. The columns stand for the target models. Darker rows correlate to the source models that generate more transferable adversarial examples. While darker columns mean that the target models are more vulnerable to the transfer attack.

former blocks.

5.2.3 Pre-training Does Not Improve Robustness

Pre-training is critical for ViTs to achieve accuracy competitive with that achieved by CNNs trained from scratch [11]. However, pre-training does not appear to be the reason of better robustness here. To illustrate this point, we include CNNs pre-trained on large datasets and fine-tuned on ImageNet-1k to check the effect of pre-training on adversarial robustness. As shown in Table 3, CNNs pre-trained on large datasets IG-1B-Targeted [25] and YFCC100M [34] that are even larger than ImageNet-21k used by ViT, ResNet50-swsl and ResNeXt-32x4d-ssl, still have similar or even higher ASR than ResNet18 and ResNet50-32x4d that are not pre-trained. This supports our observation that pre-training may not be able to improve adversarial robustness. [15] also mentioned that pre-training without adversarial training techniques can not improve the adversarial robustness in their study. The resilience of ViT against perturbations corresponds more to the transformer structure rather than the pre-training.

5.3. Robustness under Transfer Attack

We also conduct transfer attack to test the adversarial robustness in the black-box setting. We apply the one-step FGSM attack on the source model and measure the ASR on the target model. Note that FGSM is used due to its better transferability than the iterative ones. We consider attacks with ℓ_∞ norm perturbation no larger than 0.1 and present the results in Figure 4. When the ViTs serve as target models, as shown in the first three columns of each subplot, the ASR of the transfer attack is quite low. On the other hand, when the ViTs are the source models, the adversarial examples they generate have higher ASR when transferred to other target models. As a result, the first three rows and the last seven columns are darker than the others. Besides, for the diagonal lines in the figure where FGSM actually attacks the models in a white-box setting, we can observe that ViTs are less sensitive to attack with smaller radii compared to CNNs, and T2T modules make ViTs more robust to such one-step attack. In addition, adversarial examples transfer well between models with similar structures. As ViT-S/16, ViT-B/16 and ViT-L/16 have similar structures, the adversarial examples generated by them can transfer well to each other, and it is similar for T2T-ViTs and CNNs respectively.

5.4. Frequency Study

We design a frequency study to verify our hypothesis that ViTs are adversarially more robust compared with CNNs because ViTs learn less high-frequency features. For adversarial perturbations generated by PGD attack, we first project them to the frequency domain by DCT. We design three frequency filters as shown in Figure 3: the full-pass filter, the low-pass filter, and the high-pass filter. And we take 32×32 pixels in the low-frequency area out of 224×224 pixels as the low-pass filter, and 192×192 pixels in the high-frequency area as the high-pass filter. Each filter allows only the corresponding frequencies to pass through – when the adversarial perturbations go through the low-pass filter, the high-frequency components are filtered out and vice versa, and the full-pass filter makes no change. We then apply these filters to the frequencies of the perturbations, and project them back to the spacial domain with the IDCT. We test the ASR under different frequency areas, and show the results in Table 4.

The ASR of ViT is relatively low in the “High-pass” column when only the high-frequencies of the perturbations are preserved. In contrast, CNNs show significantly higher ASR in the “High-pass” column than in the “Low-pass” column. It reflects that CNNs tend to be more sensitive to high-frequency adversarial perturbations compared to ViT. We also observe that adding modules that learn low-level structures makes the models more sensitive to the high-frequency perturbations. T2T-ViT-14, T2T-ViT-24 and ViT-B/16-Res have higher ASR in the “High-pass” column and lower ASR in the “Low-pass” column compared with vanilla ViTs, which verifies our hypothesis that low-level features are less adversarially robust. Besides, when adding more transformer blocks to the T2T-ViT model, the model becomes less sensitive to the high frequencies of the adversarial perturbations, e.g., the T2T-ViT-24 has an 8.7% lower ASR than that of the T2T-ViT-14 in the “High-pass” column.

5.5. Feature Visualization

We follow the work of [44] to visualize the learned features from the first blocks of the target models in Figure 5. We resize the input images to a resolution of 224×224 for CNNs and a resolution of 1792×1792 for ViTs and T2T-ViTs such that the feature maps from the first block are in the same shape of 112×112 . Low-level features like lines and edges are highlighted in blue (obviously perceptible) and green (minorly perceptible). As shown in Figure 5, CNNs like ResNet50-sws1 and ResNet50-32x4d learn features with obvious edges and lines. Minorly perceptible low-level features are learned by T2T-ViT-24 and ViT-B/16-Res. While it is hard to observe such information in the features maps learned by ViT-B/16.

The feature visualization combined with the frequency

study shows that the model’s vulnerability under adversarial perturbations is highly relative to the model’s tendency to learn low-level high-frequency features. Techniques that help the model learn such features may improve the performance on clean data but sacrifice adversarial robustness.

5.6. Adversarial Training

Finally, we conduct a preliminary experiment on adversarial training for ViT. We use a relatively smaller dataset, CIFAR-10 [20] with $\epsilon = 8/255$, for this experiment, and we take the ViT-B/16 model. Since originally this ViT was pre-trained on ImageNet with image size 224×224 and patch size 16×16 while image size on CIFAR-10 is as small as 32×32 , we downsample the weights for patch embeddings and resize patches to 4×4 , so that there are still 8×8 patches and we name the new model as ViT-B/4. Though ViT [11] originally enlarged input images on CIFAR-10 for natural fine-tuning and evaluation, we keep the input size as 32×32 to make the attack radius comparable. For training, we use PGD-7 (PGD with 7 iterations) [24] and TRADES [46] methods respectively, with no additional data during adversarial training. And we compare ViT with two CNNs, ResNet18 [14] and WideResNet-34-10 [45]. To save training cost, we train each model for 20 epochs only, although some prior works used around hundreds of epochs [24, 26] and are very costly for large models. We use a batch size of 128, an initial learning rate of 0.1, an SGD optimizer with momentum 0.9, and the learning rate decays after 15 epochs and 18 epochs respectively with a rate of 0.1. While we use a weight decay of 5×10^{-4} for CNNs as suggested by [26] that 5×10^{-4} is better than 2×10^{-4} , we still use 2×10^{-4} for ViT as we find 5×10^{-4} causes an underfitting for ViT. We evaluate the models with PGD-10 (PGD with 10 iterations) and AutoAttack respectively.

We show the results in Table 5. The ViT model achieves higher robust accuracy compared to ResNet18, and comparable robust accuracy compared to WideResNet-34-10, while ViT achieves much better clean accuracy compared to the other two models. Here ViT does not advance the robust accuracy after adversarial training compared to large CNNs such as WideResNet-34-10. We conjecture that ViT may need larger training data or longer training epochs to further improve its robust training performance, inspired by the fact that on natural training ViT is not able to perform well either without large-scale pre-training. And although T2T-ViT improved the performance of natural training when trained from scratch, our previous results in Table 2 and Table 3 show that the T2T-ViT structure may be inherently less robust. We have also tried [38] which was proposed to mitigate the overfitting of FGSM to conduct fast adversarial training with FGSM, but we preliminarily find that it can still cause catastrophic overfitting for ViT such that the test accuracy on PGD attacks remains al-

Table 4. Frequency study. The ASR of the target models against PGD attack. In the “Low-pass” column, only low-frequency adversarial perturbations are preserved and added to the input images. In the “High-pass” column, only high-frequency perturbations can pass through the filter. The “Full-pass” mode is the same as the traditional PGD attack. We set the attack step fixed to 40 and variate the attack radius to different values as shown in the second row.

Model	Low-pass					High-pass					Full-pass				
	0.001	0.003	0.005	0.01	0.1	0.001	0.003	0.005	0.01	0.1	0.001	0.003	0.005	0.01	0.1
ViT-S/16	26.0	31.9	35.3	40.2	43.8	29.2	39.3	49.4	59.6	76.6	44.6	75.4	89.8	99.0	100.0
ViT-B/16	28.1	35.7	39.7	44.2	50.4	33.7	46.9	56.0	66.6	78.1	51.1	85.4	94.0	99.1	100.0
ViT-L/16	25.1	35.9	41.7	49.8	58.0	27.1	37.7	43.4	52.5	71.1	44.9	76.6	90.1	98.2	100.0
ViT-B/16-Res	16.9	18.6	19.6	21.0	24.9	37.1	70.8	84.0	92.7	96.7	54.5	91.6	97.7	99.9	100.0
T2T-ViT-14	22.0	22.8	24.0	24.2	25.7	50.4	79.5	90.9	96.9	98.6	62.9	93.0	98.2	100.0	100.0
T2T-ViT-24	19.8	20.8	21.6	22.3	25.6	41.7	68.9	82.3	91.8	96.9	52.3	87.7	96.6	99.8	100.0
ResNet50-swsl	21.8	25.1	26.3	28.4	27.5	54.7	87.6	95.0	97.8	96.5	75.3	97.1	98.6	99.6	100.0
ResNet50-32x4d	25.0	33.7	37.3	41.0	38.5	52.3	82.9	92.6	96.7	96.5	71.8	96.8	98.8	99.6	99.9

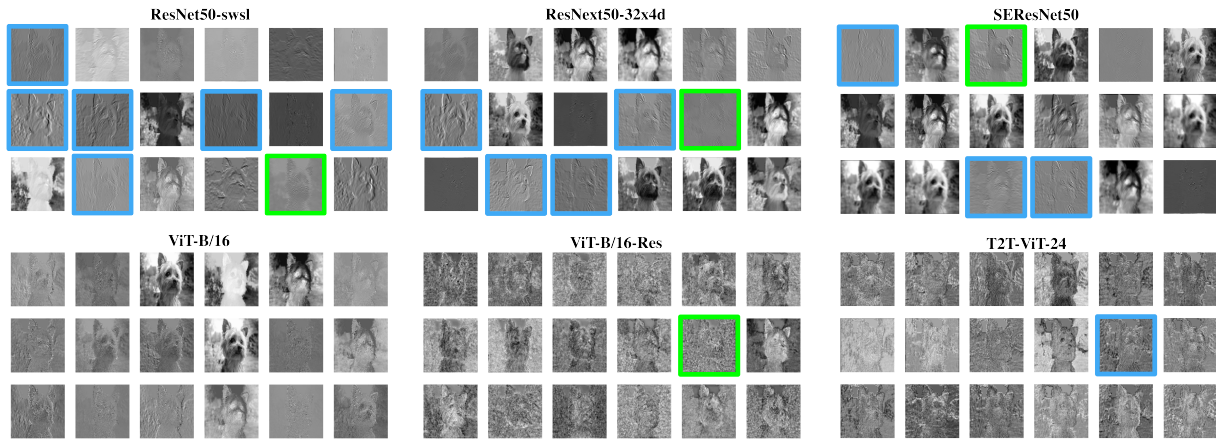


Figure 5. Feature visualization from the first blocks of ResNet50-swsl, ResNet50-32x4d, SEResNet50, ViT-B/16, ViT-B/16-Res, and T2T-ViT-24. Low-level structure features learned are highlighted in blue (obviously perceptible) and green (minorly perceptible). The CNNs in the first row learn more about the low-level features compared with the vision transformers in the second row. The ViTs pay more attention to the low-level structures and their feature maps become noisier when ResNet features are introduced (ViT-B/16-Res) or neighboring tokens are aggregated into one token recursively (T2T-ViT-24).

Table 5. Results of adversarial training of different models trained using PGD-7 and TRADES respectively on CIFAR-10. ViT-B/4 is a variant of ViT-B/16 where we downsample the patch embedding kernel from 16×16 to 4×4 to accommodate the smaller image size on CIFAR-10. We report the clean accuracy and robust accuracy evaluated with PGD-10 and AutoAttack respectively. Each model is trained using only 20 epochs to reduce the cost.

Model	Method	Clean	PGD-10	AutoAttack
ResNet18	PGD-7	77.3	48.9	44.4
	TRADES	77.6	49.4	44.9
WideResNet-34-10	PGD-7	80.3	52.2	48.4
	TRADES	81.6	53.4	49.3
ViT-B/4	PGD-7	85.9	51.7	47.6
	TRADES	85.0	53.9	49.2

most 0. We conjecture that this fast training method may be not suitable for pre-trained models or require further adjustments. Our experiments in this section demonstrate that

the adversarial training framework with PGD or TRADES is also applicable for transformers on vision tasks, and we provide baseline results and insights for future exploration and improvement.

6. Conclusion

This paper presents the first study on the robustness of ViTs against adversarial perturbations. Comparing to CNNs, our experimental results indicate that ViTs are more adversarially robust in both white-box and black-box settings. In particular, the features learned by ViTs contain less low-level information, contributing to improved robustness against adversarial perturbations that are often with high frequency. Introducing convolutional blocks in ViTs can facilitate the learning of low-level features but has a negative effect on adversarial robustness and makes the models more sensitive to high-frequency perturbations. We also ap-

ply adversarial training to ViTs. Our work provides a deep understanding of the intrinsic robustness of ViTs and can be used to inform the design of robust vision models based on the transformer structure.

References

- [1] Faisal Alamri, Sinan Kalkan, and Nicolas Pugeault. Transformer-encoder detector module: Using context to improve robustness to adversarial attacks on object detection. *arXiv preprint arXiv:2011.06978*, 2020. [2](#)
- [2] Anish Athalye, Nicholas Carlini, and David Wagner. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. *International Conference on International Conference on Machine Learning*, 2018. [4](#)
- [3] Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020. [1](#)
- [4] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European Conference on Computer Vision*, pages 213–229. Springer, 2020. [1](#)
- [5] Mark Chen, Alec Radford, Rewon Child, Jeffrey Wu, Heewoo Jun, David Luan, and Ilya Sutskever. Generative pre-training from pixels. In *International Conference on Machine Learning*, pages 1691–1703. PMLR, 2020. [1](#)
- [6] Francesco Croce, Maksym Andriushchenko, and Matthias Hein. Provable robustness of relu networks via maximization of linear regions. In *the 22nd International Conference on Artificial Intelligence and Statistics*, pages 2057–2066. PMLR, 2019. [4](#)
- [7] Francesco Croce and Matthias Hein. Minimally distorted adversarial examples with a fast adaptive boundary attack. In *International Conference on Machine Learning*, pages 2196–2205. PMLR, 2020. [4](#)
- [8] Francesco Croce and Matthias Hein. Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks. In *International Conference on Machine Learning*, pages 2206–2216. PMLR, 2020. [1](#), [4](#)
- [9] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. [3](#), [4](#)
- [10] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. [1](#), [2](#)
- [11] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. [1](#), [2](#), [3](#), [6](#), [7](#)
- [12] Siddhant Garg and Goutham Ramakrishnan. Bae: Bert-based adversarial examples for text classification. *arXiv preprint arXiv:2004.01970*, 2020. [2](#)
- [13] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014. [4](#)
- [14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. In *European conference on computer vision*, pages 630–645. Springer, 2016. [3](#), [7](#)
- [15] Dan Hendrycks, Kimin Lee, and Mantas Mazeika. Using pre-training can improve model robustness and uncertainty. In *International Conference on Machine Learning*, pages 2712–2721. PMLR, 2019. [5](#), [6](#)
- [16] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017. [3](#)
- [17] Yu-Lun Hsieh, Minhao Cheng, Da-Cheng Juan, Wei Wei, Wen-Lian Hsu, and Cho-Jui Hsieh. On the robustness of self-attentive models. In *ACL*, pages 1520–1529, 2019. [2](#)
- [18] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7132–7141, 2018. [3](#)
- [19] Di Jin, Zhijing Jin, Joey Tianyi Zhou, and Peter Szolovits. Is bert really robust? a strong baseline for natural language attack on text classification and entailment. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 8018–8025, 2020. [2](#)
- [20] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. [3](#), [4](#), [7](#)
- [21] Alexey Kurakin, Ian J. Goodfellow, and Samy Bengio. Adversarial machine learning at scale. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017. [4](#)
- [22] Linyang Li, Ruotian Ma, Qipeng Guo, Xiangyang Xue, and Xipeng Qiu. Bert-attack: Adversarial attack against bert using bert. *arXiv preprint arXiv:2004.09984*, 2020. [2](#)
- [23] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019. [1](#)
- [24] 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. [2](#), [4](#), [7](#)
- [25] Dhruv Mahajan, Ross Girshick, Vignesh Ramanathan, Kaiming He, Manohar Paluri, Yixuan Li, Ashwin Bharambe, and Laurens Van Der Maaten. Exploring the limits of weakly supervised pretraining. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 181–196, 2018. [3](#), [6](#)
- [26] Tianyu Pang, Xiao Yang, Yinpeng Dong, Hang Su, and Jun Zhu. Bag of tricks for adversarial training. *arXiv preprint arXiv:2010.00467*, 2020. [7](#)

- [27] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *arXiv preprint arXiv:1912.01703*, 2019. [2](#), [4](#)
- [28] Jonas Rauber, Roland Zimmermann, Matthias Bethge, and Wieland Brendel. Foolbox native: Fast adversarial attacks to benchmark the robustness of machine learning models in pytorch, tensorflow, and jax. *Journal of Open Source Software*, 5(53):2607, 2020. [4](#)
- [29] Zhouxing Shi and Minlie Huang. Robustness to modification with shared words in paraphrase identification. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: Findings*, pages 164–171, 2020. [2](#)
- [30] Zhouxing Shi, Huan Zhang, Kai-Wei Chang, Minlie Huang, and Cho-Jui Hsieh. Robustness verification for transformers. *arXiv preprint arXiv:2002.06622*, 2020. [2](#)
- [31] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. [3](#)
- [32] Chen Sun, Abhinav Shrivastava, Saurabh Singh, and Abhinav Gupta. Revisiting unreasonable effectiveness of data in deep learning era. In *Proceedings of the IEEE international conference on computer vision*, pages 843–852, 2017. [3](#)
- [33] Yu Sun, Shuohuan Wang, Yukun Li, Shikun Feng, Xuyi Chen, Han Zhang, Xin Tian, Danxiang Zhu, Hao Tian, and Hua Wu. Ernie: Enhanced representation through knowledge integration. *arXiv preprint arXiv:1904.09223*, 2019. [1](#)
- [34] Bart Thomee, David A Shamma, Gerald Friedland, Benjamin Elizalde, Karl Ni, Douglas Poland, Damian Borth, and Li-Jia Li. The new data and new challenges in multimedia research. *arXiv preprint arXiv:1503.01817*, 1(8), 2015. [3](#), [6](#)
- [35] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *arXiv preprint arXiv:1706.03762*, 2017. [1](#), [2](#)
- [36] Boxin Wang, Shuohang Wang, Y. Cheng, Zhe Gan, R. Jia, Bo Li, and Jing jing Liu. Infobert: Improving robustness of language models from an information theoretic perspective. *ArXiv*, abs/2010.02329, 2020. [2](#)
- [37] Ross Wightman. Pytorch image models. <https://github.com/rwightman/pytorch-image-models>, 2019. [2](#), [3](#), [4](#)
- [38] Eric Wong, Leslie Rice, and J Zico Kolter. Fast is better than free: Revisiting adversarial training. *arXiv preprint arXiv:2001.03994*, 2020. [7](#)
- [39] Kaidi Xu, Zhouxing Shi, Huan Zhang, Yihan Wang, Kai-Wei Chang, Minlie Huang, Bhavya Kailkhura, Xue Lin, and Cho-Jui Hsieh. Automatic perturbation analysis for scalable certified robustness and beyond. *Advances in Neural Information Processing Systems*, 33, 2020. [2](#)
- [40] I Zeki Yalniz, Hervé Jégou, Kan Chen, Manohar Paluri, and Dhruv Mahajan. Billion-scale semi-supervised learning for image classification. *arXiv preprint arXiv:1905.00546*, 2019. [3](#)
- [41] Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Ruslan Salakhutdinov, and Quoc V Le. Xlnet: Generalized autoregressive pretraining for language understanding. *arXiv preprint arXiv:1906.08237*, 2019. [1](#)
- [42] Mao Ye, Chengyue Gong, and Qiang Liu. Safer: A structure-free approach for certified robustness to adversarial word substitutions. *arXiv preprint arXiv:2005.14424*, 2020. [2](#)
- [43] Fan Yin, Quanyu Long, Tao Meng, and Kai-Wei Chang. On the robustness of language encoders against grammatical errors. *arXiv preprint arXiv:2005.05683*, 2020. [2](#)
- [44] Li Yuan, Yunpeng Chen, Tao Wang, Weihao Yu, Yujun Shi, Francis EH Tay, Jiashi Feng, and Shuicheng Yan. Tokens-to-token vit: Training vision transformers from scratch on imagenet. *arXiv preprint arXiv:2101.11986*, 2021. [2](#), [3](#), [5](#), [7](#)
- [45] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016. [7](#)
- [46] 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*, pages 7472–7482. PMLR, 2019. [2](#), [7](#)
- [47] Xiangyu Zhang, Xinyu Zhou, Mengxiao Lin, and Jian Sun. Shufflenet: An extremely efficient convolutional neural network for mobile devices. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6848–6856, 2018. [3](#)
- [48] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. *arXiv preprint arXiv:2010.04159*, 2020. [1](#)