

Intriguing Properties of Vision Transformers

Muzammal Naseer^{†*} Kanchana Ranasinghe^{*} Salman Khan^{*†}
 Munawar Hayat^{*} Fahad Shahbaz Khan^{*§} Ming-Hsuan Yang^{‡◦▽}

[†]Australian National University, ^{*}Mohamed bin Zayed University of AI, ^{*}Monash University
[§]Linköping University, [‡]University of California, Merced, [◦]Yonsei University [▽]Google Research
 muzammal.naseer@anu.edu.au

Abstract

Vision transformers (ViT) have demonstrated impressive performance across various machine vision problems. These models are based on multi-head self-attention mechanisms that can flexibly attend to a sequence of image patches to encode contextual cues. An important question is how such flexibility in attending image-wide context conditioned on a given patch can facilitate handling nuisances in natural images e.g., severe occlusions, domain shifts, spatial permutations, adversarial and natural perturbations. We systematically study this question via an extensive set of experiments encompassing three ViT families and comparisons with a high-performing convolutional neural network (CNN). We show and analyze the following intriguing properties of ViT: (a) Transformers are highly robust to severe occlusions, perturbations and domain shifts, e.g., retain as high as 60% top-1 accuracy on ImageNet even after randomly occluding 80% of the image content. (b) The robust performance to occlusions is not due to a bias towards local textures, and ViTs are significantly less biased towards textures compared to CNNs. When properly trained to encode shape-based features, ViTs demonstrate shape recognition capability comparable to that of human visual system, previously unmatched in the literature. (c) Using ViTs to encode shape representation leads to an interesting consequence of accurate semantic segmentation without pixel-level supervision. (d) Off-the-shelf features from a single ViT model can be combined to create a feature ensemble, leading to high accuracy rates across a range of classification datasets in both traditional and few-shot learning paradigms. We show effective features of ViTs are due to flexible and dynamic receptive fields possible via self-attention mechanisms. Code: <https://git.io/Js15X>

1 Introduction

As visual transformers (ViT) attract more interest in visual learning tasks [1], it becomes more pertinent to study characteristics of the learned representations. Specifically, from the perspective of safety-critical applications such as autonomous cars, robots and healthcare; learned representations need to be robust and generalizable to deal with real-world conditions. In this paper, we compare the performance of transformers with convolutional neural networks (CNNs) for handling nuisances (e.g., occlusions, distributional shifts, adversarial and natural perturbations) and generalization across different data distributions. Our in-depth analysis is based on three transformer families, ViT [2], DeiT [3] and T2T [4] across fourteen vision datasets. For the sake of brevity, we will be referring to all the transformer families considered in this work as ViT, unless otherwise mentioned.

We are intrigued by the fundamental differences in the operation of convolution and self-attention, that have not been extensively explored in the realm of robustness and generalization. While convolutions excel at learning local interactions between elements in the input domain e.g., edges and contour information, self-attention has been shown to learn global interactions in an effective manner e.g., the

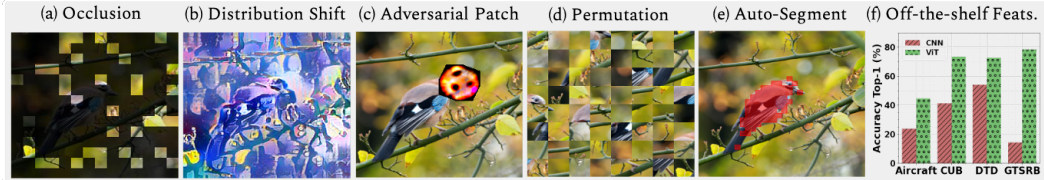


Figure 1: We notice intriguing properties of ViT including impressive robustness to (a) severe occlusions, (b) distributional shifts (e.g., stylization to remove texture cues), (c) adversarial perturbations, and (d) patch permutations. Moreover, our ViT models trained to focus on shape cues offer foreground segmentation without any pixel-level supervision (e). Finally, off-the-shelf features from ViT models generalize better than CNNs (f).

relations between distant object parts [5, 6]. Additionally, given a query embedding, self-attention finds its interactions with the other embeddings in the sequence, thereby conditioning on the local content while modeling relationships [7]. In comparison, convolutions are content-independent, as the same filter weights are applied to all inputs regardless of their distinct nature. Given the content-dependent long-range interaction modeling, our analysis shows that ViTs are able to flexibly adjust their receptive field to cope with nuisances in the data such as occlusions and perturbations.

Our systematic experiments and novel design choices lead to the following interesting observations:

- ViTs demonstrate strong robustness against severe occlusions for foreground objects, non-salient background and random patch locations; in comparison to state-of-the-art CNNs. For instance, with a significant random occlusion of up to 80%, DeiT model can maintain top-1 accuracy on ImageNet val. set as high as $\sim 60\%$ where the CNN has zero accuracy.
- In comparison to CNNs that exhibit a strong texture-bias (presented with texture and shape of the same object, they often make their decisions based on texture [8]), ViT are relatively less biased towards texture information and when robustly trained, they deliver a very strong shape-bias, better than CNNs and comparable to remarkable human performance on shape recognition. This highlights robustness of ViTs to deal with significant distribution shifts e.g., recognizing sketches and paintings from just the shape cues.
- ViTs show better robustness to other nuisance factors such as patch permutations, adversarial perturbations and common natural corruptions (e.g., noise, blur, contrast and pixelation artefacts) compared to CNNs. However, similar to CNNs, a shape-focused training renders them vulnerable against adversarial attacks and common corruptions.
- Apart from their promising robustness properties, off-the-shelf ViT features from ImageNet pretrained models generalize exceptionally well to new domains e.g., few-shot learning, fine-grained recognition, scene classification and long-tail recognition settings.

In addition to the extensive experimental analysis and new insights we have developed, this paper introduces several novel and intuitive design choices to highlight the remarkable potential of ViTs. For example, we propose an architectural modification to DeiT to encode shape-information via a dedicated token that demonstrates how seemingly contradictory cues can be modeled with different tokens within the same architecture, leading to favorable implications such as automated segmentation without pixel-level supervision. Furthermore, our off-the-shelf feature transfer approach utilizes an ensemble of representations derived from a single architecture, leading to state-of-the-art generalization performance with a pre-trained ViT (see Fig. 1).

2 Related Works

CNNs have shown excellent performance in independent and identically distributed (i.i.d) settings but remain highly sensitive to distributional shifts; adversarial noise [9], common image corruptions [10], and domain shifts such as from RGB to sketches or cartoons [11]. It is natural to ask if ViT that processes pixels/features based on self-attention offers any advantages in comparison to CNN. Shao *et al.* [12] analyzed ViTs against adversarial noise with a conclusion that ViTs are more robust to high frequency changes. Similarly, Bhojanapalli *et al.* [13] further study ViT against spatial perturbations [12] and its robustness to removal of any single layer. Since ViTs process image patches, we focus on how ViT offers robustness against localized adversarial patches [14, 15] and common corruptions. We note that a concurrent work from Paul and Chen [16] also develops similar insights on robustness of ViTs but with a somewhat different set of experiments.

Geirhos *et al.* [8] provide evidence that CNN mainly exploit texture to make a decision and give less importance to global shape. This is further backed-up by CNN ability to only use local features [17]. Recently, Islam *et al.* [18] quantify mutual information [19] between shape and texture features. Our analysis indicates that large ViT models have less texture bias and give relatively higher emphasis to shape information. ViT’s shape bias approaches human-level performance when directly trained on stylized ImageNet [8]. Our insights match with a concurrent recent work that demonstrates the importance of this trend on human behavioural understanding and bridging the gap between human and machine vision [20]. Another work on self-supervision [21] shows that ViT trained without labels attend to can automatically segment foreground objects from the background. In comparison, we show how shape-focused learning can impart similar capability in the image-level supervised models, without pixel supervision.

Zeiler *et al.* [22] introduced a method to visualize CNN features at different layers. They also study the performance of off-the-shelf features. In similar spirit, we study the generalization of off-the-shelf features of ViT in comparison to CNN. Receptive field is an indication of model’s ability to model long range dependencies. Mao *et al.* [23] compared receptive fields of Transformer based models with CNN. The receptive field of Transformer based models covers the entire input space, a property that resembles handcrafted features [24], but they have higher representative capacity. This allows Transformers to model global context and preserve the structural information compared to CNN [23]. In summary, our work is an effort to demonstrate the effectiveness of flexible receptive field and content-based context modeling towards robustness and generalization of ViT features.

3 Intriguing Properties of Vision Transformers

3.1 Are Vision Transformers Robust to Occlusions?

Vision Transformers’ receptive field spans over the entire image [23] and it models the interaction between the sequence of image patches using self-attention. We study if Vision Transformers provide robustness in occluded environments, where some or most of the image content is missing.

Occlusion Modeling: Consider a network f , that processes an input image x to predict a label y , where x can be represented as a sequence of patches i.e., $x = \{x_i\}_{i=1}^N$, and N is the total number of image patches [2]. While there can be infinite ways to model occlusion, we adopt a simple masking strategy, where we select a subset of the total image patches, $M < N$ and set pixel values of these patches to zero to create an occluded image, x' . We refer to this approach as PatchDrop. The objective is then to observe robustness such that $f(x')_{\text{argmax}} = y$. We experiment with three variants of our occlusion approach, **(a)** Random PatchDrop, **(b)** Salient (foreground) PatchDrop, and **(c)** Non-salient (background) PatchDrop.

–*Random PatchDrop:* A subset of M patches is randomly selected and dropped (Fig. 2). Recent Vision Transformers [2, 3, 4] usually divide an image into 196 patches belonging to a 14x14 spatial grid; i.e. an image of size 224x224x3 is split into 196 patches, each of size 16x16x3. As an example, dropping 100 such patches from the input will be equivalent to losing 51% of the image content.

–*Salient (foreground) PatchDrop:* Not all the pixels hold equal value for a classifier. Therefore, it is important to study the robustness of Transformers against occlusions of highly salient regions. To estimate salient regions, we leverage a self-supervised ViT model DINO [21], that uses attention to segment salient objects within an image. In particular, the spatial positions of information flowing into the final feature vector (class token) within the last attention block (of the ViT model) are exploited to locate the salient pixels for a given image. This allows to control the amount of salient information captured within the selected pixels by thresholding the quantity of attention flow. Following this approach, we select a subset of patches (from the same 196) containing the top N-percent of foreground information and drop them. Note that this N-percent does not always correspond to the pixel percentage, e.g. 50% of the foreground information of an image may be contained within only 10% of its pixels.

–*Non-salient (background) PatchDrop:* The least salient regions of the image are selected following the same approach as above [21]. The patches containing the lowest N-percent of foreground information are selected and dropped here. Again, this does not always correspond to the pixel percentage, e.g. 80% of the pixels may only contain 20% of the non-salient information for an image.

Figure 2: An example showing studied occlusion types (random, salient (foreground) and non-salient (background)). The occluded images are correctly classified by ViT (DeiT-S [3]) but mis-classified by CNN (ResNet50 [25]). Pixel values in occluded (black) regions are set to zero.

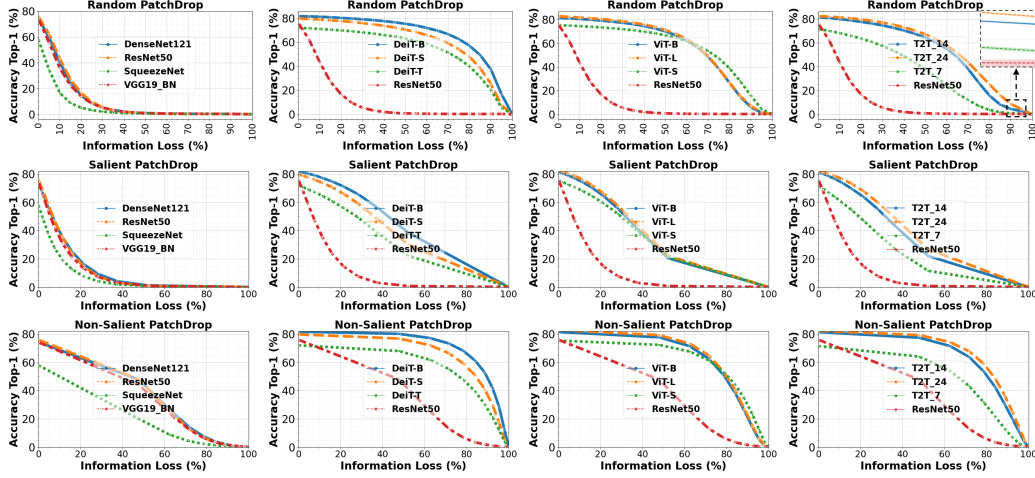
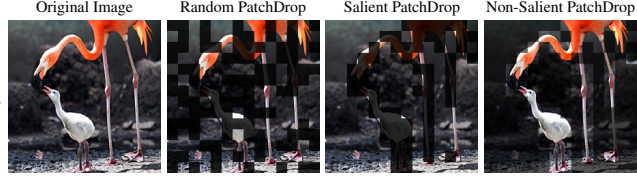


Figure 3: Robustness against object occlusion in images is studied under three PatchDrop settings (see Sec 3.1). (left) We study the robustness of CNN models to occlusions, and identify ResNet50 as a strong baseline. (mid-left) We compare the DeiT model family against ResNet50 exhibiting their superior robustness to object occlusion. (mid-right) Comparison against ViT model family. (right) Comparison against T2T model family.

Astounding Robustness of Transformers Against Occlusions: The effect of occlusion is studied on ImageNet [2] validation set (50k images). We define the ratio of dropped patches to total patches (M/N) as information loss (IL), and vary this value to obtain a range of occlusion levels for each PatchDrop methodology. Results (Top-1 %) reported in Fig. 3 present an astonishing robustness of ViT models in comparison to CNN. In the case of random PatchDrop, we report the mean and standard-deviation of accuracy across 5 runs. For Salient and Non-Salient Patchdrop, since the obtained occlusion mask is deterministic, we report the accuracy values over a single run. Randomly dropping 50% of image information almost completely collapses CNN’s recognition ability. For example, ResNet50 (23 Million parameters) has 0.1% accuracy in comparison to DeiT-S (22 Million parameters) which has 70% accuracy when 50% of the image content is removed. An extreme example can be observed when 90% of the image information is lost but DeiT-B still exhibits 37% recognition accuracy. This behavior is consistent within different ViT architectures [2, 3, 4]. Similarly, Vision Transformers show excellent robustness to the foreground (salient) and background (non-salient) content removal.

Class Token Preserves Information: In order to better understand model behavior against such occlusions, we visualize the attention (Fig. 4) from each head of different layers. While, initial layers pay more attention to occluded areas, deeper layers tend to focus more on the leftover information in an image. We then study if such changes from initial to deeper layers lead to token invariance against occlusion which is important for classification. We measure the correlation coefficient between features/tokens of original and occluded images. In the case of ResNet50, we consider features before the logit layer and for ViT models, class tokens are extracted from the last transformer block. Class tokens from transformers are significantly more robust and do not suffer much information loss as compared to ResNet50 features (Table 1). Further, we visualize the correlation coefficient across the 12 selected superclasses within ImageNet hierarchy and note that the trend holds across different class types, even for relatively small object types such as insects, food items and birds.

Given the intriguing robustness of transformer models and discriminability preserving behaviour of the learned tokens, an immediate question is if the learned representations in vision transformers are

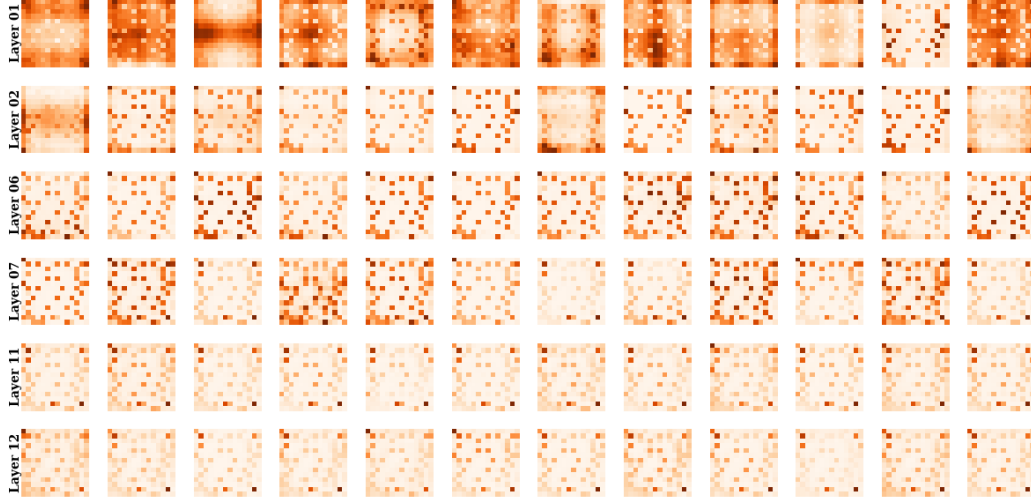


Figure 4: Attention maps (averaged over the entire ImageNet val. set) relevant to each head in multiple layers of an ImageNet pre-trained DeiT-B model. All images are occluded (RandomPatchDrop) with the same mask (bottom right). Observe how later layers clearly attend to non-occluded regions of images to make a decision, an evidence of the model’s highly dynamic receptive field.

Model	Correlation Coefficient: Random PatchDrop		
	25% Dropped	50% Dropped	75% Dropped
ResNet50	0.32±0.16	0.13±0.11	0.07±0.09
TnT-S	0.83±0.08	0.67±0.12	0.46±0.17
ViT-L	0.92±0.06	0.81±0.13	0.50±0.21
DeiT-B	0.90±0.06	0.77±0.10	0.56±0.15
T2T-24	0.80±0.10	0.60±0.15	0.31±0.17

Table 1: Correlation coefficient b/w features/final class tokens of original and occluded images for Random PatchDrop. Averaged across the ImageNet val. set.

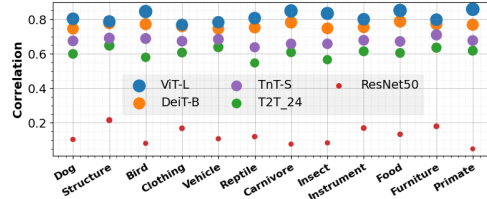


Table 2: Correlation b/w features/final tokens of original and occluded images for 50% Random Drop. Results are averaged across classes for each superclass.

biased towards texture? One can expect a biased model focusing only on texture to still perform well when the spatial structure for an object is partially lost. We investigate this question below.

3.2 Shape Vs. Texture: Can Transformer Model Both Characteristics?

Geirhos *et al.* [8] introduced a shape vs. texture hypothesis and proposed a training framework to increase shape-bias in convolutional neural networks (CNNs). *Firstly*, we run a similar analysis for ViT models, leading to a shape-bias much stronger than that of a CNN and comparable to the remarkable ability of human visual system to recognize shapes. However, this approach results in a significant drop in accuracy on the natural images. To resolve this behaviour, in a *second* approach, we introduce a shape token into the transformer architecture that dedicatedly learns to focus on shapes, thereby modeling both shape and texture related features within the same architecture using a distinct set of tokens. To this end, we distill the shape information from a pretrained CNN model with high shape bias [8]. Our distillation approach offers a balance between maintaining a reasonable classification accuracy while offering better shape-bias compared to the original ViT model.

We outline both approaches below. Note that the measure introduced in [8] is used to quantify shape-bias within ViT models and compare against their CNN counterparts.

Training without Local Texture: In this training approach, we first remove local texture cues from the training data by creating a stylized version of ImageNet [8] named SIN. We train tiny and small DeiT models [3] on this dataset. Typically, ViTs apply heavy data augmentations during training [3]. However, learning with SIN is a difficult task due to less texture details and applying further

augmentations on stylized samples destroys shape information and makes the training unstable. Therefore, we train models on SIN without applying any augmentation, label smoothing or mixup.

We note that ViT models trained on ImageNet exhibit higher shape bias in comparison to similar capacity CNN models e.g., DeiT-S with 22-Million parameters performs better than ResNet50 (23-Million parameters) (Fig. 5, right plot). When the SIN trained models are compared, ViT models consistently perform better than CNNs. Interestingly, DeiT-S [3] reaches human-level performance when trained on a SIN dataset (Fig. 5, left plot).

Shape Distillation: Knowledge distillation allows to compress large teacher models into smaller student models [26] as the teacher provides guidance to the student through soft labels. We introduce a new shape token and adapt attentive distillation [3] to distill shape knowledge from a CNN trained on SIN dataset (ResNet50-SIN [8]). We observe that ViT features are dynamic in nature and can be controlled by auxiliary tokens to focus on the desired characteristics. This means that a single ViT model can exhibit both a high shape and texture bias at the same time with separate tokens (Table 3). We obtain a more balanced performance for classification as well as shape-bias measure when the shape token is introduced (Fig. 6). In order to demonstrate that these distinct tokens (for classification and shape) indeed model unique features, we compute cosine similarity (averaged over ImageNet val. set) between class and shape tokens of our distilled models, DeiT-T-SIN and DeiT-S-SIN, which turns out to be 0.35 and 0.68, respectively. This is significantly lower than the similarity between class and distillation tokens [3]; 0.96 and 0.94 for DeiT-T and DeiT-S, respectively. This confirms our hypothesis about modeling distinct features with separate tokens within Transformers, a unique capability that cannot be straightforwardly achieved with CNNs. Further, it offers additional advantages as we explain next.

Shape-biased ViT Offers Automated Object Segmentation: Interestingly, training without local texture or with shape distillation allows a ViT to concentrate on foreground objects in the scene and ignore the background (see Fig. 4). This offers an automated semantic segmentation for an image, although the model was never shown pixel-wise object labels. This shows that promoting shape-bias within ViT acts as a self-supervision signal for the model to learn about distinct shape-related features that help successfully localise the right foreground object. Remarkably, a ViT trained without such an emphasis on shape performs much lower (Table 3).

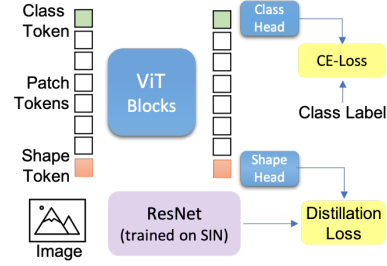


Figure 6: Shape Distillation.

In conclusion, the remarkable robustness of ViT against occlusions is not due to any biasness of ViT towards textures. Instead, we find that appropriately trained ViT models offer shape-bias nearly as high as the human ability to recognize shapes. This leads us to another related question: *Is positional*

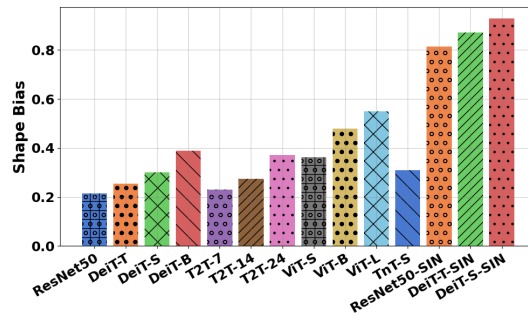
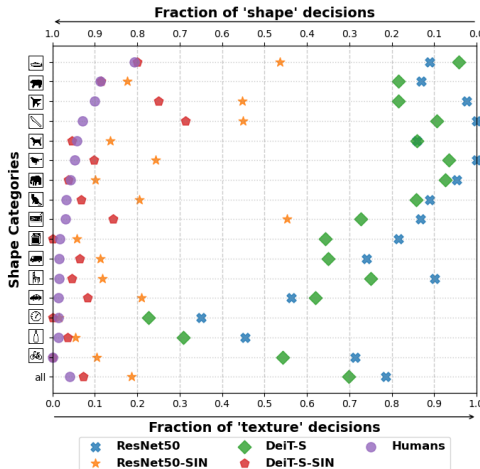


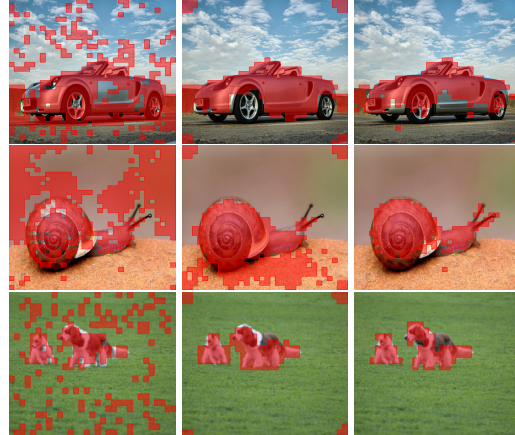
Figure 5: *Shape-bias Analysis:* ViTs perform better than a SOTA CNN model (ResNet50). Shape bias is defined as the fraction of correct decisions based on shape information presented in the images. Their shape bias increases significantly when trained on stylized ImageNet (SIN).

Model	Distilled	Token Type	ImageNet top-1 (%)	Shape Bias
DeiT-T-SIN	✗	cls	40.5	0.87
DeiT-T-SIN	✓	cls	71.8	0.35
DeiT-T-SIN	✓	shape	63.4	0.44
DeiT-S-SIN	✗	cls	52.5	0.93
DeiT-S-SIN	✓	cls	75.3	0.39
DeiT-S-SIN	✓	shape	67.7	0.47

Table 3: Performance comparison of models trained on SIN. ViT produces dynamic features that can controlled by auxiliary token. ‘cls’ represents the class token. During distillation cls and shape tokens converged to vastly different solution using the same features as compared to [3].

Model	Distilled	Token Type	Jaccard Index
DeiT-T-Random	✗	cls	19.6
DeiT-T	✗	cls	32.2
DeiT-T-SIN	✗	cls	29.4
DeiT-T-SIN	✓	cls	40.0
DeiT-T-SIN	✓	shape	42.2
DeiT-S-Random	✗	cls	22.0
DeiT-S	✗	cls	29.2
DeiT-S-SIN	✗	cls	37.5
DeiT-S-SIN	✓	cls	42.0
DeiT-S-SIN	✓	shape	42.4

Table 4: We compute the Jaccard similarity between ground truth and masks generated from the attention maps of ViT models (similar to [21] with threshold 0.9) over the PASCAL-VOC12 validation set. Only class level ImageNet labels are used for training these models. Our results indicate that supervised models can be used for an automated segmentation and perform comparably to self-supervised method DINO [21].



DeiT-S DeiT-S-SIN DeiT-S-SIN (Distilled)

Segmentation maps for different models. Shape distillation performs better than standard supervised models.

encoding the key that holds a high ViT performance under severe occlusions?, thereby allowing the later layers to perhaps recover the missing information with just a few image patches given their spatial ordering. We investigate this question next.

3.3 Does Positional Encoding preserve the Global Image Context?

Transformers’ ability to process long-range sequences in parallel using self-attention [27] (instead of a sequential design as in RNN) is invariant to sequence ordering. The obvious drawback is that it ignores the ordering of input sequence elements, which can be important. In visual domain, the order of patches represents the overall image structure and global composition. Since vision transformers takes sequence of images patches, changing the order of sequence e.g., shuffling the image patches can destroy the image structure. Current vision transformers [2, 3, 4, 23] use positional encoding to preserve this context. Here, we question if the sequence order modeled by positional encoding allows ViT to excel under occlusion handling? Our analysis however indicates that transformers show high permutation invariance to the patch positions. The effect of positional encoding towards injecting structural information of images to ViT models is limited. This observation is also consistent to findings in the language domain [28]. We outline our evaluation below.

Sensitivity to Spatial Structure: We eliminate the structural information within images (spatial relationships) as illustrated in Fig. 7 by defining a shuffling operation on input image patches. We observe that DeiT models [3] retain accuracy to a much higher extent in comparison to their CNN counterparts when spatial structure of input images is disturbed. This also indicates that positional encoding is not really crucial towards making the right classification decision, and the model is not ‘recovering’ global image context using the patch sequence information preserved in the positional encodings. Even when no such encoding, the ViT is able to maintain its performance and demonstrates better permutation invariance compared to a ViT using positional encodings (Fig. 8). Finally, when the patch size is varied during ViT training, the permutation invariance property is also degraded along with the accuracy on unshuffled natural images (Fig. 9). Overall, we attribute the permutation invariance behaviour of ViTs to their dynamic receptive field that is dependent on the input patch and can adjust attention with the other sequence element such that altering the order of patches does not degrade the performance significantly at moderate shuffling rates.



Figure 7: An illustration of shuffle operation applied on images used to eliminate their structural information.

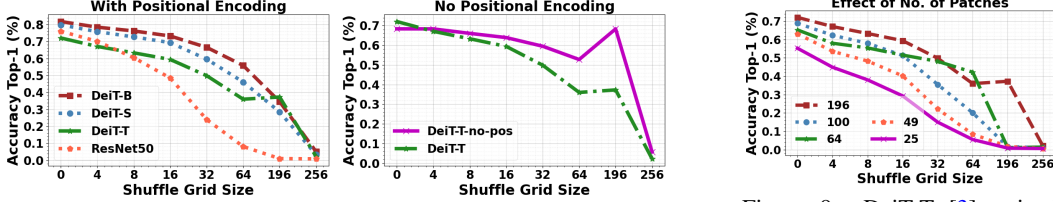


Figure 8: Models trained on 196 image patches. Top-1 (%) accuracy over ImageNet val. set when patches are shuffled. Note the performance peaks when shuffle grid size is equal to the original number of patches used during training, since it equals to only changing the position of input patch (and not disturbing the patch content).

Figure 9: DeiT-T [3] trained on different number of image patches. Reducing patch size decreases the overall performance but also increases sensitivity to shuffle grid size.

The above analysis shows that just like the texture-biasness hypothesis was wrong, the dependence on positional encodings to perform well under occlusions also is inaccurate. This leads us to the conclusion that such robustness is perhaps only due to the flexible and dynamic receptive field of ViT (see Fig. 4), that is dependent on the content of an input image. We now delve further into the robustness of ViT, and study its performance under adversarial perturbations and common corruptions.

3.4 Robustness of Vision Transformers to Adversarial & Natural Perturbations

The ability of ViT to encode shape information (Sec. 3.2) generates a fundamental question: *Does higher shape-bias means a higher robustness?* We investigate this by calculating mean corruption error (mCE) proposed by [10] on a variety of synthetic common corruptions such as rain, fog, snow and noise, etc. A ViT with similar parameters as CNN (e.g., DeiT-S) is more robust to image corruptions than ResNet50 trained with augmentations (Augmix [29]). Interestingly, convolutional and transformer models trained without augmentations on ImageNet or SIN are more vulnerable to image corruptions (Table 6). These findings are consistent with [30] which shows that data augmentations are necessary to improve robustness against common corruptions.

Trained with Augmentations						Trained without Augmentation			
DeiT-B	DeiT-S	DeiT-T	T2T-24	TnT-S	Augmix	ResNet50	ResNet50-SIN	DeiT-T-SIN	DeiT-S-SIN
48.5	54.6	71.1	49.1	53.1	65.3	76.7	77.3	94.4	84.0

Table 6: mean Corruption Error (mCE) across common corruptions [10] (lower the better) indicates that while ViTs have better robustness compared to CNNs, training to achieve a higher shape-bias makes both CNNs and ViTs more vulnerable to natural distribution shifts. All models trained with augmentations (ViT or CNN) have lower mCE in comparison to models trained without augmentations on ImageNet or SIN.

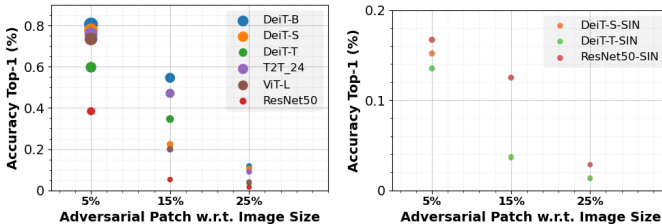


Figure 10: Robustness against adversarial patch attack. ViTs exhibit a higher robustness than CNN (even with less parameters). Models trained on ImageNet are more robust than the ones trained on SIN. Results are averaged across five runs of patch attack over ImageNet val. set.

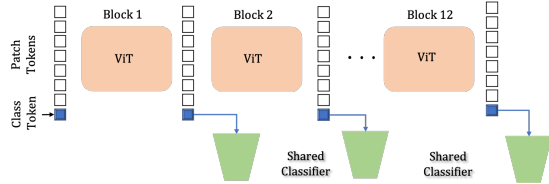


Figure 11: A single ViT model can provide a features ensemble since class token from each block can be processed by the classifier independently. This allows us to identify the most discriminative tokens useful for transfer learning.

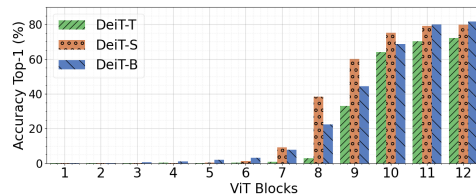


Figure 12: Top-1 (%) for ImageNet val. set for class tokens produced by each ViT block. Class tokens from the last few layers exhibit highest performance indicating the most discriminative tokens.

Blocks	Class Token	Patch Tokens	Top-1 (%)
Only 12 th (last block)	✓	✗	68.16
	✓	✓	70.66
From 1 st to 12 th	✓	✗	72.90
	✓	✓	73.16
From 9 th to 12 th	✓	✗	73.58
	✓	✓	73.37

Table 7: Ablative Study for off-the-shelf feature transfer on CUB [32] using ImageNet pre-trained DeiT-S [3]. A linear classifier is learned on only a concatenation of class tokens or the combination of class and averaged patch tokens at various blocks. We note class token from blocks 9-12 are most discriminative (Fig. 12) and have the highest transferability.

We observe similar behavior against adversarial patch attack [14]. ViTs shows higher robustness than CNN against untargeted, universal adversarial patch in white-box setting (full knowledge of model parameters). ViT and CNN trained on SIN are significantly more vulnerable to adversarial patch attack than models trained on ImageNet (Fig. 10), due to the shape-bias vs. robustness trade-off [30].

Given the remarkable robustness properties of ViT as well as their representation capability in terms of shape-biasness, automated segmentation and flexible receptive field, we question their utility as an off-the-shelf feature extractor to replace state-of-the-art CNN models as the default feature extraction mechanism [31]. Our analysis on the generalization of ViT features is outlined below.

3.5 Optimal Off-the-shelf Tokens for Vision Transformer

A unique characteristic of ViT models is that each block within the model produces a class token which can be processed by the classification head separately (Fig. 11). This allows us to measure the discriminative ability of each individual block of an ImageNet pre-trained ViT as shown in Fig. 12. Class tokens produced by the deeper blocks are more discriminative and we use this insight to identify the optimal ensemble of blocks whose tokens have the best downstream transferability.

Transfer Methodology: As illustrated in Fig. 12, we analyse the block-wise classification accuracy of DeiT models and discover the optimal discriminative information is captured within the class tokens of the last few blocks. To validate if this information can be combined to perform better, we conduct an ablation study for off-the-shelf transfer learning on fine-grained classification dataset (CUB [32]) using DeiT-S [3] as reported in Table 7. Here, we concatenate the class tokens (optionally combined with average patch tokens) from different blocks and train a linear classifier to transfer the features to downstream tasks. Note that a patch token is generated by averaging along the patch dimension. Concatenated class tokens from the last four blocks show best transfer learning performance. We refer to this transfer methodology as ‘DeiT-S (ensemble)’. Concatenation of both class and averaged patch tokens from all blocks exhibit similar performance compared to the tokens from the last four blocks but requires significantly large parameters to train. We conduct further experimentation with DeiT-S (ensemble) across a broader range of tasks to validate our hypothesis. In the following experiments, a CNN baseline is also compared against that uses the features extracted before the logit layer of a pre-trained ResNet50.

General Classification: We study the transferability of off-the-shelf features across several datasets including Aircraft [33], CUB [32], DTD [34], GTSRB [35], Fungi [36], Places365 [37], and iNaturalist [38]. These datasets aim at fine-grained recognition, texture classification, traffic sign recognition, fungi specie classification and scene recognition with 100, 200, 47, 43, 1394, 365, and 1010 classes respectively. We train a linear classifier on top of the extracted features over the train split of each dataset, and evaluate the performance on their test splits. ViT features show clear

improvements over the CNN baseline (Fig. 13). In fact, DeiT-T which has around 5 times less parameters than ResNet50 performs better. Further, the best results are achieved by our proposed ensemble strategy across all datasets.

Few-Shot Learning: In the case of few-shot learning, we consider meta dataset [39] which is a large-scale benchmark containing a diverse set of datasets covering multiple domains. We use the extracted features to learn a linear classifier on the support set for each query (similar to [40]), and evaluate using the standard FSL protocol defined in [39]. ViT features transfer better across these diverse domains (Fig. 13). We also highlight an improvement in QuickDraw, dataset containing hand-drawn sketches, which aligns with our findings on improved shape-bias of ViT models in contrast to CNN models (see Sec. 3.2). The best results are obtained using our ensemble strategy.

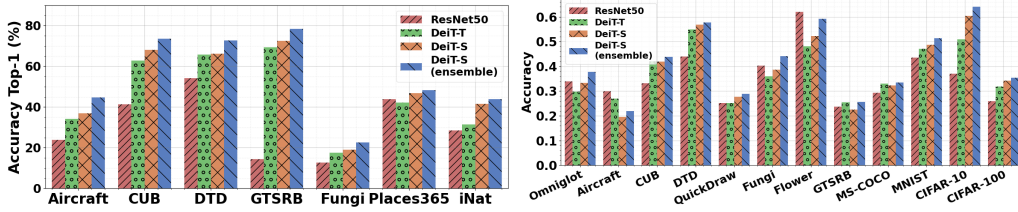


Figure 13: Off-the-shelf features of ViT transfer better than CNNs on various tasks.

4 Conclusion

This paper studies interesting properties of ViTs in terms of robustness and generalizability. Overall, we find favorable advantages of ViTs over CNNs for occlusion handling, robustness to distributional shifts and patch permutations, automatic segmentation with pixel supervision, and robustness against adversarial patches and common corruptions. Finally, we note impressive transferability of off-the-shelf ViT features to a number of downstream tasks with our proposed feature ensemble from a single ViT model. Going forward, an interesting direction is to explore if robustness to nuisance factors can be achieved while maintaining a strong shape bias in ViTs.

References

- [1] Salman Khan, Muzammal Naseer, Munawar Hayat, Syed Waqas Zamir, Fahad Shahbaz Khan, and Mubarak Shah. Transformers in vision: A survey. *arXiv preprint arXiv:2101.01169*, 2021. 1
- [2] 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, 3, 4, 7
- [3] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. *arXiv preprint arXiv:2012.12877*, 2020. 1, 3, 4, 5, 6, 7, 8, 9
- [4] 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. 1, 3, 4, 7
- [5] Prajit Ramachandran, Niki Parmar, Ashish Vaswani, Irwan Bello, Anselm Levskaya, and Jonathon Shlens. Stand-alone self-attention in vision models. *arXiv preprint arXiv:1906.05909*, 2019. 2
- [6] Han Hu, Zheng Zhang, Zhenda Xie, and Stephen Lin. Local relation networks for image recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3464–3473, 2019. 2
- [7] Ashish Vaswani, Prajit Ramachandran, Aravind Srinivas, Niki Parmar, Blake Hechtman, and Jonathon Shlens. Scaling local self-attention for parameter efficient visual backbones. *arXiv preprint arXiv:2103.12731*, 2021. 2
- [8] 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. 2, 3, 5, 6

- [9] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013. 2
- [10] Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *arXiv preprint arXiv:1903.12261*, 2019. 2, 8
- [11] Da Li, Yongxin Yang, Yi-Zhe Song, and Timothy M Hospedales. Deeper, broader and artier domain generalization. In *Proceedings of the IEEE international conference on computer vision*, pages 5542–5550, 2017. 2
- [12] Rulin Shao, Zhouxing Shi, Jinfeng Yi, Pin-Yu Chen, and Cho-Jui Hsieh. On the adversarial robustness of visual transformers. *arXiv preprint arXiv:2103.15670*, 2021. 2
- [13] Srinadh Bhojanapalli, Ayan Chakrabarti, Daniel Glasner, Daliang Li, Thomas Unterthiner, and Andreas Veit. Understanding robustness of transformers for image classification. *arXiv preprint arXiv:2103.14586*, 2021. 2
- [14] Tom B Brown, Dandelion Mané, Aurko Roy, Martín Abadi, and Justin Gilmer. Adversarial patch. *arXiv preprint arXiv:1712.09665*, 2017. 2, 9
- [15] Muzammal Naseer, Salman H Khan, Harris Khan, Fahad Shahbaz Khan, and Fatih Porikli. Cross-domain transferability of adversarial perturbations. *Advances in Neural Information Processing Systems*, 2019. 2
- [16] Sayak Paul and Pin-Yu Chen. Vision transformers are robust learners, 2021. 2
- [17] Wieland Brendel and Matthias Bethge. Approximating cnns with bag-of-local-features models works surprisingly well on imagenet. *arXiv preprint arXiv:1904.00760*, 2019. 3
- [18] Md Amirul Islam, Matthew Kowal, Patrick Esser, Sen Jia, Bjorn Ommer, Konstantinos G Derpanis, and Neil Bruce. Shape or texture: Understanding discriminative features in cnns. *arXiv preprint arXiv:2101.11604*, 2021. 3
- [19] David V Foster and Peter Grassberger. Lower bounds on mutual information. *Physical Review E*, 83(1):010101, 2011. 3
- [20] Shikhar Tuli, Ishita Dasgupta, Erin Grant, and Thomas L. Griffiths. Are convolutional neural networks or transformers more like human vision?, 2021. 3
- [21] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. *arXiv preprint arXiv:2104.14294*, 2021. 3, 7
- [22] Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *European conference on computer vision*, pages 818–833. Springer, 2014. 3
- [23] Yuxin Mao, Jing Zhang, Zhexiong Wan, Yuchao Dai, Aixuan Li, Yunqiu Lv, Xinyu Tian, Deng-Ping Fan, and Nick Barnes. Transformer transforms salient object detection and camouflaged object detection. *arXiv preprint arXiv:2104.10127*, 2021. 3, 7
- [24] Chuan Yang, Lihe Zhang, Huchuan Lu, Xiang Ruan, and Ming-Hsuan Yang. Saliency detection via graph-based manifold ranking. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3166–3173, 2013. 3
- [25] 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. 4
- [26] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015. 6
- [27] 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. 7
- [28] Kazuki Irie, Albert Zeyer, Ralf Schlüter, and Hermann Ney. Language modeling with deep transformers. *arXiv preprint arXiv:1905.04226*, 2019. 7
- [29] Dan Hendrycks, Norman Mu, Ekin D. Cubuk, Barret Zoph, Justin Gilmer, and Balaji Lakshminarayanan. AugMix: A simple data processing method to improve robustness and uncertainty. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2020. 8
- [30] Chaithanya Kumar Mummadi, Ranjitha Subramaniam, Robin Huttmacher, Julien Vitay, Volker Fischer, and Jan Hendrik Metzen. Does enhanced shape bias improve neural network robustness to common corruptions? In *International Conference on Learning Representations*, 2021. 8, 9
- [31] A. Razavian, Hossein Azizpour, J. Sullivan, and S. Carlsson. Cnn features off-the-shelf: An astounding baseline for recognition. pages 512–519, 2014. 9
- [32] P. Welinder, S. Branson, T. Mita, C. Wah, F. Schroff, S. Belongie, and P. Perona. Caltech-UCSD Birds 200. Technical Report CNS-TR-2010-001, California Institute of Technology, 2010. 9

- [33] S. Maji, J. Kannala, E. Rahtu, M. Blaschko, and A. Vedaldi. Fine-grained visual classification of aircraft. Technical report, 2013. [9](#)
- [34] M. Cimpoi, S. Maji, I. Kokkinos, S. Mohamed, , and A. Vedaldi. Describing textures in the wild. In *Proceedings of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2014. [9](#)
- [35] Sebastian Houben, Johannes Stallkamp, Jan Salmen, Marc Schlipsing, and Christian Igel. Detection of traffic signs in real-world images: The German Traffic Sign Detection Benchmark. In *International Joint Conference on Neural Networks*, number 1288, 2013. [9](#)
- [36] Brigit Schroeder and Yin Cui. Fgvcx fungi classification challenge 2018. In github.com/visipedia/fgvcx_fungi_comp, 2018. [9](#)
- [37] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017. [9](#)
- [38] Grant Van Horn, Oisin Mac Aodha, Yang Song, Yin Cui, C. Sun, Alexander Shepard, Hartwig Adam, P. Perona, and Serge J. Belongie. The inaturalist species classification and detection dataset. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8769–8778, 2018. [9](#)
- [39] Eleni Triantafillou, Tyler Zhu, Vincent Dumoulin, Pascal Lamblin, Kelvin Xu, Ross Goroshin, Carles Gelada, Kevin Swersky, Pierre-Antoine Manzagol, and Hugo Larochelle. Meta-dataset: A dataset of datasets for learning to learn from few examples. <http://arxiv.org/abs/1903.03096>, abs/1903.03096, 2019. [10](#)
- [40] Yonglong Tian, Yue Wang, Dilip Krishnan, Joshua B Tenenbaum, and Phillip Isola. Rethinking few-shot image classification: a good embedding is all you need? *arXiv preprint arXiv:2003.11539*, 2020. [10](#)

A Additional Qualitative Results

Here, we show some qualitative results, e.g., Figure 14 show the examples of our occlusion (random, foreground, and background) method. The performance of our shape models to segment the salient image is shown in Fig. 15. We show the dynamic behavior of ViT’s receptive field by visualizing the attention in the case of information loss in Fig. 16. Finally, we show adversarial patches optimized to fool different ViT models (Fig. 18).

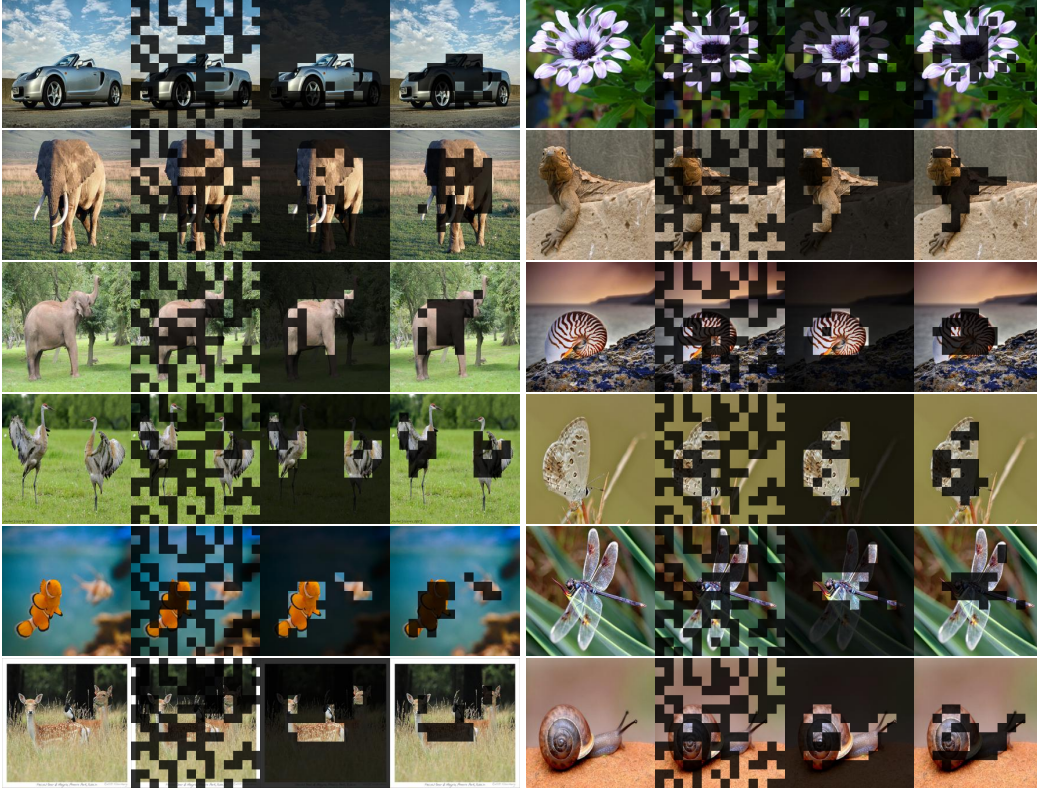


Figure 14: Visualizations for our three PatchDrop occlusion strategies: original, random (50% w.r.t the image), non-salient (50% w.r.t the foreground predicted by DINO), and salient (50% of the background as predicted by DINO) PatchDrop (shown from *left to right*). DeiT-B model achieves accuracies of 81.7%, 75.5%, 68.1%, and 71.3% across the ImageNet val. set for each level of occlusion illustrated from *left to right*, respectively.



Figure 15: Automatic segmentation of images using class-token attention for a DeiT-S model. Original, SIN trained, and SIN distilled model outputs are illustrated from *top to bottom*, respectively.

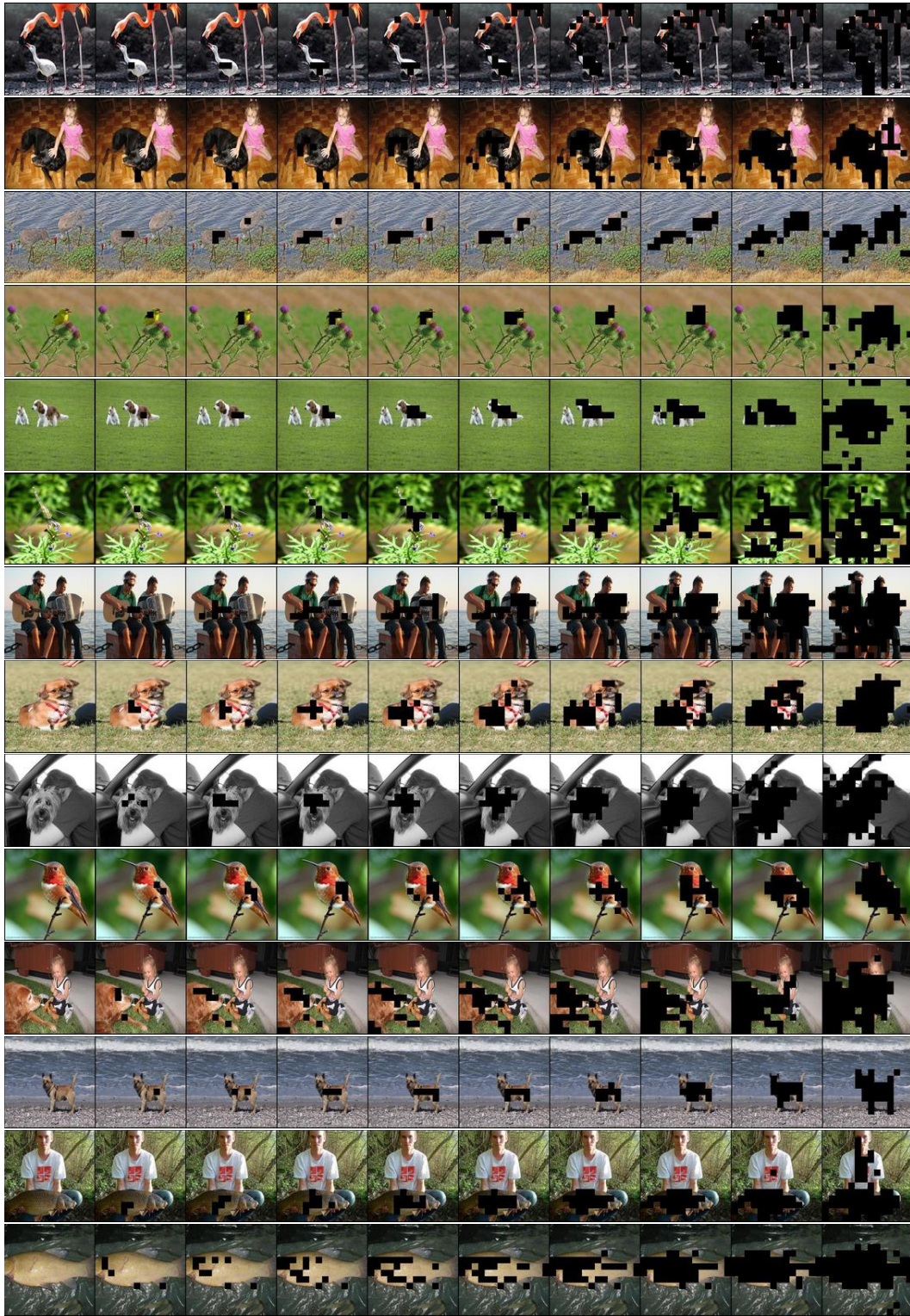


Figure 16: The variation (level increasing from *left to right*) of Salient PatchDrop on different images.

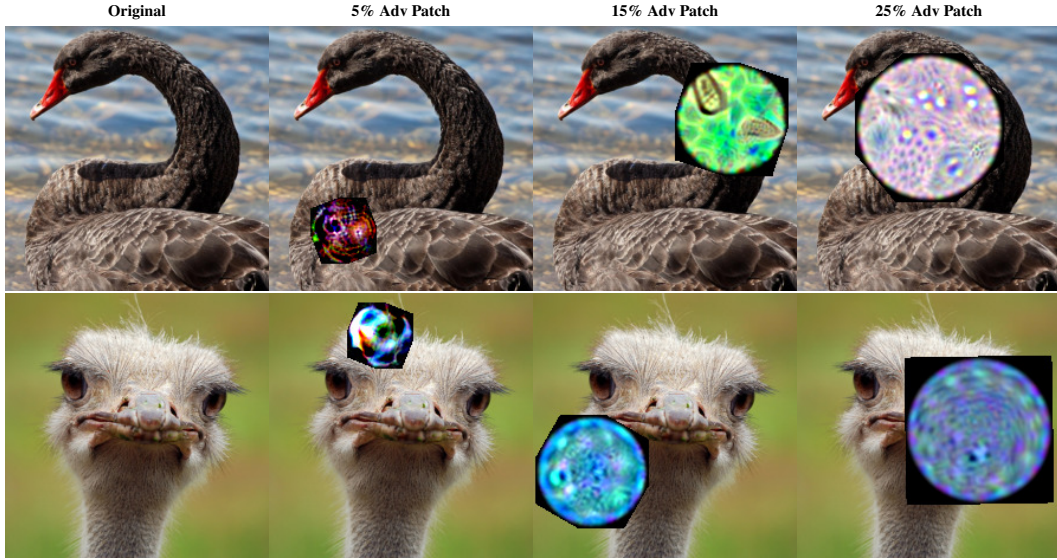


Figure 17: Adversarial patch (universal and untargeted) visualizations. *Top* row shows adversarial patches optimized to fool DeiT-S trained on ImageNet, while *bottom* row shows patches for DeiT-S-SIN. DeiT-S performs significantly better than DeiT-S-SIN. On the other hand, DeiT-SIN has higher shape-bias than DeiT-S.

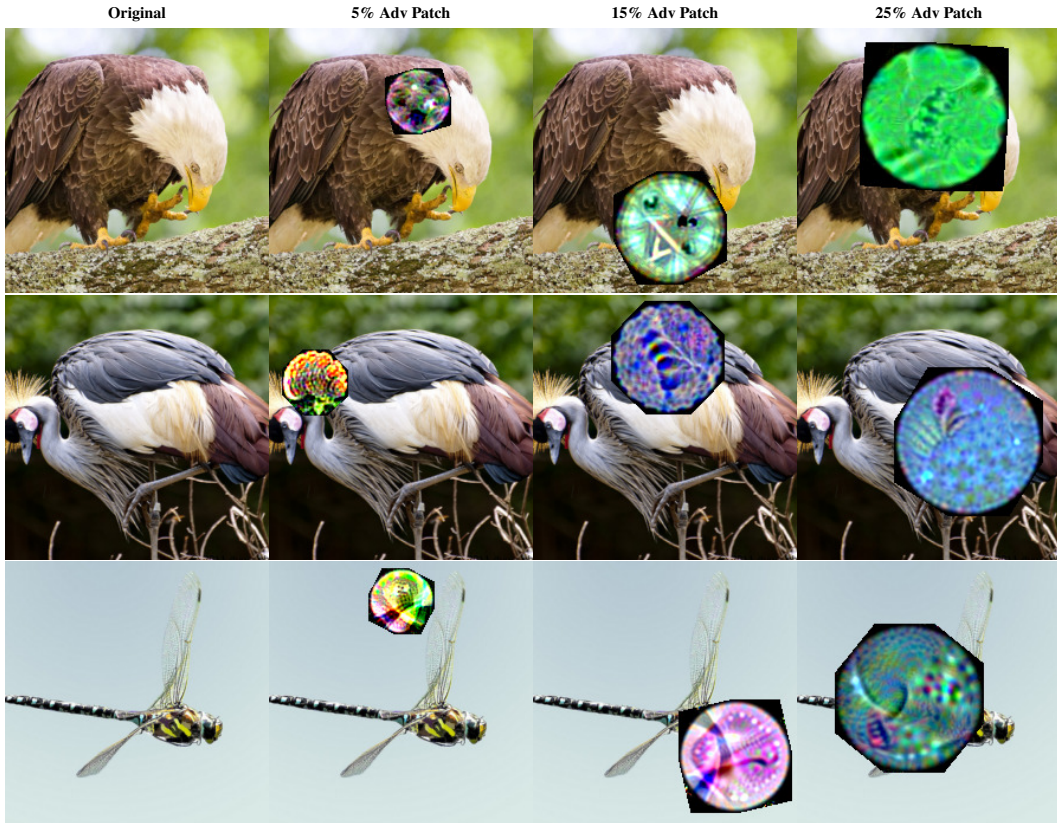


Figure 18: Adversarial patches (universal and untargeted) optimized to fool DeiT-T, DeiT-B, and T2T-24 models from *top* to *bottom*. These ViT models are more robust to such adversarial patterns than CNN (e.g., ResNet50).