

BACKGROUND & COMPARISON OF PERTURBATIONS

Vision transformers (ViTs) are often more robust than CNNs but still remain very vulnerable against corruptions and perturbations.

- Since ViTs are inherently patch-based, we explicitly study the sensitivity to patch corruptions/perturbations.
- We randomly sample a small number of patches to be perturbed/corrupted (10%, keeping the mask fixed) and introduce different perturbations and corruptions into the selected patches.



Clean Image (with Patch Mask) 63.8% Confidence



3.1% Confidence



Adversarial Perturbation Patch Corruption with on Patches (PGD-5) Random Noise (severity=5) 61.4% Confidence

Figure 1: Impact of different corruptions and perturbations on ViTs.

Observations:

- Transformers are very sensitive to adversarial patch perturbations.
- Directly introducing corruptions yields marginal degradation but is much more efficient.
- Occluding patches with noise can significantly hamper the prediction.

Idea: Occluding patches with noise is a good proxy of adversarial patch perturbations.

SENSITIVITY OF VITS TO PATCH-BASED CORRUPTIONS

mechanism against patch-based corruptions.



Observations: The self-attention mechanism is very sensitive to patch-based corruptions, which could be a major reason for the lack of robustness.

Improving Robustness of Vision Transformers by Reducing Sensitivity to Patch Corruptions

Yong Guo, David Stutz, Bernt Schiele



Patch Occlusion with Random Noise 17.3% Confidence

CONTRIBUTIONS

- reduce the sensitivity to them by aligning the intermediate features.
- When constructing patch corruptions, we develop a patch corruption model to find particularly different from adversarial training methods
- stable attention mechanism across layers.



Figure 2: Overview of our reducing sensitivity to patch corruptions (RSPC) training procedure.

Key Components of the proposed RSPC:

• Finding Vulnerable Patches to be Corrupted:

• Reducing Sensitivity via Feature Alignment:

 $\min_{\mathcal{T}} \max_{x \sim \mathcal{D}} \mathbb{E}_{x \sim \mathcal{D}} [\mathcal{L}_{ce}(x) + \lambda \mathcal{L}_{align}(x, \hat{x})],$

• We propose a new training method that **improves robustness by reducing sensitivity to patch** corruptions (RSPC). To this end, we first construct effective patch-based corruptions and then

vulnerable patches that severely distract intermediate attention layers. In practice, the corruption model is trained adversarially to the classification model, which, however, is essentially

• In experiments, we demonstrate that the robustness improvement against patch corruptions can generalize well to diverse architectures on various robustness benchmarks. More critically, we can show, both qualitatively and quantitatively, that these improvements stem from the **more**

> $\max \mathbb{E}_{x \sim \mathcal{D}} \mathcal{L}_{\text{align}}(x, \hat{x}),$ where $\mathcal{L}_{align}(x, \hat{x}) = \frac{1}{L} \sum_{l=1}^{L} \|\mathcal{F}_l(x) - \mathcal{F}_l(\hat{x})\|^2$.

Results: Improving Corruption Robustness



	#FLOPs (G)	#Params (M)	ImageNet	Robustness Benchmarks			
				IN-A	IN-C \downarrow	IN-C w/o Noise ↓	IN-P \downarrow
	1.3	5.7	72.2	7.3	71.1	72.9	56.7
	1.4	5.7	73.3	8.9	68.4	70.4	53.7
	1.3	10.9	79.2	14.6 (+0.0)	57.0 (-0.0)	58.9 (-0.0)	39.1 (-0.0)
	1.3	10.9	79.5	16.5 (+1.9)	55.7 (-1.3)	57.5 (-1.4)	38.0 (-1.1)
	3.5	7.5	80.1	21.9 (+0.0)	58.3 (-0.0)	59.8 (-0.0)	38.3 (-0.0)
	3.5	7.5	80.3	23.6 (+1.7)	57.2 (-1.1)	58.4 (-1.4)	37.3 (-1.0)
	4.6	22.1	79.9	6.3	54.6	56.6	36.9
	5.4	27.8	81.5	18.9	49.8	52.1	35.8
	4.7	23.3	81.9	25.7 (+0.0)	49.4 (-0.0)	51.6 (-0.0)	35.2 (-0.0)
	4.7	23.3	82.2	27.9 (+2.2)	48.4 (-1.0)	50.4 (-1.2)	34.3 (-0.9)
	6.7	25.7	83.5	33.9 (+0.0)	48.5 (-0.0)	50.7 (-0.0)	34.5 (-0.0)
	6.7	25.7	83.6	36.8 (+2.9)	47.5 (-1.0)	49.4 (-1.3)	33.5 (-1.0)
	17.6	86.6	83.6	35.9	51.7	-	-
	17.6	86.6	82.0	27.4	48.5	50.9	32.1
	17.7	86.5	82.4	29.0	46.9	49.3	32.2
	17.7	91.8	82.6	28.5 (+0.0)	46.8 (-0.0)	49.8 (-0.0)	31.9 (-0.0)
	17.7	91.8	82.8	32.1 (+3.6)	45.7 (-1.1)	48.5 (-1.3)	31.0 (-0.8)
	11.3	50.5	83.9	39.6 (+0.0)	46.1 (-0.0)	48.1 (-0.0)	31.3 (-0.0)
	11.3	50.5	84.2	41.1 (+1.5)	44.5 (-1.6)	46.8 (-1.3)	30.0 (-1.2)

