

# Exploring Adversarial Attacks and Defenses in Vision Transformers trained with DINO

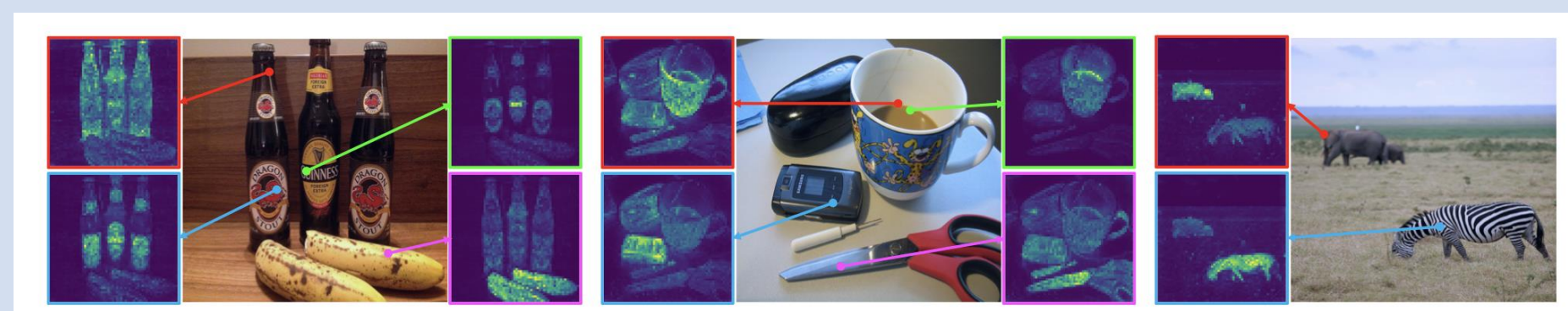


Read the paper

Javier Rando, Nasib Naimi, Thomas Baumann, Max Mathys  
ETH Zurich

## 1 Introduction

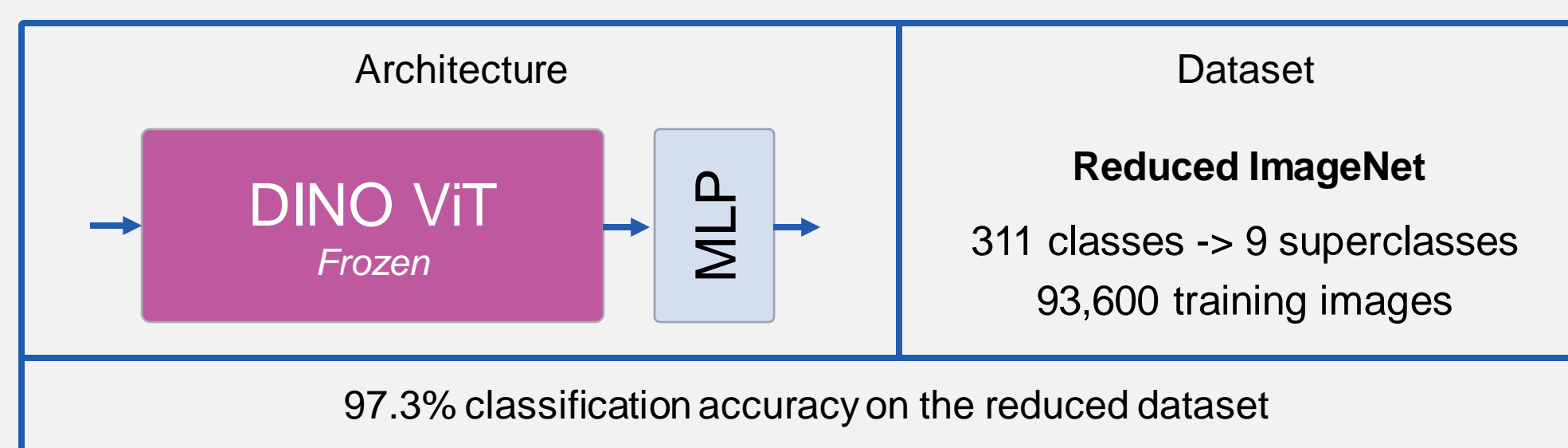
Self-supervised Vision Transformers trained with DINO learn human interpretable features<sup>1</sup>.



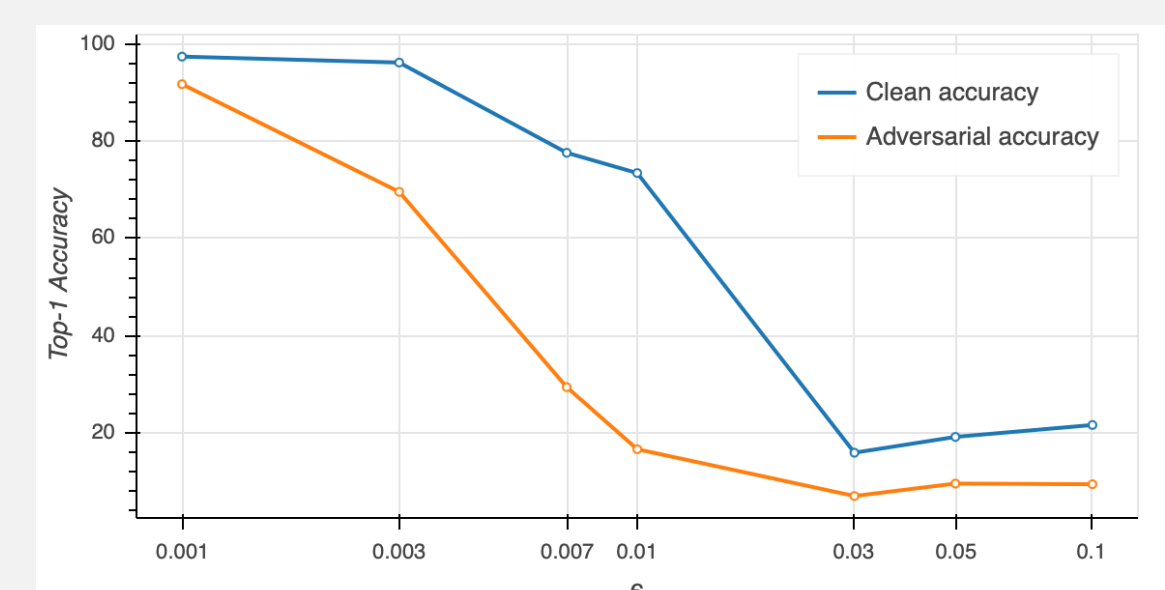
Are they more robust to Adversarial Attacks that preserve input similarity?

## 5 Defenses Under Limited Computational Power

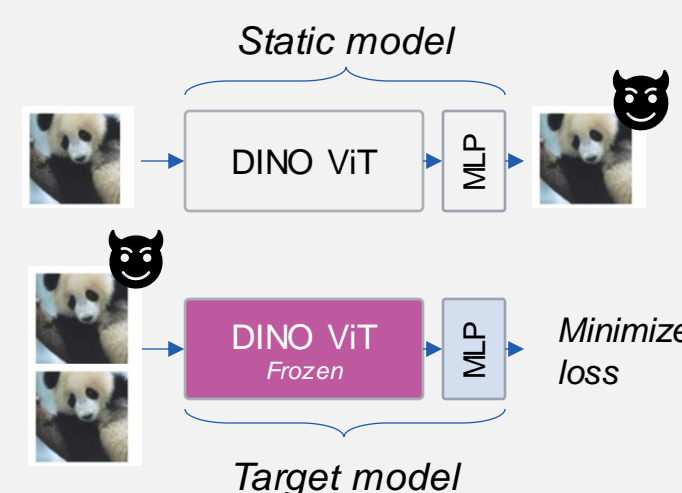
Does it have enough representational power to increase robustness only through fine-tuning?



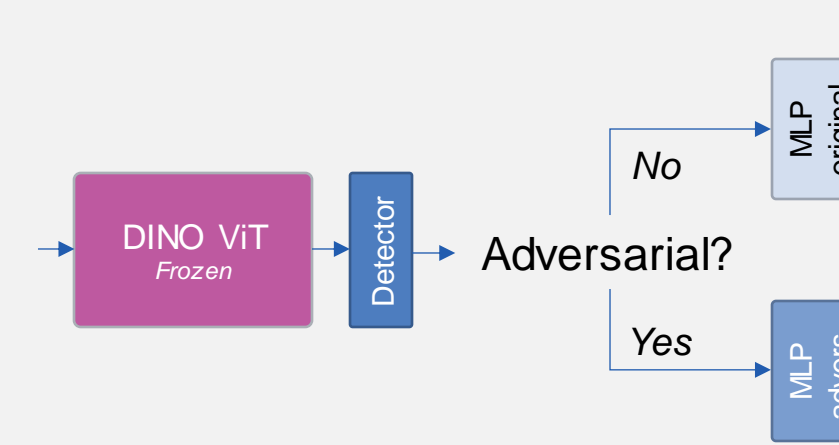
Adversarial Training<sup>2</sup> collapses for large perturbations if the ViT is frozen.



### Ensemble Adv. Training<sup>3</sup>



### Ensemble Specialized Networks



## 2 DINO ViTs Are Not More Robust to Adv. Attacks

DINO		FGSM			PGD			C&W	Clean
		$\epsilon = 0.001$	$\epsilon = 0.03$	$\epsilon = 0.1$	$\epsilon = 0.001$	$\epsilon = 0.03$	$\epsilon = 0.1$	$c = 50$	
DINO	ViT-S/16	52.4%	0.9%	1.1%	49.6%	0.0%	0.0%	0.2%	76.8%
	ViT-B/16	<b>58.9%</b>	1.8%	1.5%	<b>56.8%</b>	0.0%	0.0%	0.4%	77.9%
Superv.	ViT-B/16	55.1%	<b>17.3%</b>	14.5%	47.7%	<b>0.7%</b>	<b>0.1%</b>	0.8%	<b>80.2%</b>
ResNet-50		47.8%	8.0%	<b>24.3%</b>	43.9%	0.1%	0.0%	<b>7.2%</b>	75.7%

Accuracy of ViT models trained using DINO and supervised (Superv.) learning against different white-box adversarial attacks. Metrics are computed on the whole ImageNet validation set. Last column represents accuracy on the original samples from which adversarial images were generated.

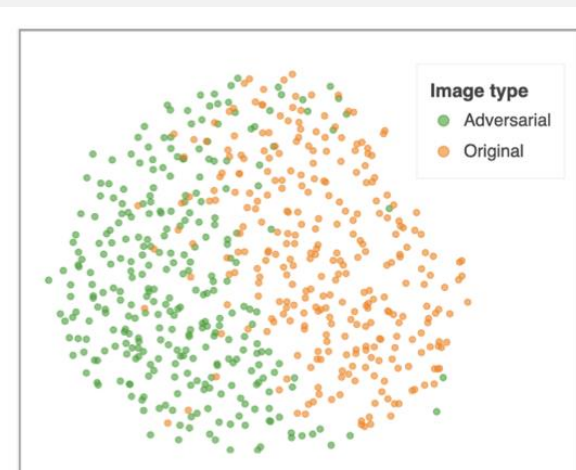
## 3 Self-Supervision Increases Attack Transferability

Self-supervision increases attacks transferability among ViTs and CNNs

Crafted on	Evaluated on			
	ViT-S (D)	ViT-B (D)	ViT-B (S)	ResNet-50
ViT-S (D)	0.0%	12.5%	41.1%	51.3%
ViT-B (D)	5.2%	0.0%	31.4%	48.5%
ViT-B (S)	47.1%	47.8%	0.8%	59.7%
ResNet-50	65.2%	68.4%	74.9%	0.1%

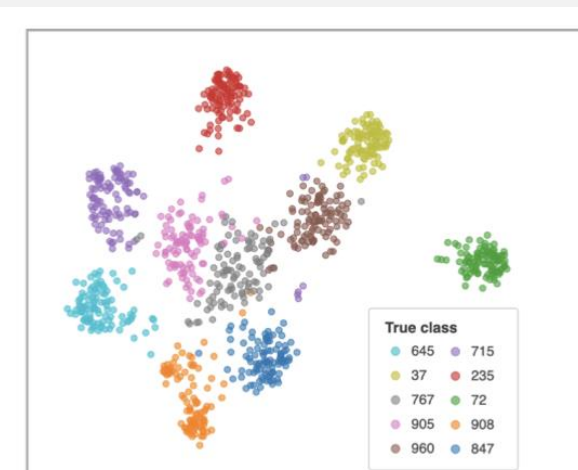
Classification accuracy of adv. samples transferred across architectures. All attacks were crafted using PGD ( $\epsilon = 0.03$ ). Rows represent generation setups and columns, the network used for evaluation. Computed on all validation images from ImageNet-1k. (S) and (D) indicate Supervised and DINO training respectively.

## 4 Latent Space Analysis



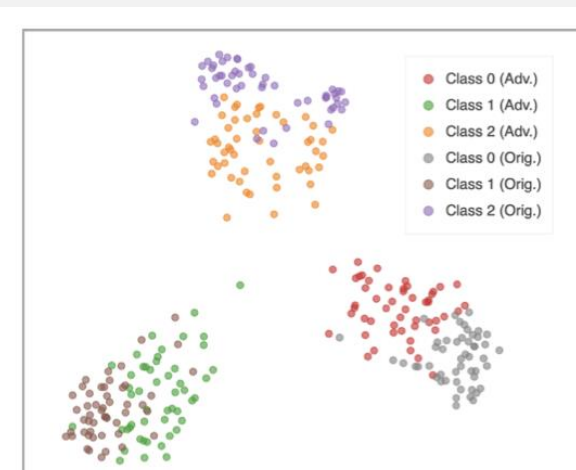
(a) Random images and their adversarial perturbations. Color shows whether a sample is adversarial or original.

Adversarial inputs are linearly separable in the latent space.



(b) Original and adversarial images from 10 random classes. Color represents true label.

They stay close to original samples from their true class.



(c) Adversarial and original images from 3 random classes. Color indicates true class, and if it is an attack.

However, they remain separable within classes.

The latent space may comprise enough information to linearly separate adversarial samples without retraining the ViT.

## 5 Conclusion

- DINO provides **no additional robustness** against white-box adversarial attacks.
- It **increases attacks transferability** among ViTs and CNNs.
- Several defenses may provide **robustness against black-box attacks under limited computational power** only through fine-tuning.

## References

- Caron, M., et al. *Emerging properties in self-supervised vision transformers*. 2021.
- Madry, A., et al. *Towards deep learning models resistant to adversarial attacks*. 2017.
- Tramèr, F., et al. *Ensemble adversarial training: Attacks and defenses*. 2017.