

DINomite: Adversarial Robustness of DINOv3 Vision Transformers

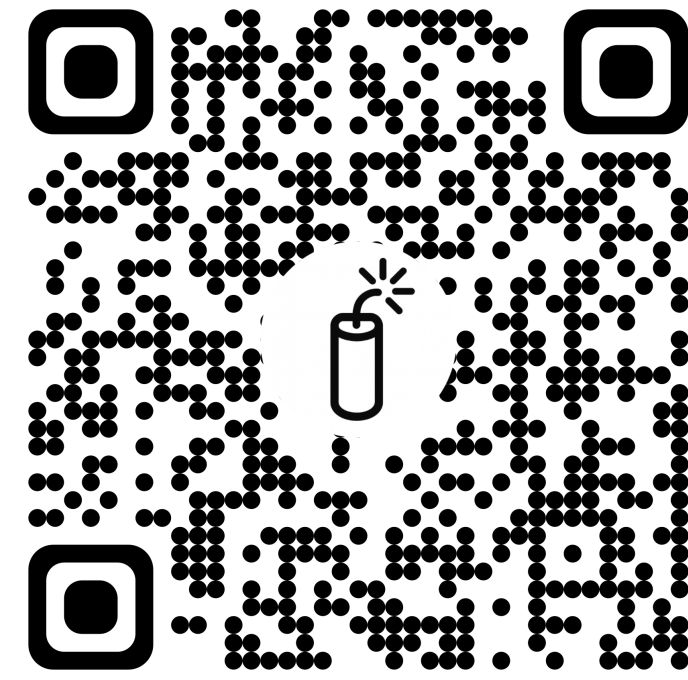


Wojciech Aszkiełowicz* Hubert Berlicki* Jan Burdzicki* Igor Jakus*

w.aszkielowicz@proton.me berlickihubert@gmail.com janburdzicki@gmail.com igorjakus@protonmail.com

Institute of Computer Science, University of Wrocław

*Equal contribution



Abstract

Problem: DINOv3 Vision Transformers achieve state-of-the-art accuracy but remain vulnerable to adversarial attacks.
Contribution: We present a comprehensive evaluation framework for DINOv3, testing robustness against gradient-based and optimization-based attacks. We compare standard fine-tuning against advanced adversarial defense strategies.

Motivation

Why It Matters

- Safety-Critical Systems:** Autonomous driving, medical imaging.
- Physical Threats:** Adversarial patches work in reality.

Research Questions

- Does self-supervised DINOv3 offer inherent robustness?
- Which defense strategy minimizes the accuracy trade-off?

Problem Definition

We formulate the attack as a constrained optimization problem:

$$\max_{\delta} L(f(x + \delta), y) \quad \text{s.t.} \quad \|\delta\|_{\infty} \leq \epsilon$$

Threat Model

- White-box:** Full access to gradients.
- Constraint:** $\epsilon = 8/255$.
- Goal:** Untargeted misclassification.

Experimental Setup

Model Architecture

- Backbone:** DINOv3 ViT-S/16 (Frozen for Classification).
- Head:** Linear Classification Head (Trainable).

Datasets

- CIFAR-10:** Primary benchmark.

Evaluation

- Metric:** Clean vs. Robust Accuracy (PGD-40).
- Epsilon:** $\epsilon \in \{0, \dots, 8/255\}$.
- Attacks:** FGSM, PGD, C&W.
- Defenses:** PGD-AT, TRADES.

Conclusions

Main Takeaways

- Adversarial Training Effectiveness:** DINO-pretrained ViTs are vulnerable to attacks like PGD. However, defenses like PGD-AT or TRADES significantly enhance robustness, mirroring Madry et al.'s findings for CNNs.
- Trade-offs:** Increased robustness introduces a performance penalty on clean data, causing a noticeable drop in standard accuracy.
- Visual Differences:** C&W attacks generate perturbations that are less perceptible to the human eye than those produced by PGD and FGSM.

Future Work

- Evaluating Certified Robustness.
- Scaling to ImageNet-1k and GTSRB.
- Testing diverse attack methods (e.g., AutoAttack).
- Exploring advanced defenses (e.g., MART).

Attack Methods

We focus on three primary attack vectors representing different threat levels:

FGSM (Fast Gradient Sign Method)

A single-step attack assuming linear loss surface.

$$x_{adv} = x + \epsilon \cdot \text{sign}(\nabla_x L(\theta, x, y))$$

Note: Fast, but weak against iterative training.

PGD (Projected Gradient Descent)

The "Universal First-Order Adversary". Iterative FGSM with projection.

$$x^{t+1} = \Pi_{x+\mathcal{S}}(x^t + \alpha \cdot \text{sign}(\nabla_x L(\theta, x^t, y)))$$

Note: The strongest defense benchmark. We use 40 steps.

C&W (Carlini & Wagner)

Optimization-based attack minimizing distance and loss term.

$$\min_{\delta} \|\delta\|_2 + c \cdot f(x + \delta)$$

Note: Finds adversarial examples with minimal visible perturbation.

Adversarial Examples

C&W (L_2) - Minimal Perturbation

Original (left) vs Adversarial (right)



Logits: -3.50 -0.42 5.69 -1.71 2.82 -0.34 -4.72 -0.32 -3.19 -1.11
Classes: airplane automobile bird cat deer dog frog horse ship truck



Logits: -6.59 -4.46 -1.97 1.48 2.11 2.38 0.14 1.57 -0.08 -2.79
Classes: airplane automobile bird cat deer dog frog horse ship truck

FGSM ($\epsilon = 8/255$) - Patterned Noise

Original (left) vs Adversarial (right)



Logits: -4.96 -3.32 -1.14 0.11 0.07 0.46 0.14 1.10 0.84 -0.07
Classes: airplane automobile bird cat deer dog frog horse ship truck



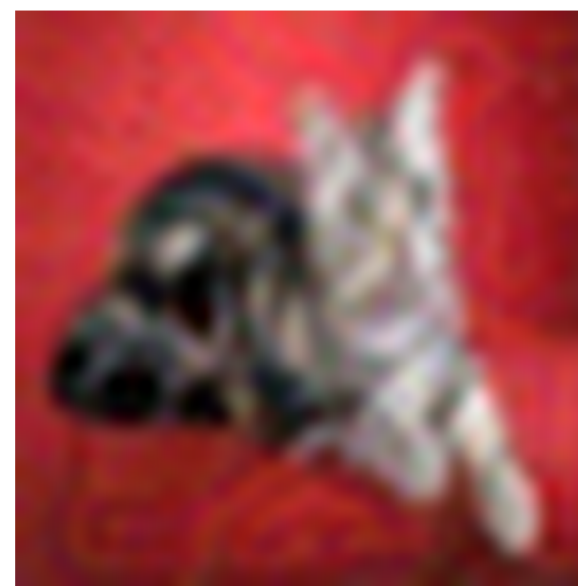
Logits: -4.81 -3.10 -1.27 3.44 0.04 1.38 -0.47 0.88 -0.84 -0.00
Classes: airplane automobile bird cat deer dog frog horse ship truck

PGD ($\epsilon = 8/255$) - High Noise

Original (left) vs Adversarial PGD (right)



Logits: -3.13 -3.53 -3.81 -0.14 -4.74 3.30 -4.78 -2.40 -4.89 -5.08
Classes: airplane automobile bird cat deer dog frog horse ship truck



Logits: -4.46 -4.04 -2.02 -7.29 -4.09 10.25 -1.93 1.11 -6.51 -0.68
Classes: airplane automobile bird cat deer dog frog horse ship truck

Defense Strategies

Standard training yields 0% robustness. We evaluate two state-of-the-art defenses:

PGD-AT (Adversarial Training)

A min-max game formulation (Madry et al.). The inner loop generates strong adversarial examples, while the outer loop updates the model to resist them.

$$\min_{\theta} \mathbb{E}_{(x,y) \sim D} \left[\max_{\|\delta\|_{\infty} \leq \epsilon} L(f_{\theta}(x + \delta), y) \right]$$

Key Idea: Train directly on the worst-case examples.

TRADES

Separates loss into natural accuracy and robustness regularization (Zhang et al.). It minimizes KL-divergence between predictions on clean and adversarial inputs.

$$\min_{\theta} \mathbb{E} \left[\underbrace{L(f(x), y)}_{\text{Natural}} + \beta \cdot \underbrace{\text{KL}(f(x) \| f(x + \delta))}_{\text{Robustness}} \right]$$

Key Idea: Enforce smoothness of the decision boundary.

Defense Results

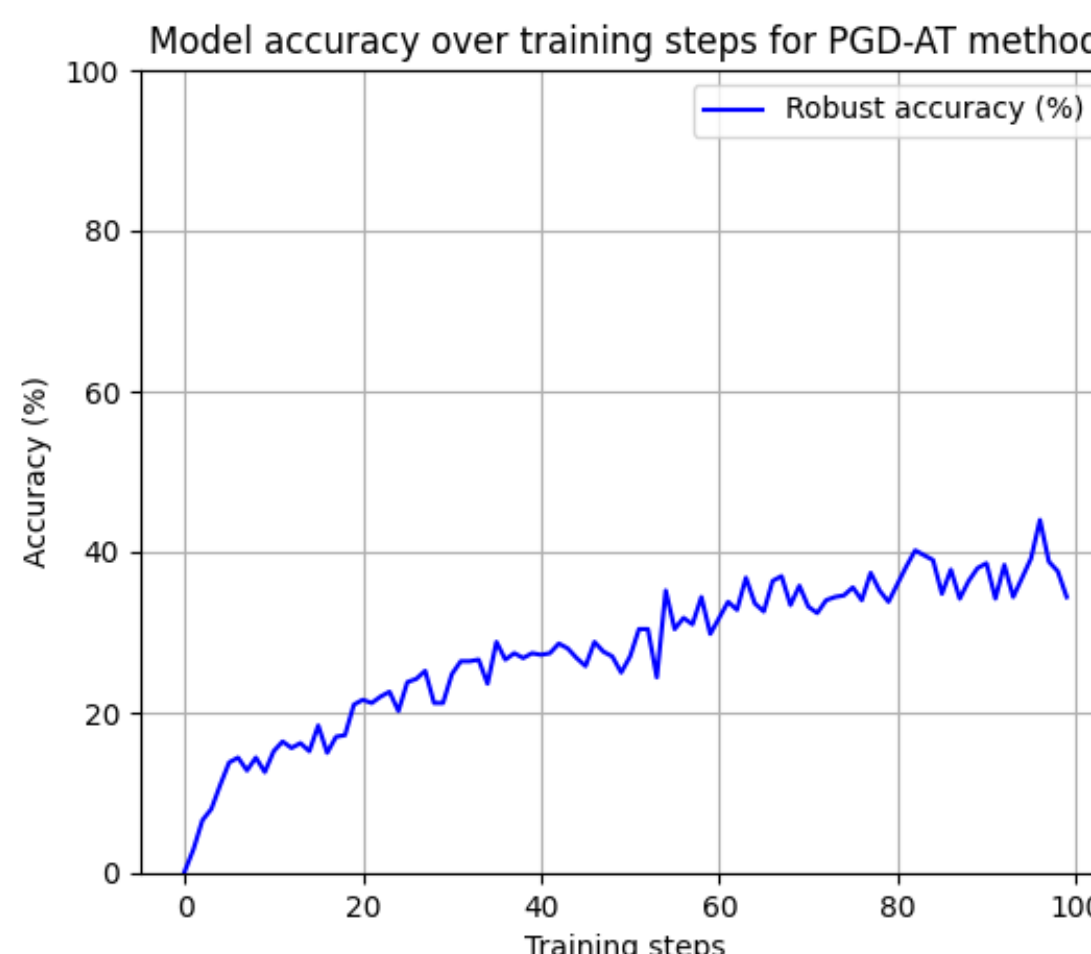
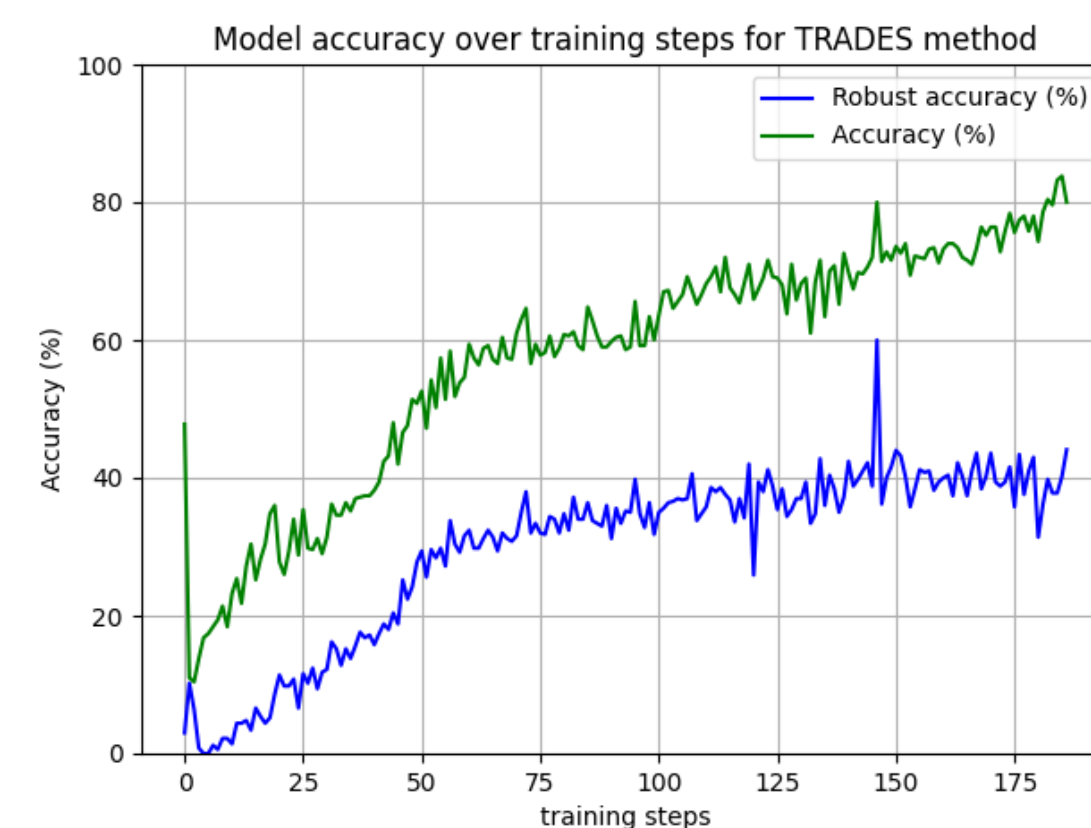


Figure 1: Robustness and accuracy across defense methods.

Quantitative Results

Comparison on CIFAR-10 ($\epsilon = 8/255$).

Model	Clean	PGD
Standard	96.56%	0.00%
PGD-AT	77.19%	40.62%
TRADES	81.25%	35.94%

References

- Goodfellow et al. (2015). *Explaining and Harnessing Adversarial Examples*. arXiv:1412.6572
- Moosavi-Dezfooli et al. (2016). *DeepFool: A Simple and Accurate Method to Fool Deep Neural Networks*. arXiv:1511.04599
- Madry et al. (2018). *Towards Deep Learning Models Resistant to Adversarial Attacks*. arXiv:1706.06083
- Carlini & Wagner (2017). *Towards Evaluating the Robustness of Neural Networks*. arXiv:1608.04644
- Zhang et al. (2019). *Theoretically Principled Trade-off between Robustness and Accuracy*. arXiv:1901.08573
- Kurakin et al. (2017). *Adversarial Examples in the Physical World*. arXiv:1607.02533
- Wang et al. (2020). *Improving Adversarial Robustness Requires Revisiting Data Augmentation and Training*. OpenReview: rklOg6EFwS