# ASTrA: Adversarial Self-supervised Training with Adaptive-Attack

ICLR 2025 Singapore

*https://prakashchhipa.github.io/projects/ASTrA*

Prakash Chandra Chhipa[1]*, Gautam Vashishtha[2]*, Settur Jithamanyu[3]*, Rajkumar Saini[1], Mubarak Shah[4], Marcus Liwicki[1]

[1]Machine Learning Group, Luleå Tekniska Universitet, Sweden
[2]Indian Institute of Technology, Gandinagar
[3]Indian Institute of Technology, Madras
[4]Center For Research in Computer Vision, University of Central Florida, USA

*equal contribution
presenter

# Self-supervised adversarial attacks – limitation

o Networks **vulnerability** to adversarial examples.



| cat | δ=8/255. | airliner |

o Limitation: **Hand-crafted** adversarial attack strategy fail to adapt dynamically in Self-supervised adversarial training (Self-AT).

- o Does not align with model's learning dynamics
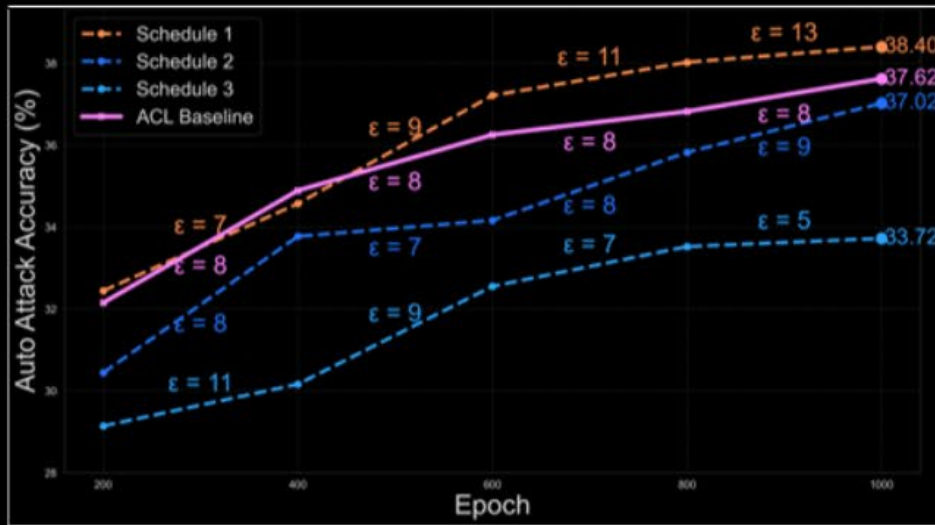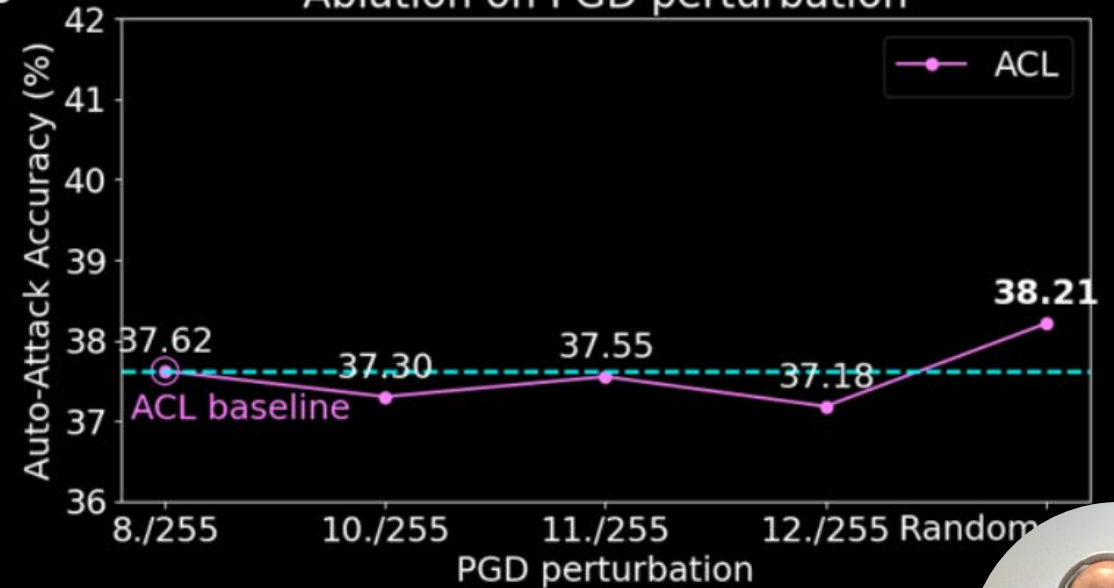- o No correspondence between training examples and attack strategy parameters

# Towards goal

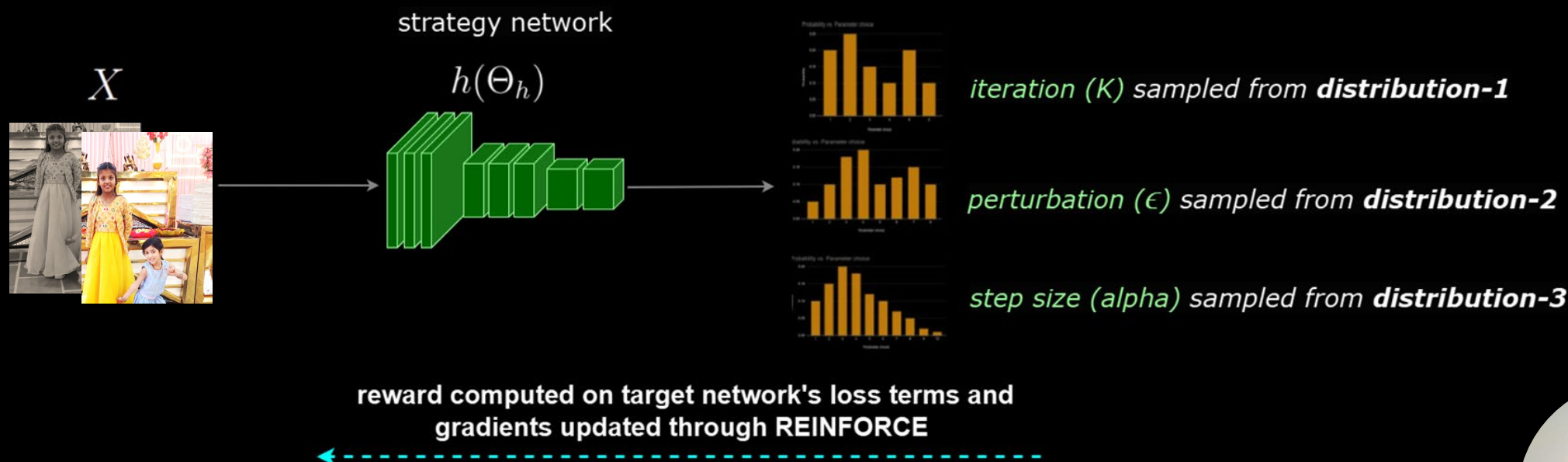*ACL Method - Robust Pre-Training by Adversarial Contrastive Learning, NeurIPS 2020 (on CIFAR10)*



Develop adaptive, self-supervised adversarial attack strategy

# Learnable attacks in ASTrA

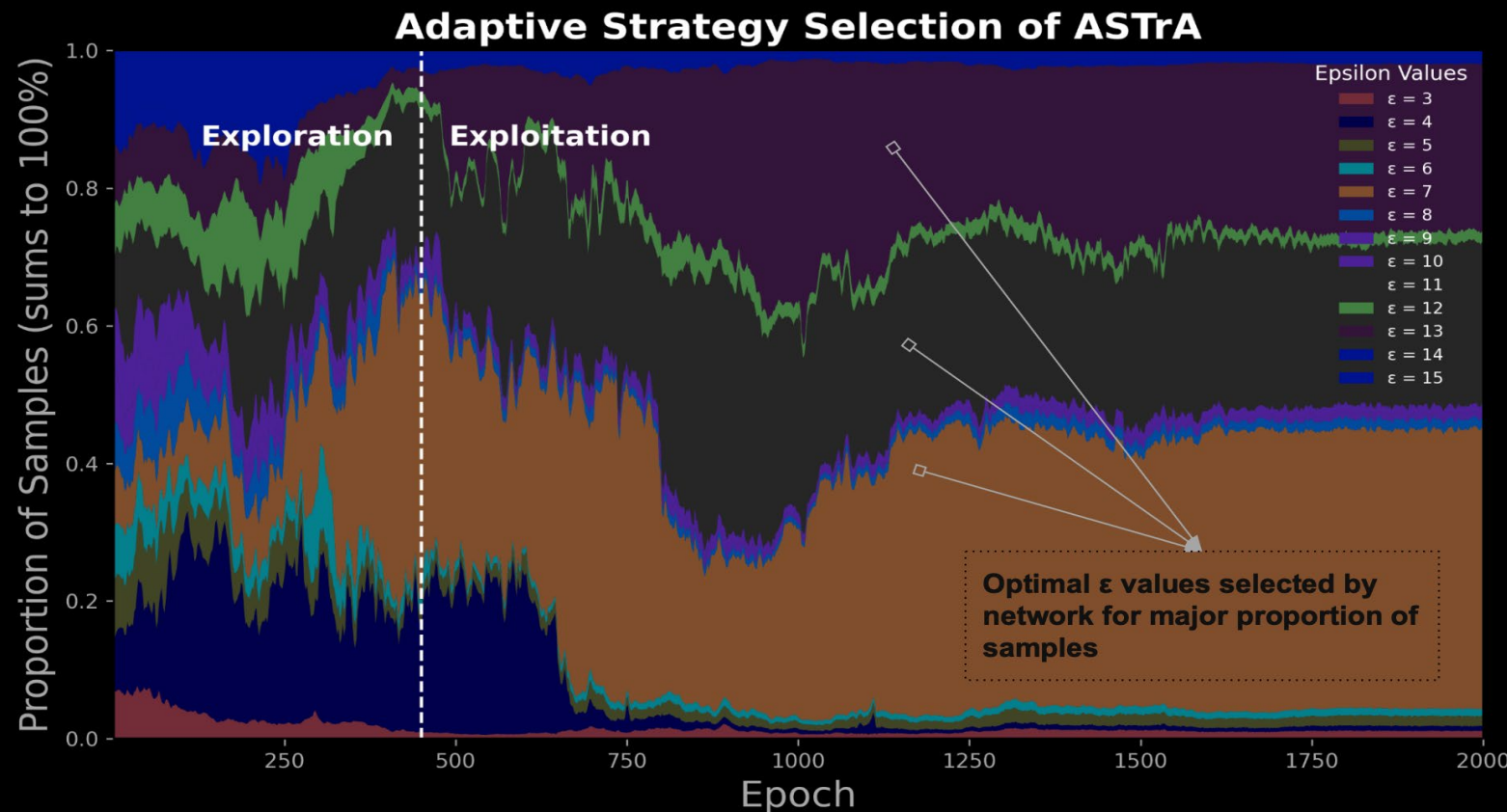✓ **Learnable strategy** network autonomously finds optimal attacks.

strategy network

$h(\Theta_h)$

$X$

*iteration (K) sampled from* **distribution-1**

*perturbation (ε) sampled from* **distribution-2**

*step size (alpha) sampled from* **distribution-3**

reward computed on target network's loss terms and gradients updated through REINFORCE

*ASTrA: Adversarial Self-supervised Training with Adaptive-Attack, ICLR 2025*

*https://prakashchhipa.github.io/projects/ASTrA*

ASTrA framework

✓ Exploration-Exploitation using SSL contrastive reward and REINFORCE optimization

*ASTrA: Adversarial Self-supervised Training with Adaptive-Attack, ICLR 2025*

*https://prakashchhipa.github.io/projects/ASTrA*
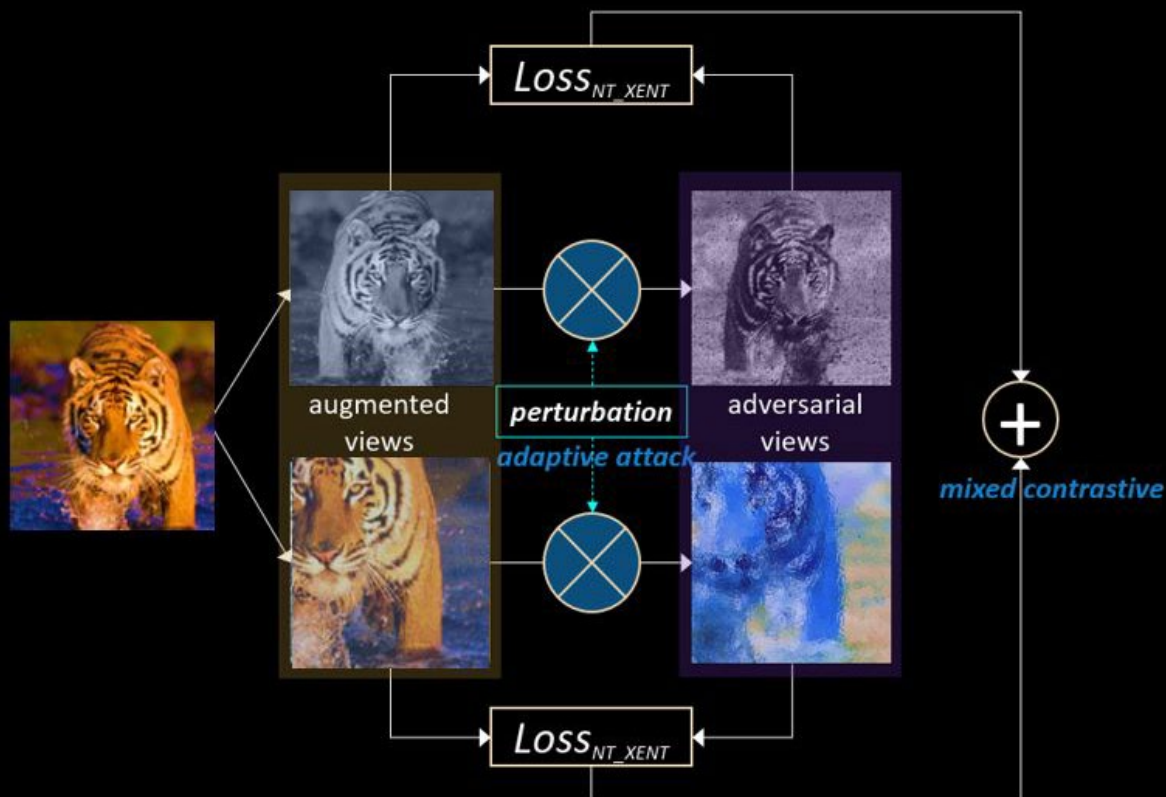
# ASTrA framework

✓ **Exploration-Exploitation** using SSL contrastive reward and **REINFORCE** optimization



**Self-AT Loss** for Target

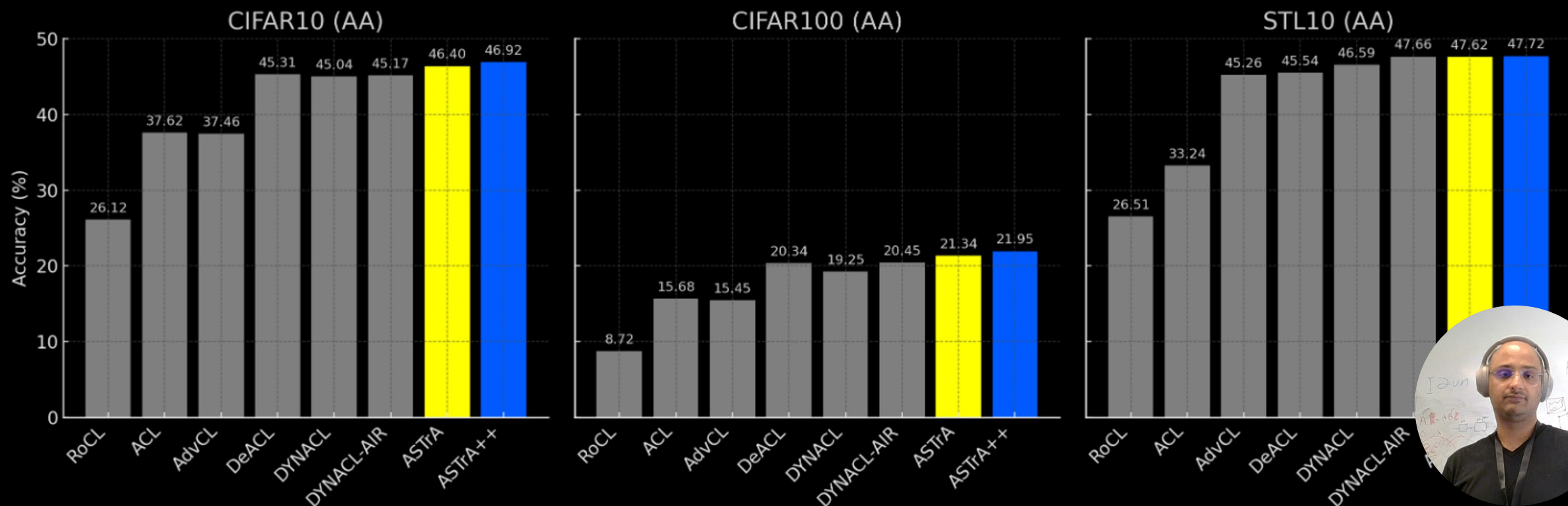$$Loss_{NT\_CLEAN} + Loss_{NT\_MIXED} + Loss_{NT\_ADV}$$

*minimization*

**SSL Reward** for Strategy

$$Loss_{NT\_ADV} - Loss_{NT\_CLEAN}$$

*maximization*

*ASTrA: Adversarial Self-supervised Training with Adaptive-Attack, ICLR 2025*

https://prakashchhipa.github.io/projects/ASTrA

# Mixed contrastive loss

✓ Align representations using of clean view to corresponding (adaptively attacked) perturbed view.

Results *public benchmarks*

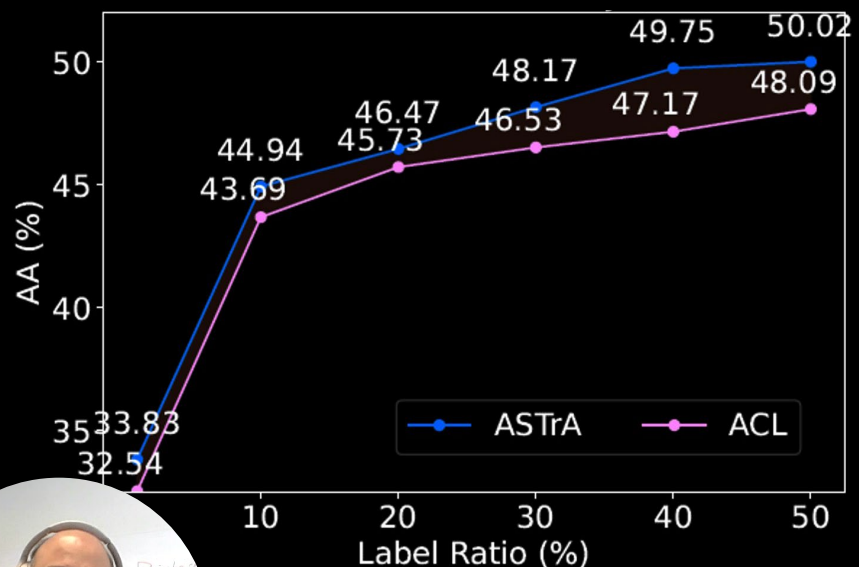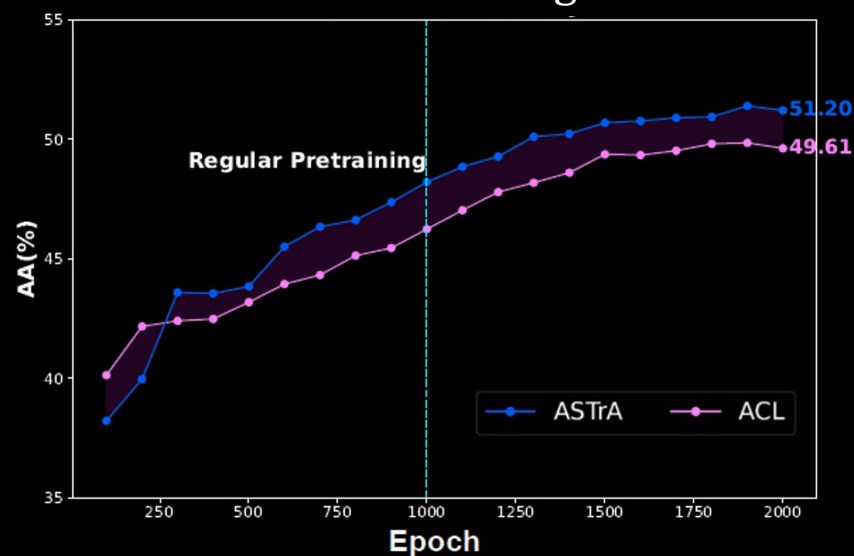**Standard Linear Finetuning Performance – ASTrA vs. other Self-AT method**

*ASTrA: Adversarial Self-supervised Training with Adaptive-Attack, ICLR 2025*

*https://prakashchhipa.github.io/projects/ASTrA*

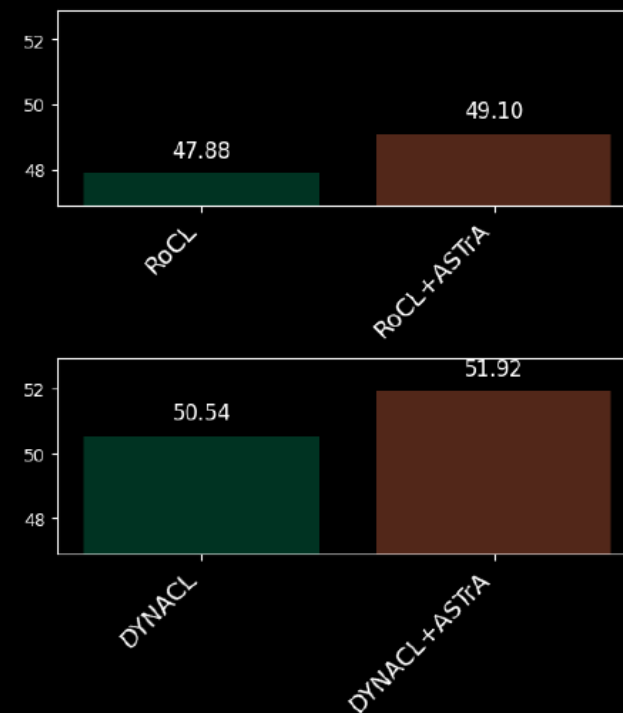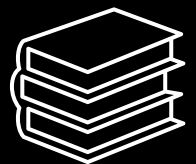# Results *label efficiency, improved robustness and modularity*

**semi-supervised setting**

**robust overfitting**

**plug-and-play with self-ATs**

🙏 Thank you

*ASTrA: Adversarial Self-supervised Training with Adaptive-Attack, ICLR 2025*

*https://prakashchhipa.github.io/projects/ASTrA*