# On Adversarial Training without Perturbing All Examples

MAX PLANCK INSTITUTE FOR INFORMATICS

Max Losch[MPII], Mohamed Omran[MPII], David Stutz[MPII], Mario Fritz[CISPA], Bernt Schiele[MPII]

CISPA HELMHOLTZ CENTER FOR INFORMATION SECURITY

## In a Nutshell



**Subset Adv. Training (SAT)**

**Evaluation**

**Robust accuracy**

Vanilla adversarial training (AT) and most its variants perturb **every training example**. To what extent is that necessary? We split the training set into subsets $A$ and $B$, train on $A \cup B$ but construct adv. examples only for examples in $A$.

## Contributions

We propose an analytical tool *Subset Adversarial Training* (SAT) **1** to investigate robustness when only a training subset has been attacked.

**2** Adv. robustness transfers to never attacked classes

**3 & 4** Harder examples tend to provide best robustness transfer

**4** Attacking $50\%$ of training data is sufficient to recover baseline robust accuracy

**5** $30\%$ reach baseline robust accuracy after transfer to downstream tasks

**6** Can be combined with single-step attack training

## Paper and Code: github.com/mlosch/SAT

[1] Dan Hendrycks et al. "Natural adversarial examples". In: *CVPR* (2021).
[2] Eric Wong, Leslie Rice, and J Zico Kolter. "Fast is better than free: Revisiting adversarial training". In: *ICLR* (2020).
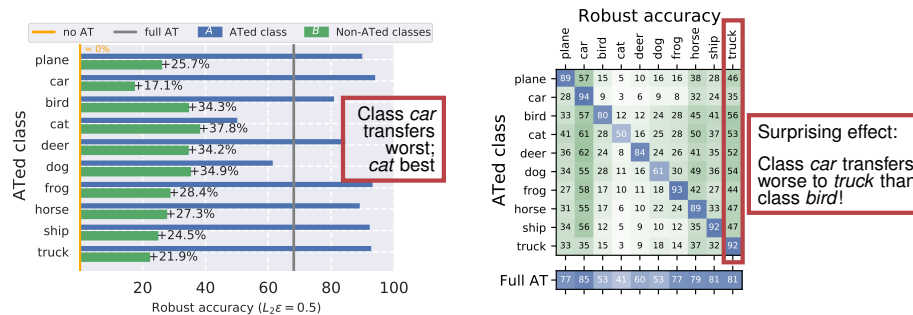
## 1 Subset Adversarial Training (SAT)

**Questions:** Does robustness transfer to unseen classes/examples? Does it depend on particular classes/examples?

$$\text{SAT:} \quad \min_{\theta} \mathbb{E}_{(x,y) \in \mathcal{D}_{\text{train}}} \left[ \underbrace{w_A \mathbb{1}_{(x,y) \in A} \max_{||\delta||_2 \leq \epsilon} \mathcal{L}(x+\delta, y; \theta)}_{\text{Adv. Training loss } \forall x \in A} + \underbrace{w_B \mathbb{1}_{(x,y) \in B} \mathcal{L}(x, y; \theta)}_{\text{Vanilla cross-entropy } \forall x \in B} \right]$$

**Note:** $A$ and $B$ are *fixed* pre training    Examples in $B$ are never attacked

## 2 Class-subset Splits (CSAT)



Class *car* transfers worst; *cat* best

Robust accuracy

Surprising effect: Class *car* transfers worse to *truck* than class *bird*!
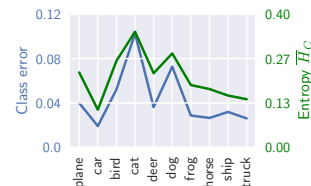
▶ Non-trivial robust accuracy on classes in $B$
▶ Characteristics correlate strongly with class *difficulty*
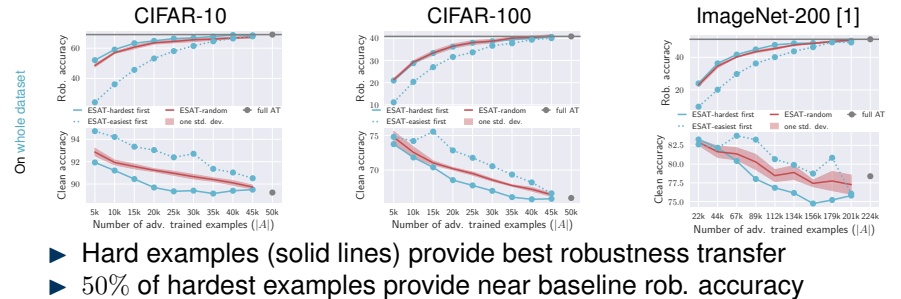
## 3 Measuring Class Difficulty

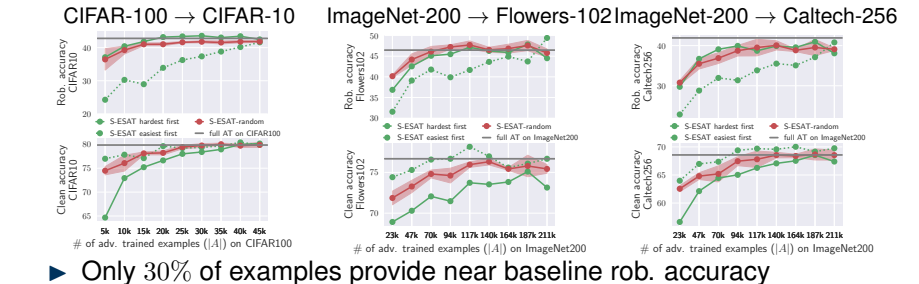As class difficulty metric, we utilize entropy $\mathcal{H}$ over softmax $\sigma$. We rank examples once before training.

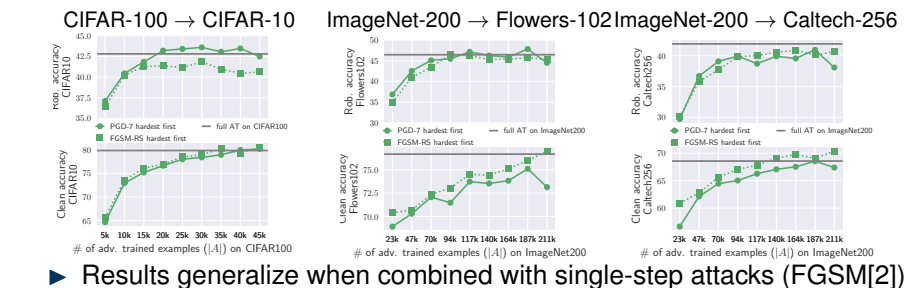$$\mathcal{H}(f(x)) = -\sum_{i=1}^{N} \sigma_i(f(x)) \cdot \log \sigma_i(f(x))$$



## 4 Example-subset Splits (ESAT)



▶ Hard examples (solid lines) provide best robustness transfer
▶ $50\%$ of hardest examples provide near baseline rob. accuracy

## 5 Transfer to Downstream Tasks (S-ESAT)



▶ Only $30\%$ of examples provide near baseline rob. accuracy

## 6 Single-step S-ESAT



▶ Results generalize when combined with single-step attacks (FGSM[2])