

# Towards Robust Foundation Models

## Efficient and Effective Adversarial Contrastive Learning

Xilie Xu

Ph.D. student, School of Computing, NUS

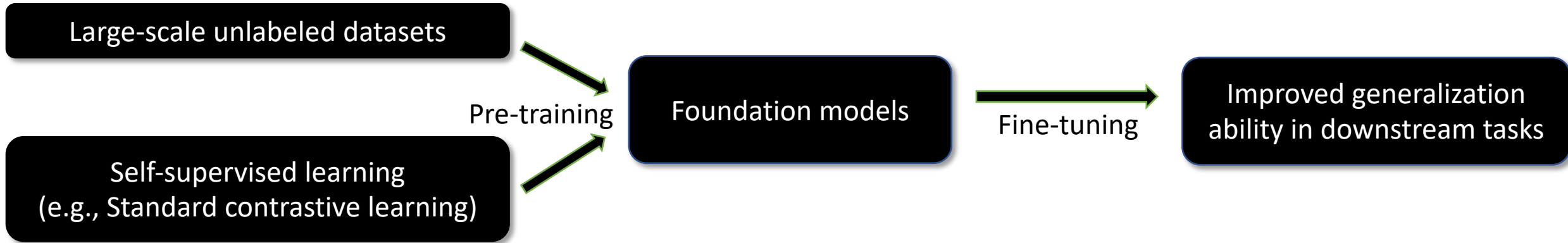
28<sup>th</sup> Nov 2023

Advisor: Prof. Mohan Kankanhalli

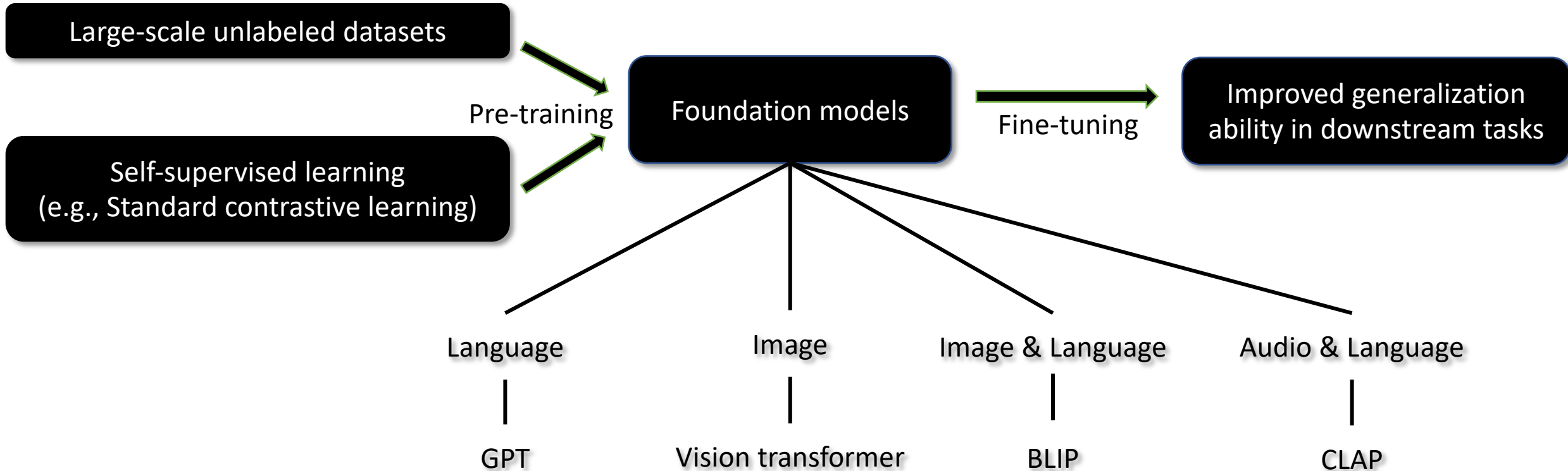
# Towards Robust Foundation Models

Efficient and Effective Adversarial Contrastive Learning

# Foundation Model



# Foundation Model



# Towards Robust Foundation Models

Efficient and Effective Adversarial Contrastive Learning

# Security Risk--Adversarial Attack

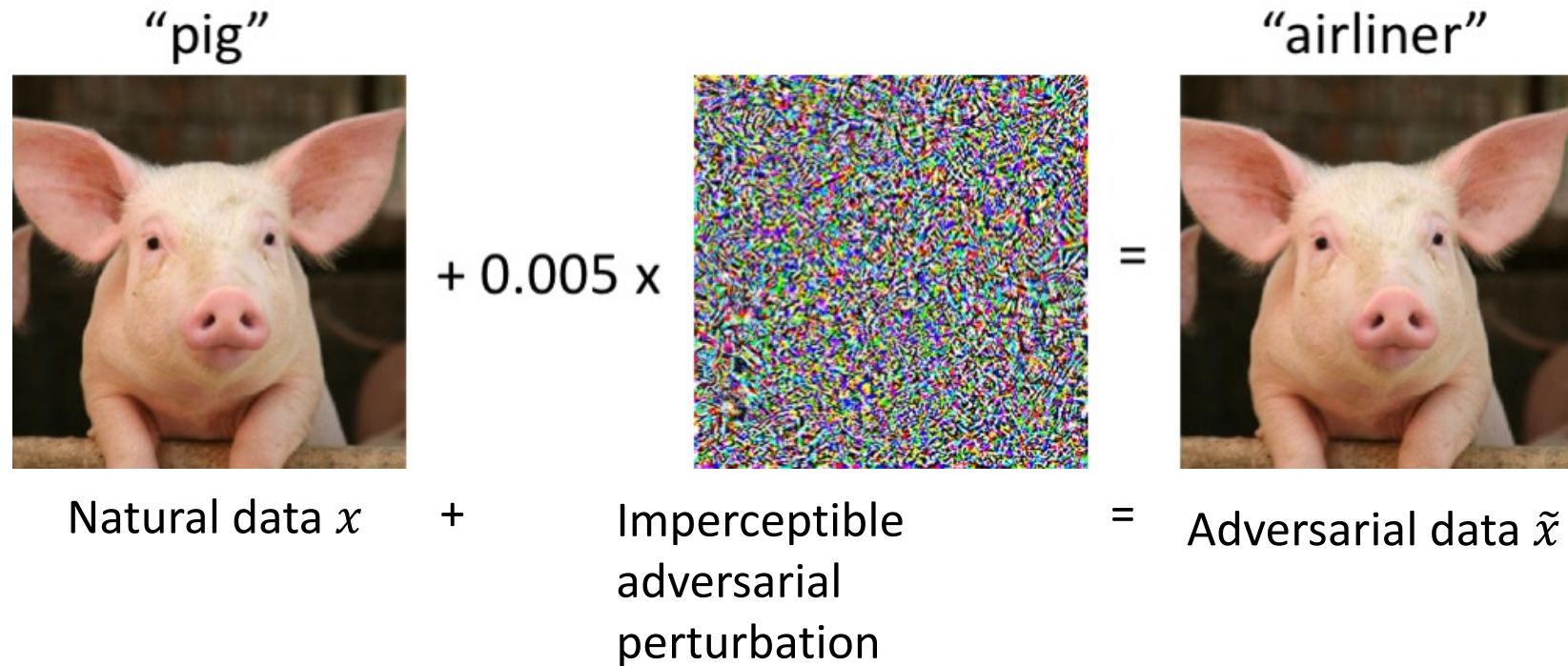
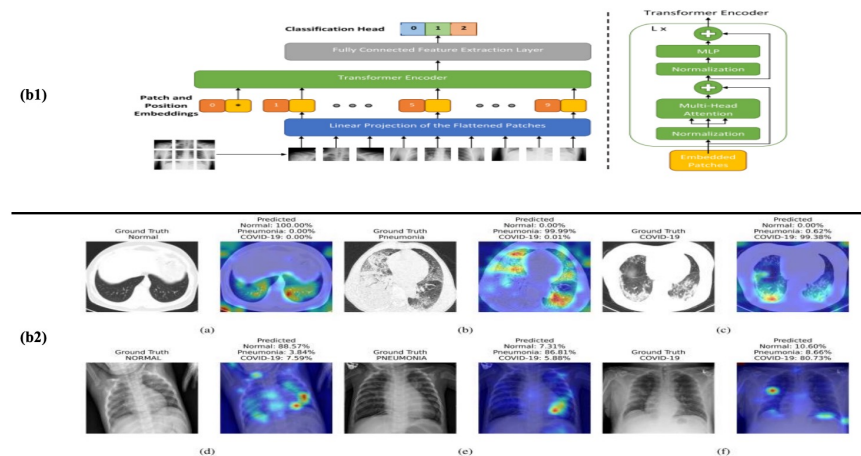


Image from [https://gradientscience.org/intro\\_adversarial/](https://gradientscience.org/intro_adversarial/)

# Security Risk--Adversarial Attack

Potential security risks when applying foundation models to safety-critical tasks

Medical diagnosis



[Henry et al., ArXiv 2022]

Traffic sign recognition

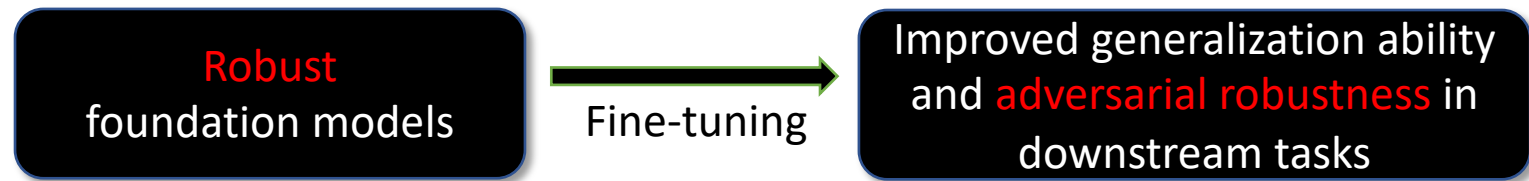


Image from <https://www.google.com/url?sa=i&url=https%3A%2F%2Fwww.itm-p.com%2Fprotect-iot-applications-from-adversarial-evasion-attacks&psig=AOvVaw2ltdTWXQ51FptX2EFu9gKr&ust=1700842535801000&source=images&cd=vfe&opi=89978449&ved=0CBQQjhXqFwoTCJidm6vC2oiDFQAAAAAdAAAAABAJ>

# Risks Urge **Robust** Foundation Models

Robust foundation models should:

1. be generalizable to downstream tasks;
2. **be robust against adversarial attacks.**

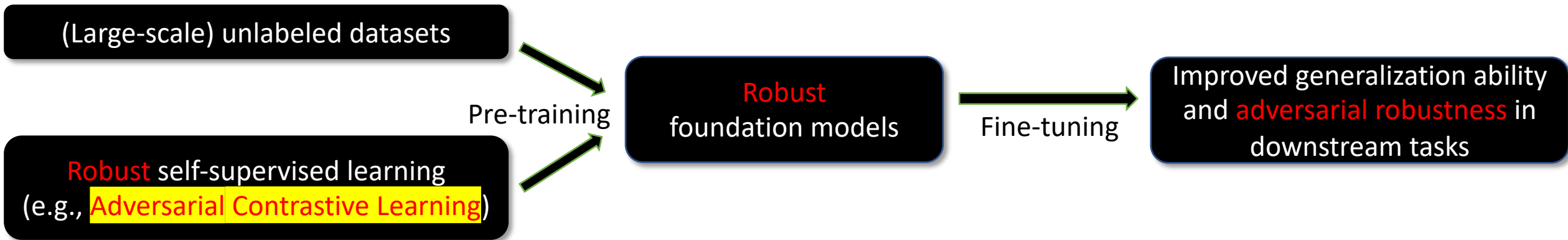




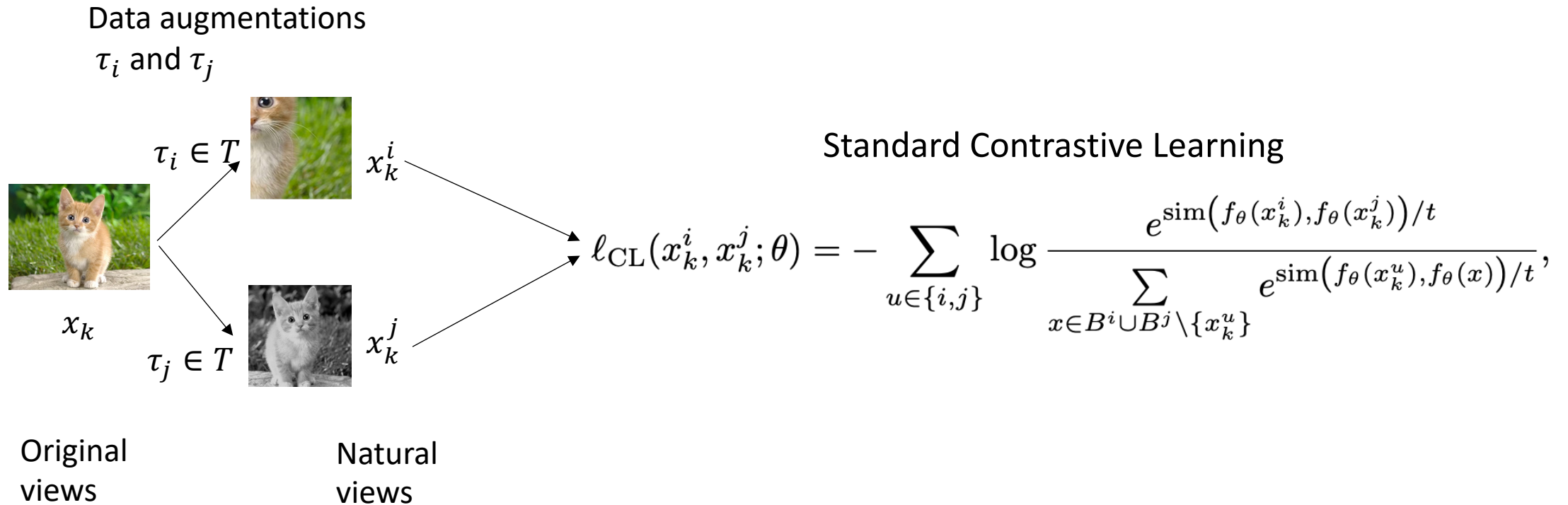
# Risks Urge **Robust** Foundation Models

Robust foundation models should:

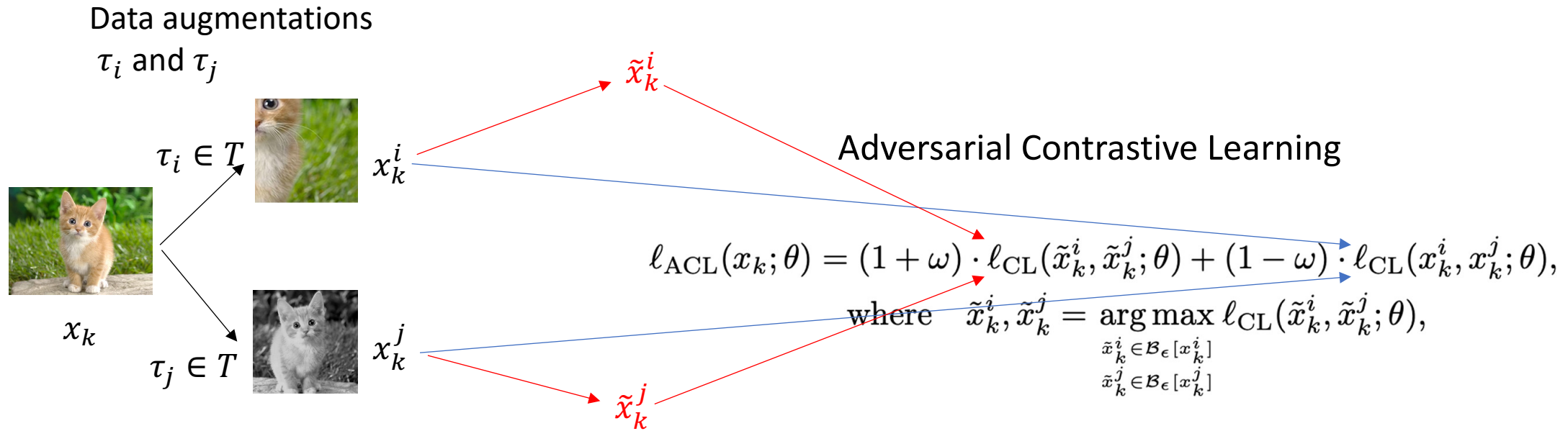
1. be generalizable to downstream tasks;
2. **be robust against adversarial attacks.**



# Adversarial Contrastive Learning (ACL)



# Adversarial Contrastive Learning (ACL)



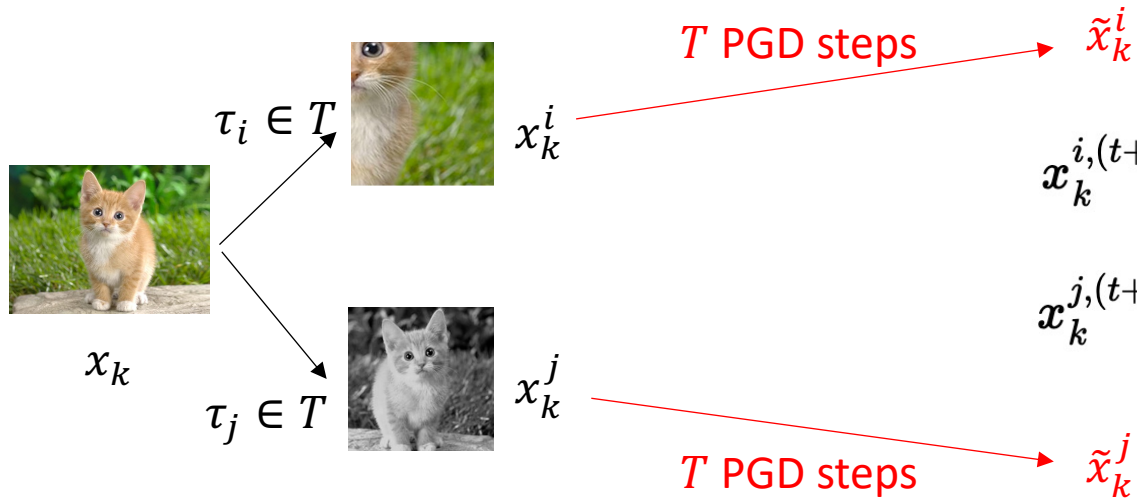
At each epoch, the ACL alternatively conducts Step (1) and (2):

- Step (1): Generating adversarial views
- Step (2): Updating parameters via minimizing the contrastive loss on the natural views and adversarial views.

# Towards Robust Foundation Models

Efficient and Effective Adversarial Contrastive Learning

# Motivation: ACL is inefficient due to $T$ PGD steps



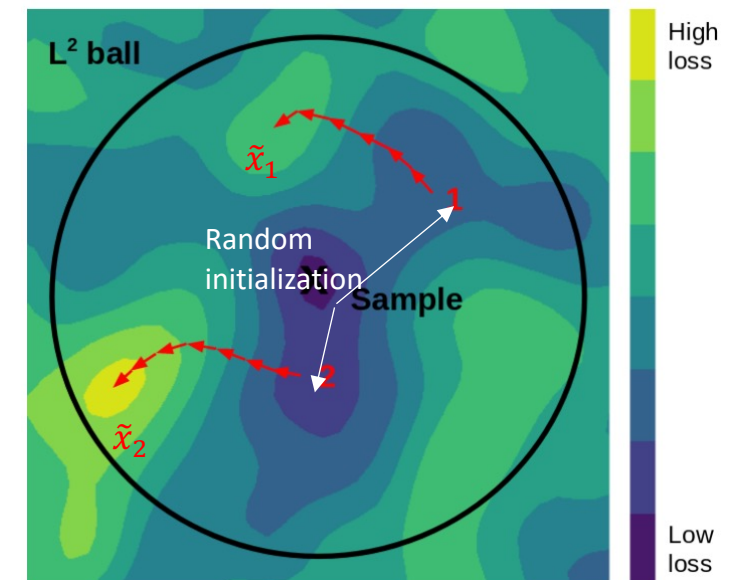
$$\mathbf{x}_k^{i,(t+1)} = \Pi_{\mathcal{B}_\epsilon[\mathbf{x}_k^{i,(0)}]} \left( \mathbf{x}_k^{i,(t)} + \rho \cdot \text{sign}(\nabla_{\mathbf{x}_k^{i,(t)}} \ell_{\text{CL}}(\mathbf{x}_k^{i,(t)}, \mathbf{x}_k^{j,(t)})) \right)$$

$$\mathbf{x}_k^{j,(t+1)} = \Pi_{\mathcal{B}_\epsilon[\mathbf{x}_k^{j,(0)}]} \left( \mathbf{x}_k^{j,(t)} + \rho \cdot \text{sign}(\nabla_{\mathbf{x}_k^{j,(t)}} \ell_{\text{CL}}(\mathbf{x}_k^{i,(t)}, \mathbf{x}_k^{j,(t)})) \right)$$

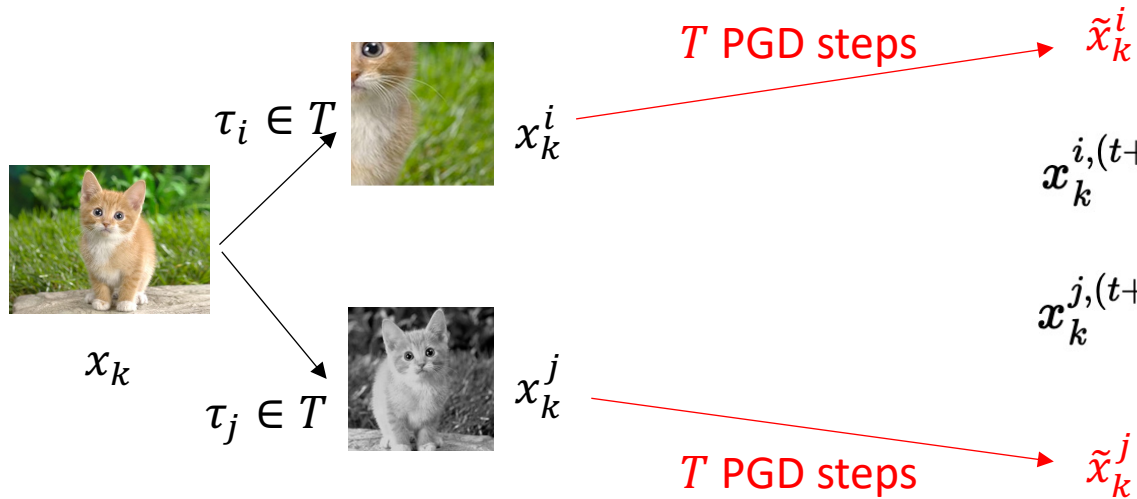
ACL on the entire training set is extremely time-consuming.

- CIFAR-10: about *43 hours*
- ImageNet-1K: about *650 hours*

Project Gradient Descent (PGD)



# Motivation: ACL is inefficient due to $T$ PGD steps



$$\mathbf{x}_k^{i,(t+1)} = \Pi_{\mathcal{B}_\epsilon[\mathbf{x}_k^{i,(0)}]} \left( \mathbf{x}_k^{i,(t)} + \rho \cdot \text{sign}(\nabla_{\mathbf{x}_k^{i,(t)}} \ell_{\text{CL}}(\mathbf{x}_k^{i,(t)}, \mathbf{x}_k^{j,(t)})) \right)$$

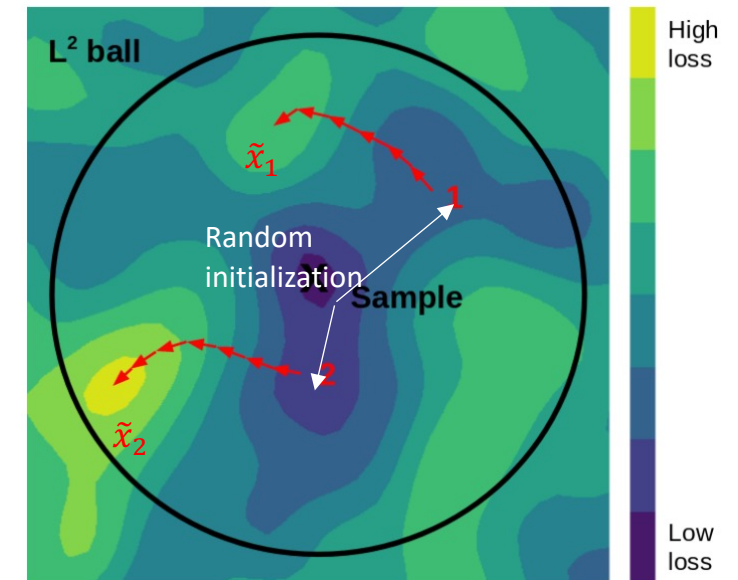
$$\mathbf{x}_k^{j,(t+1)} = \Pi_{\mathcal{B}_\epsilon[\mathbf{x}_k^{j,(0)}]} \left( \mathbf{x}_k^{j,(t)} + \rho \cdot \text{sign}(\nabla_{\mathbf{x}_k^{j,(t)}} \ell_{\text{CL}}(\mathbf{x}_k^{i,(t)}, \mathbf{x}_k^{j,(t)})) \right)$$

ACL on the entire training set is extremely time-consuming.

- CIFAR-10: about *43 hours*
- ImageNet-1K: about *650 hours*

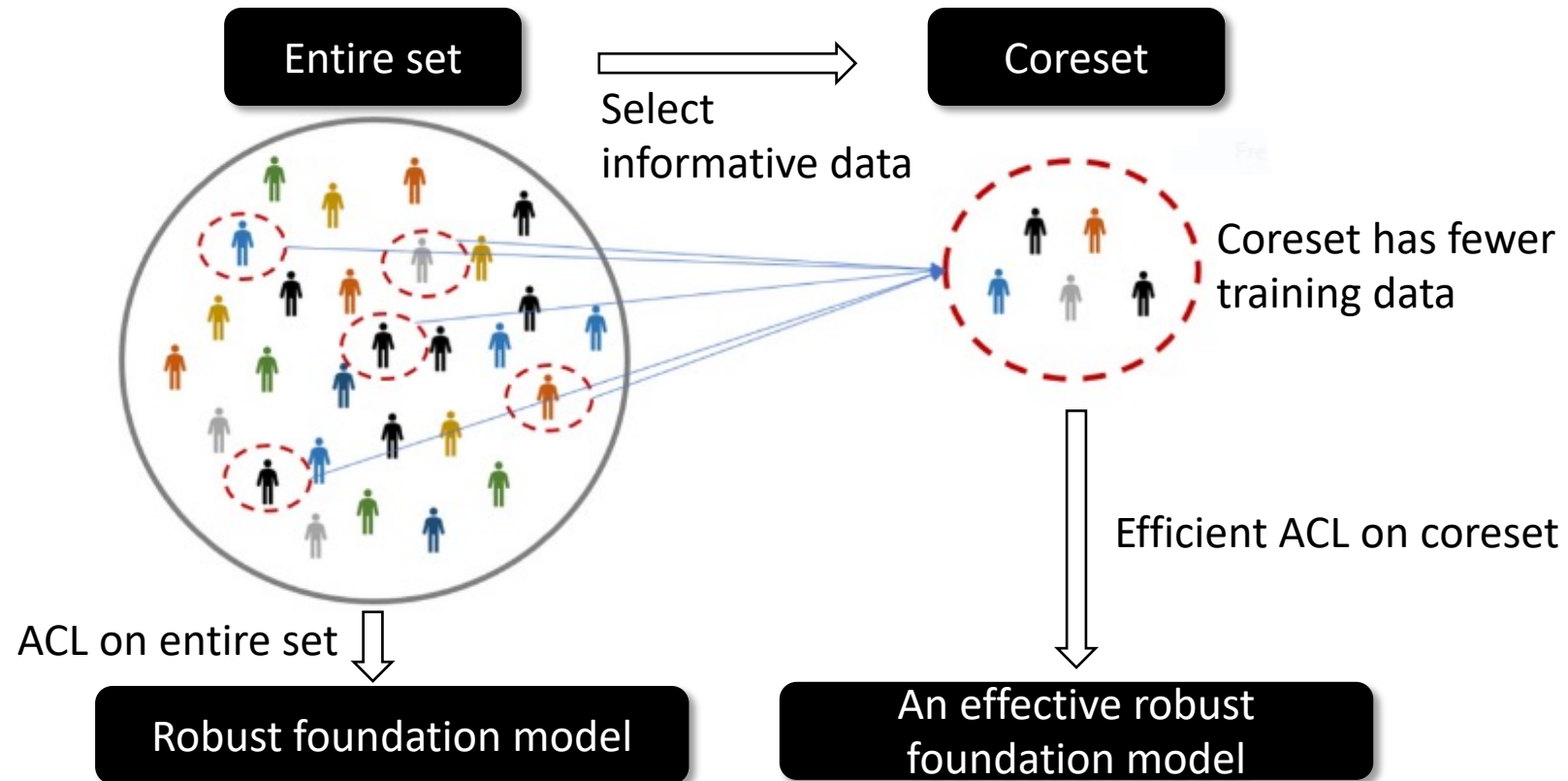
## How can we speed up ACL?

Project Gradient Descent (PGD)



# Robustness-Aware Coreset Selection (RCS)

- Intuitive idea: Find an informative training subset (called “coreset”)
  - Decreasing the number of training samples
  - Guaranteeing the model to effectively learn robust representations



# Robustness-Aware Coreset Selection (RCS)

- Intuitive idea: Find an informative training subset (called “coreset”)
- Objective function of RCS

ACL on CIFAR-10	SLF on CIFAR-10	SLF on CIFAR-10	
	RD loss (lower is better)	SA (%)	RA (%)
Entire	0.1243	78.87	39.19
Random-0.05	0.3357	67.45	22.96

Representational  
divergence (RD)

$$\ell_{\text{RD}}(x; \theta) = d(g \circ f_{\theta}(\tilde{x}), g \circ f_{\theta}(x)) \quad \text{s.t.} \quad \tilde{x} = \arg \max_{x' \in \mathcal{B}_{\epsilon}[x]} d(g \circ f_{\theta}(x'), g \circ f_{\theta}(x))$$

RD, without using labels, measures the adversarial robustness.



# Robustness-Aware Coreset Selection (RCS)

- Intuitive idea: Find an informative training subset (called “coreset”)
- Objective function of RCS

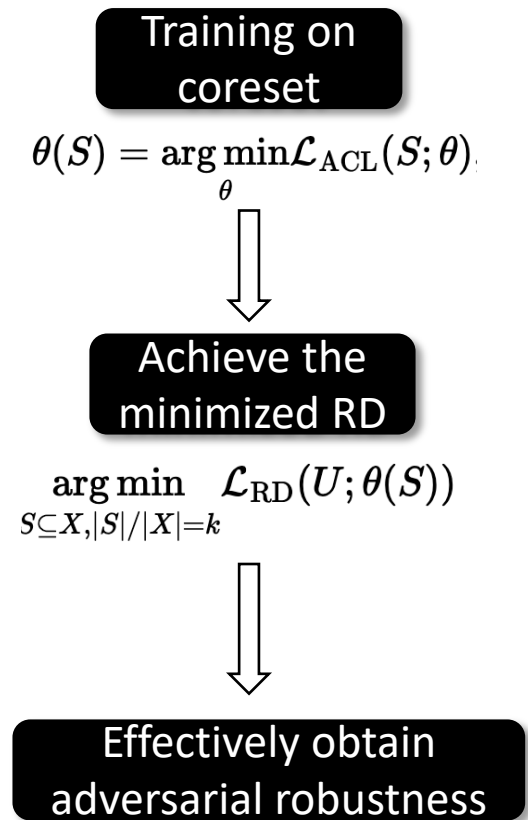
Coreset  $S^* = \arg \min_{S \subseteq X, |S|/|X|=k} \mathcal{L}_{RD}(U; \theta(S))$ 
Subset fraction
Unlabeled validation set

Representational divergence (RD)

$\theta(S) = \arg \min_{\theta} \mathcal{L}_{ACL}(S; \theta),$ 
Adversarial contrastive loss

$\ell_{RD}(x; \theta) = d(g \circ f_{\theta}(\tilde{x}), g \circ f_{\theta}(x)) \quad \text{s.t.} \quad \tilde{x} = \arg \max_{x' \in \mathcal{B}_{\epsilon}[x]} d(g \circ f_{\theta}(x'), g \circ f_{\theta}(x))$

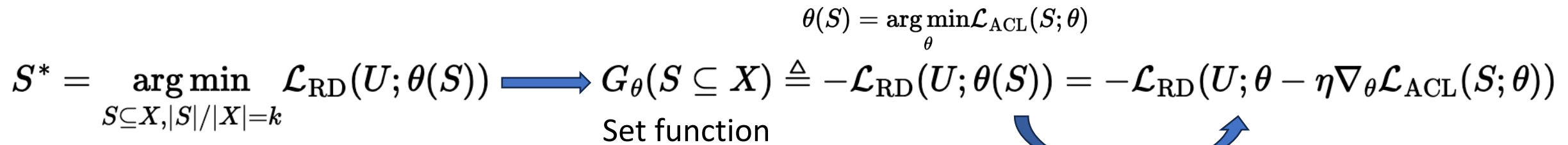
RD, without using labels, measures the adversarial robustness.



# RCS is a problem of set function maximization

$$\theta(S) = \arg \min_{\theta} \mathcal{L}_{\text{ACL}}(S; \theta)$$
$$S^* = \arg \min_{S \subseteq X, |S|/|X|=k} \mathcal{L}_{\text{RD}}(U; \theta(S)) \longrightarrow G_{\theta}(S \subseteq X) \triangleq -\mathcal{L}_{\text{RD}}(U; \theta(S)) = -\mathcal{L}_{\text{RD}}(U; \theta - \eta \nabla_{\theta} \mathcal{L}_{\text{ACL}}(S; \theta))$$

Set function



One-step gradient approximation

$$S^* = \arg \max_{S \subseteq X, |S|/|X|=k} G_{\theta}(S)$$

Transform into a problem of **maximizing a set function**  
**subject to a constraint on the size of the set**

# RCS via greedy search

RCS greedily finds and adds the data which has the largest marginal gain into the coreset.



$$G_{\theta}(x | S) = G_{\theta}(S \cup \{x\}) - G_{\theta}(S)$$

Derived via Taylor expansion

$$\begin{aligned}
 &= -\mathcal{L}_{RD}(U; \theta - \eta \nabla_{\theta} \mathcal{L}_{ACL}(S; \theta) - \eta \nabla_{\theta} \mathcal{L}_{ACL}(\{x\}; \theta)) \\
 &\quad + \mathcal{L}_{RD}(U; \theta - \eta \nabla_{\theta} \mathcal{L}_{ACL}(S; \theta)) \\
 &\approx -(\mathcal{L}_{RD}(U; \theta - \eta \nabla_{\theta} \mathcal{L}_{ACL}(S; \theta)) - \eta \nabla_{\theta} \mathcal{L}_{RD}(U; \theta_S)^{\top} \nabla_{\theta} \mathcal{L}_{ACL}(\{x\}; \theta) + \xi) \\
 &\quad + \mathcal{L}_{RD}(U; \theta - \eta \nabla_{\theta} \mathcal{L}_{ACL}(S; \theta)) \\
 &\approx \eta \nabla_{\theta} \mathcal{L}_{RD}(U; \theta - \eta \nabla_{\theta} \mathcal{L}_{ACL}(S; \theta))^{\top} \nabla_{\theta} \mathcal{L}_{ACL}(\{x\}; \theta)
 \end{aligned}$$

Marginal gain = Similarity (validation loss gradient, training loss gradient)

More adversarially robust



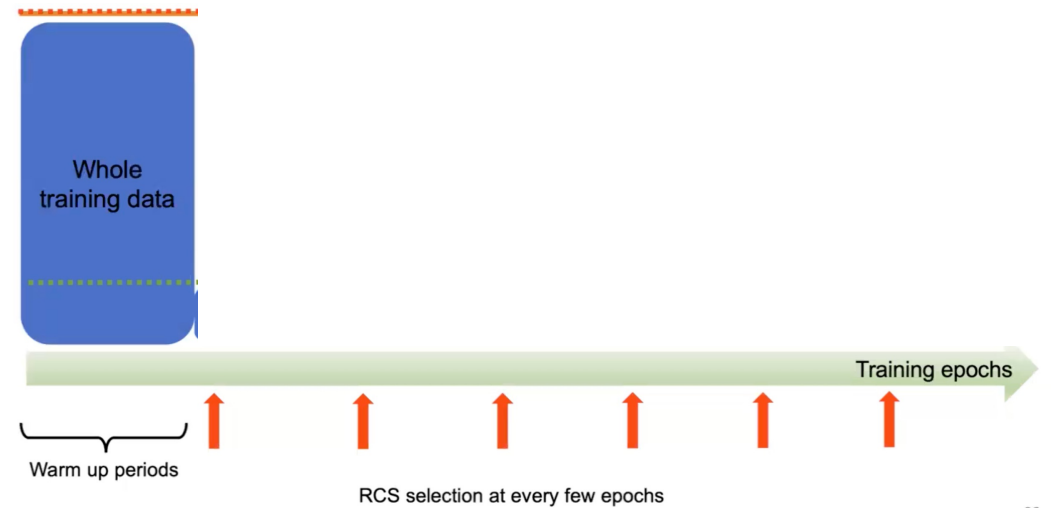
Most beneficial in  
optimizing RD



Training on the  
selected data

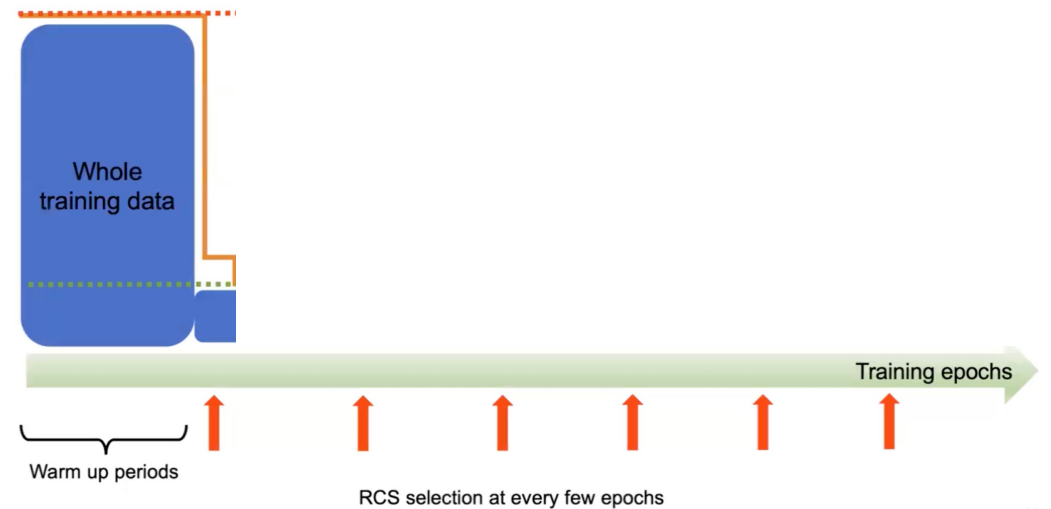
# Algorithm: Efficient ACL via RCS

- Step 1 (Warm up): Warm up training on entire training set to find a better starting point  $f_{\theta}$ .



# Algorithm: Efficient ACL via RCS

- Step 1 (Warm up): Warm up training on entire training set to find a better starting point  $f_\theta$ .
- **Step 2.1 (RCS):**  $S \leftarrow \emptyset$ .  $\theta' \leftarrow \theta$ . Compute gradients  $Q \leftarrow \{q_k = \nabla_{\theta} \mathcal{L}_{\text{ACL}}(x_k; \theta) \mid \forall x_k \in X\}$  on unlabeled training dataset  $X$ .
- **Step 2.2 (RCS):** Compute gradients  $q_U \leftarrow \nabla_{\theta} \mathcal{L}_{\text{RD}}(U; \theta')$  on unlabeled validation dataset  $U$ .
- **Step 2.3 (RCS):** Select a data  $x_k$ , whose gradient  $q_k$  matches best with  $q_U$ , i.e.,  $\arg \max_k \{q_k^\top q_U\}$ .  $G_\theta(x \mid S)$
- **Step 2.4 (RCS):**  $S \leftarrow S \cup \{x_k\}$ ,  $X \leftarrow X \setminus \{x_k\}$ ,  $\theta' \leftarrow \theta' - \eta' q_k$ .
- **Step 2.5 (RCS):** Repeat Step 2.2-2.3 until  $|S| / |X| = k$ .



# Algorithm: Efficient ACL via RCS

- Step 1 (Warm up): Warm up training on entire training set to find a better starting point  $f_\theta$ .
- **Step 2.1 (RCS):**  $S \leftarrow \emptyset$ .  $\theta' \leftarrow \theta$ . Compute gradients  $Q \leftarrow \{q_k = \nabla_{\theta} \mathcal{L}_{\text{ACL}}(x_k; \theta) \mid \forall x_k \in X\}$  on unlabeled training dataset  $X$ .
- **Step 2.2 (RCS):** Compute gradients  $q_U \leftarrow \nabla_{\theta} \mathcal{L}_{\text{RD}}(U; \theta')$  on unlabeled validation dataset  $U$ .
- **Step 2.3 (RCS):** Select a data  $x_k$ , whose gradient  $q_k$  matches best with  $q_U$ , i.e.,  $\arg \max_k \{q_k^\top q_U\}$ .  $G_\theta(x \mid S)$
- **Step 2.4 (RCS):**  $S \leftarrow S \cup \{x_k\}$ ,  $X \leftarrow X \setminus \{x_k\}$ ,  $\theta' \leftarrow \theta' - \eta' q_k$ .
- **Step 2.5 (RCS):** Repeat Step 2.2-2.3 until  $|S| / |X| = k$ .
- Step 3 (ACL training): Update parameters  $\theta \leftarrow \theta - \eta \nabla_{\theta} \mathcal{L}_{\text{ACL}}(S; \theta)$
- Step 4: Every  $I$  epochs, go to Step 2.1 to generate a new coreset; otherwise go to Step 3 to update model parameters. The algorithm stops when reaches the final training epoch.



# Algorithm: Efficient ACL via RCS

- Step 1 (Warm up): Warm up training on entire training set to find a better starting point  $f_\theta$ .
- **Step 2.1 (RCS):**  $S \leftarrow \emptyset$ .  $\theta' \leftarrow \theta$ . Compute gradients  $Q \leftarrow \{q_k = \nabla_{\theta} \mathcal{L}_{\text{ACL}}(x_k; \theta) \mid \forall x_k \in X\}$  on unlabeled training dataset  $X$ .
- **Step 2.2 (RCS):** Compute gradients  $q_U \leftarrow \nabla_{\theta} \mathcal{L}_{\text{RD}}(U; \theta')$  on unlabeled validation dataset  $U$ .
- **Step 2.3 (RCS):** Select a data  $x_k$ , whose gradient  $q_k$  matches best with  $q_U$ , i.e.,  $\arg \max_k \{q_k^\top q_U\}$ .  $G_\theta(x \mid S)$
- **Step 2.4 (RCS):**  $S \leftarrow S \cup \{x_k\}$ ,  $X \leftarrow X \setminus \{x_k\}$ ,  $\theta' \leftarrow \theta' - \eta' q_k$ .
- **Step 2.5 (RCS):** Repeat Step 2.2-2.3 until  $|S| / |X| = k$ .
- Step 3 (ACL training): Update parameters  $\theta \leftarrow \theta - \eta \nabla_{\theta} \mathcal{L}_{\text{ACL}}(S; \theta)$
- Step 4: Every  $I$  epochs, go to Step 2.1 to generate a new coreset; otherwise go to Step 3 to update model parameters. The algorithm stops when reaches the final training epoch.

Whether the greedy search provide any optimality guarantee theoretically?



# Theoretical analysis: greedy search has an optimality guarantee

Step 2.3 (RCS): Select a data  $x_k$ , whose gradient  $q_k$  matches best with  $q_U$ , i.e.,  $\arg \max_k \{q_k^\top q_U\}$ .  $G_\theta(x | S)$

- Proof sketch using a proxy set problem  $\hat{S}^* = \arg \max_{S \subseteq X, |S|/|X|=k} \hat{G}_\theta(S) = \arg \max_{S \subseteq X, |S|/|X|=k} G_\theta(S) + |S|\sigma$

## 1. Monotonicity

$$\hat{G}(x | X) = \hat{G}(S \cup \{x\}) - \hat{G}(S) \geq 0 \text{ for any } S \subseteq X \text{ and } x \in X \setminus S.$$

More data, better representation in terms of robustness

## 2. $\gamma^*$ -submodularity---diminishing returns

$$\forall_{A, B | A \subseteq B} G_\theta(x | A) \geq (1 - \gamma^*) G_\theta(x | B)$$

$$\text{where } \gamma^* = \frac{1}{2\sigma-1} \in (0, 1) \text{ and } A \subseteq B \subseteq X$$

More data have diminishing gains for learning representations

Proxy  
coreset

$$\hat{S}^*$$



## Guaranteed lower bound of the original set problem using the proxy coreset

$$G_\theta(\hat{S}^*) \geq G_\theta^* - (G_\theta^* + kN\sigma) \cdot e^{-\gamma^*}$$

Proxy coreset selected via greedy search based on the proxy marginal gain  $\hat{G}_\theta(x | S)$

Guaranteed Lower bound



# Theoretical analysis: greedy search has an optimality guarantee

- Proof sketch using a proxy set problem  $\hat{S}^* = \arg \max_{S \subseteq X, |S|/|X|=k} \hat{G}_\theta(S) = \arg \max_{S \subseteq X, |S|/|X|=k} G_\theta(S) + |S|\sigma$

**Step 2.3 (RCS):** Select a data  $x_k$ , whose gradient  $q_k$  matches best with  $q_U$ , i.e.,

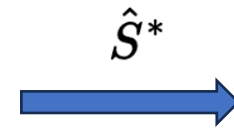
$$\arg \max_k \{q_k^\top q_U\}. \quad G_\theta(x | S)$$

Our greedy search via  
the original marginal gain

Approximate



the proxy marginal gain



$$G_\theta(\hat{S}^*) \geq G_\theta^* - (G_\theta^* + kN\sigma) \cdot e^{-\gamma^*}$$

The proxy coreset provides  
a guaranteed lower bound!

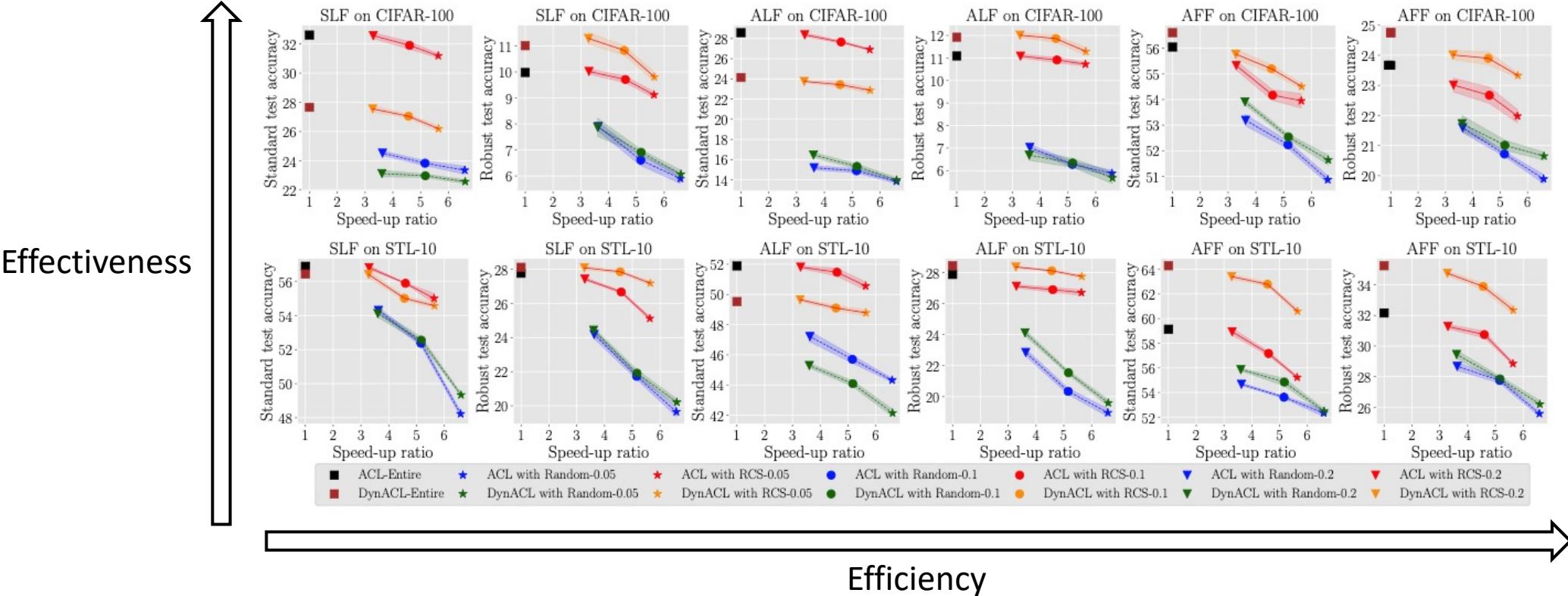
$$\begin{aligned} \hat{G}_\theta(x | S) &= \hat{G}_\theta(S \cup \{x\}) - \hat{G}_\theta(S) \\ &\approx \eta \nabla_\theta \mathcal{L}_{\mathcal{RD}}(U; \theta - \eta \nabla_\theta \mathcal{L}_{\text{ACL}}(S; \theta))^\top \nabla_\theta \mathcal{L}_{\text{ACL}}(\{x\}; \theta) + \sigma \\ &= \boxed{G_\theta(x | S)} + \sigma. \end{aligned}$$

Can be discarded since it is a constant

# Empirical Results

RCS is more efficient (higher speed-up ratio) compared to ACL on the entire set.

RCS is more effective (higher test accuracy) compared to random selection.



The upper-right (ours) is better!

# Empirical Results

For the first time to conduct ACL on ImageNet-1K using WRN-28-10

Table 1: Cross-task adversarial robustness transferability from ImageNet-1K to CIFAR-10.

Pre-training	Runing time (hours)	SLF		ALF		AFF	
		SA (%)	RA (%)	SA (%)	RA (%)	SA (%)	RA (%)
Standard CL	147.4	<b>84.36</b> $\pm 0.17$	0.01 $\pm 0.01$	10.00 $\pm 0.00$	10.00 $\pm 0.00$	<b>86.63</b> $\pm 0.12$	49.71 $\pm 0.16$
ACL on entire set	650.2	-	-	-	-	-	-
ACL with Random	94.3	68.75 $\pm 0.06$	15.89 $\pm 0.06$	59.57 $\pm 0.28$	27.14 $\pm 0.19$	84.75 $\pm 0.18$	50.12 $\pm 0.21$
ACL with RCS	111.8	70.02 $\pm 0.12$	<b>22.45</b> $\pm 0.13$	<b>63.94</b> $\pm 0.21$	<b>31.13</b> $\pm 0.17$	85.23 $\pm 0.23$	<b>52.21</b> $\pm 0.14$

Table 2: Cross-task adversarial robustness transferability from ImageNet-1K to CIFAR-100.

Pre-training	Runing time (hours)	SLF		ALF		AFF	
		SA (%)	RA (%)	SA (%)	RA (%)	SA (%)	RA (%)
Standard CL	147.4	<b>57.34</b> $\pm 0.23$	0.01 $\pm 0.01$	9.32 $\pm 0.01$	0.06 $\pm 0.01$	<b>61.33</b> $\pm 0.12$	25.11 $\pm 0.15$
ACL on entire set	650.2	-	-	-	-	-	-
ACL with Random	94.3	38.53 $\pm 0.15$	10.50 $\pm 0.13$	28.44 $\pm 0.23$	11.93 $\pm 0.21$	59.63 $\pm 0.33$	25.46 $\pm 0.26$
ACL with RCS	111.8	40.28 $\pm 0.17$	<b>14.55</b> $\pm 0.10$	<b>33.15</b> $\pm 0.26$	<b>14.89</b> $\pm 0.16$	60.25 $\pm 0.18$	<b>28.24</b> $\pm 0.13$

We prove the possibility of applying ACL on large-scale datasets.

# Empirical Results

RCS for speeding up supervised adversarial training (SAT) on ImageNet-1K while maintaining robustness transferability.

Table 17: Cross-task adversarial robustness transferability of adversarially pre-trained WRN-28-10 from ImageNet-1K to CIFAR-10. Here, “RA” stands for robust test accuracy under PGD-20 attacks following the setting of Hendrycks et al. [54]. The number after the dash line denotes subset fraction  $k \in \{0.05, 0.1, 0.2\}$ .

Pre-training	Runing time (hours)	ALF		AFF	
		SA (%)	RA (%)	SA (%)	RA (%)
Standard training on entire set	66.7	10.12	10.04	84.73	51.91
SAT [54] on entire set	341.7	85.90	50.89	89.35	59.68
SAT with Random-0.05	53.6	69.59	31.58	85.55	53.53
SAT with RCS-0.05	<b>68.6</b>	<b>79.72</b>	<b>44.44</b>	<b>87.99</b>	<b>56.87</b>
SAT with Random-0.1	70.2	73.28	33.86	86.78	54.95
SAT with RCS-0.1	<b>81.9</b>	<b>81.92</b>	<b>45.10</b>	<b>88.87</b>	<b>57.69</b>
SAT with Random-0.2	103.4	75.46	39.62	86.64	56.46
SAT with RCS-0.2	<b>121.9</b>	<b>83.94</b>	<b>46.88</b>	<b>89.54</b>	<b>58.13</b>

**RCS is also applicable to robust supervised pre-training!**

# Towards Robust Foundation Models

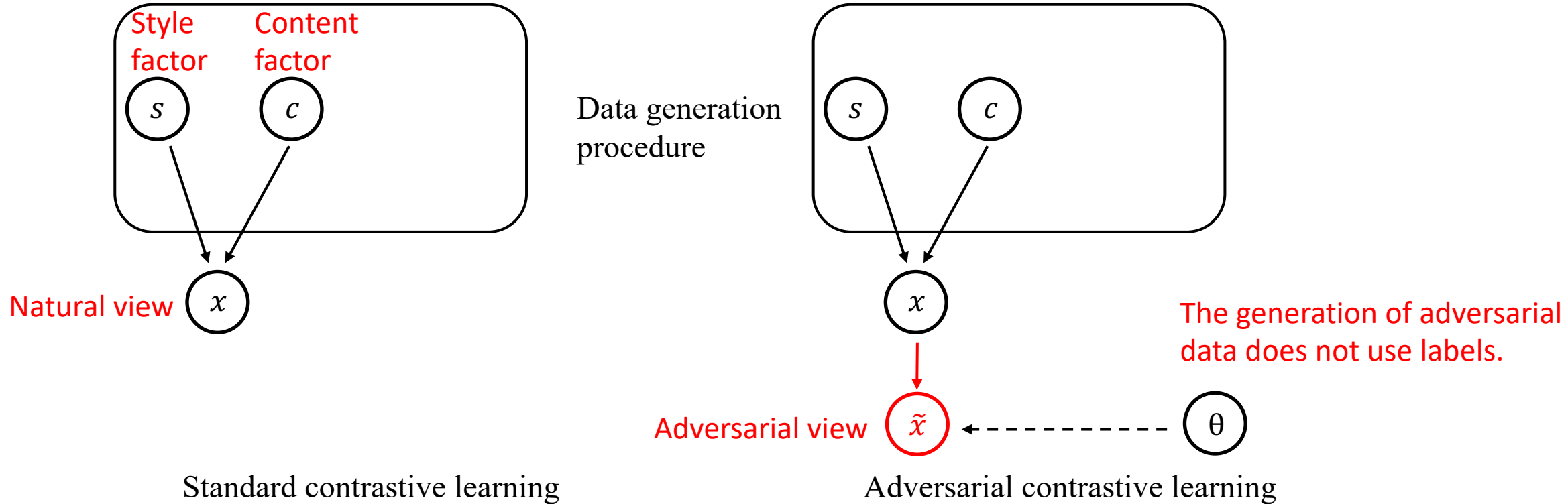
Efficient and **Effective** Adversarial Contrastive Learning

# Motivation

- *Limited robustness transferability* to downstream tasks

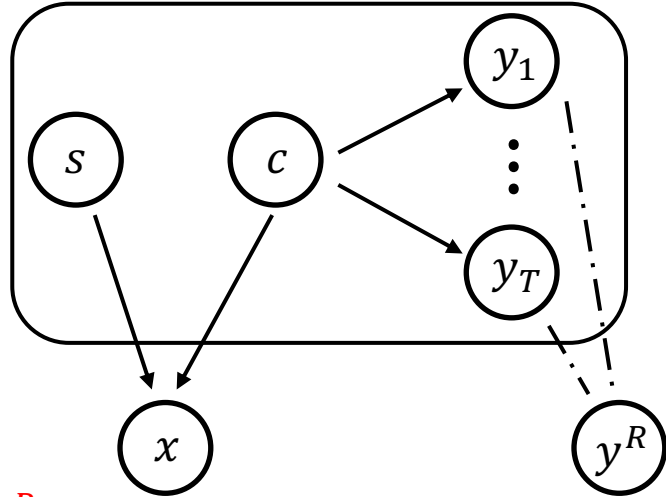
How can we improve ACL's robustness transferability?

# Causal View of ACL

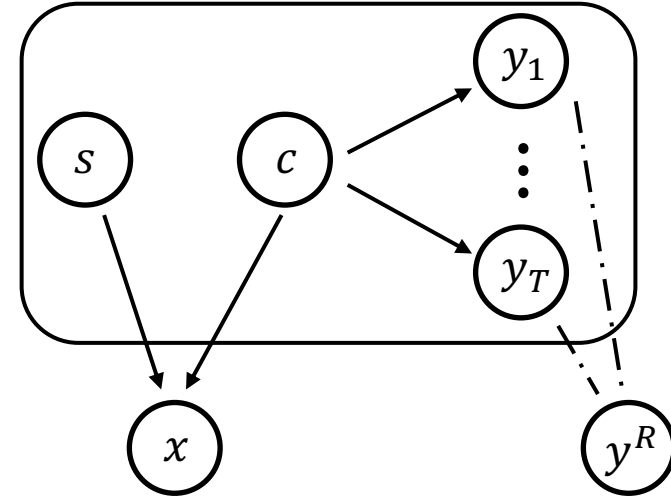


# Causal View of ACL

Target labels from an unknown downstream task



Data generation procedure



Learning procedure



Adversarial contrastive learning

The proxy label  $y^R$  is the refinement of  $y_t$ .

Standard contrastive learning

Proxy label  $y^R$

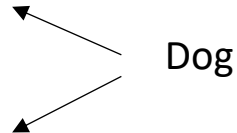
Target label  $y_t$



Golden Retriever with yellow hair

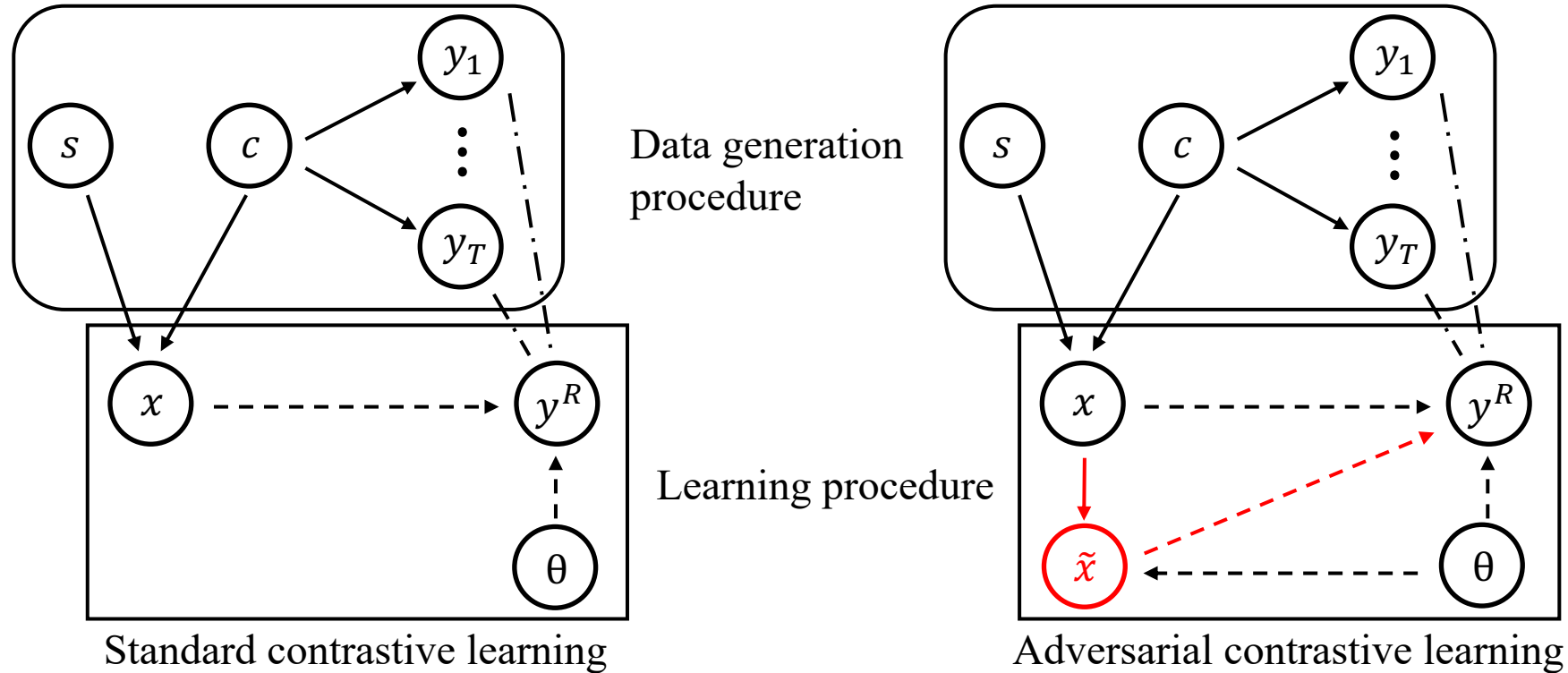


Labrador Retriever with black hair

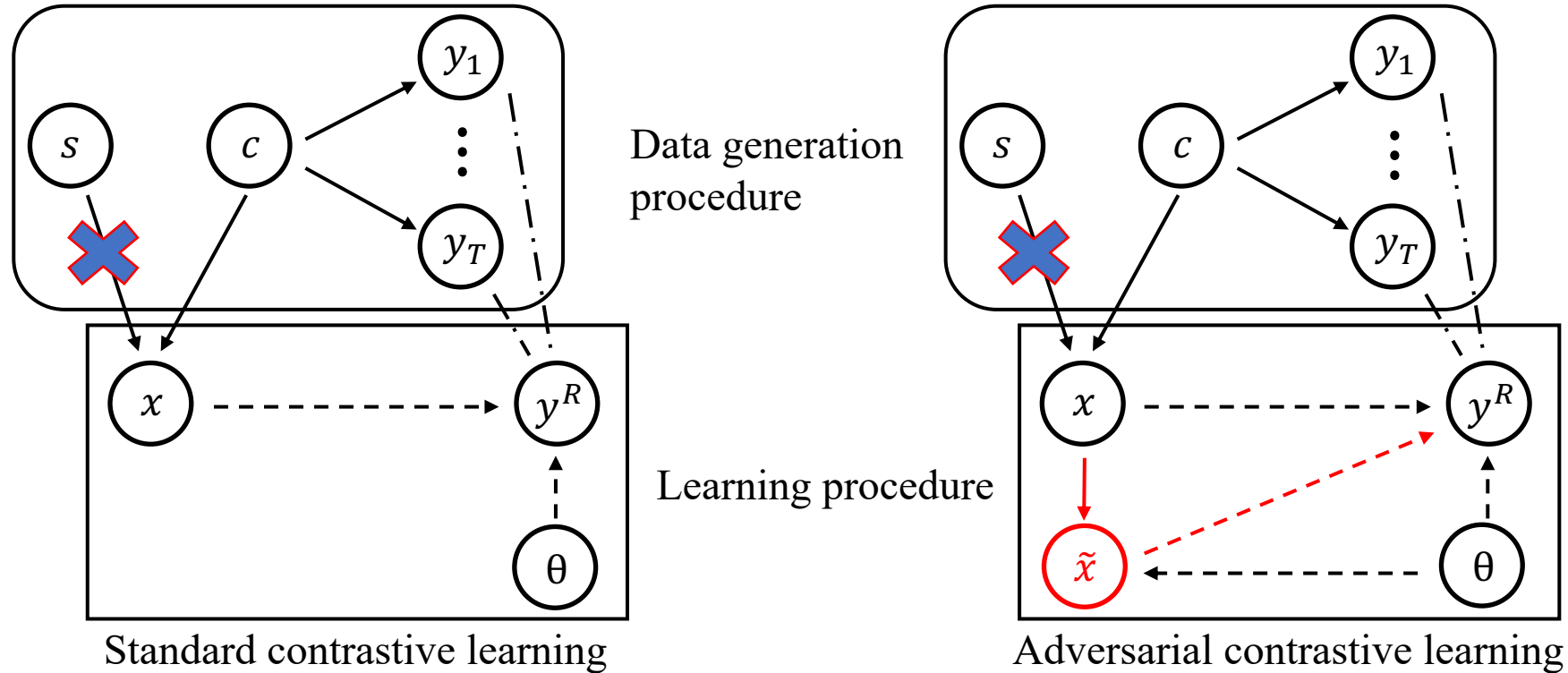




# Causal View of ACL



# Causal View of ACL



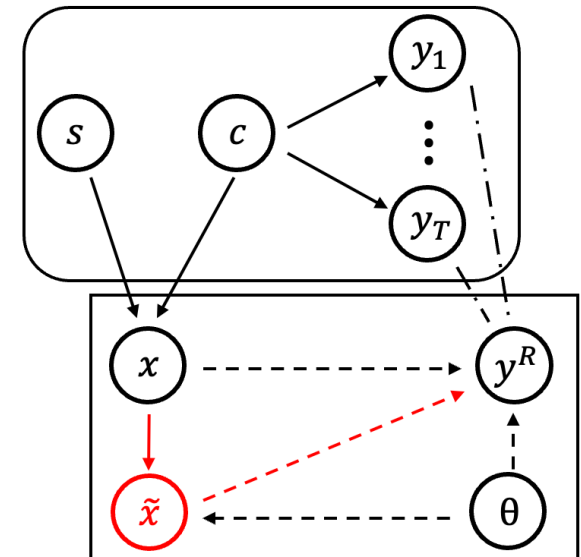
Style-invariant criterion: The intervention on the style factor should not affect the conditional probability

$$p^{do(\tau_i)}(y^R | x) = p^{do(\tau_j)}(y^R | x)$$

# Adversarial Invariant Regularization (AIR)

- The conditional probability learned via ACL the mild assumption of Markov condition

$$p(y^R|x) = p(y^R|\tilde{x})p(\tilde{x}|x)$$

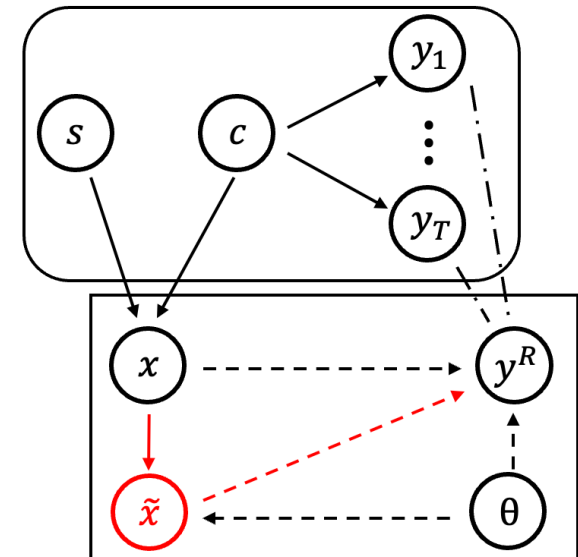


Adversarial contrastive learning

# Adversarial Invariant Regularization (AIR)

- The conditional probability learned via ACL  $p(y^R|x) = p(y^R|\tilde{x})p(\tilde{x}|x)$
- Style-independent criterion: The intervention on the style factor should not affect the conditional probability

$$p^{do(\tau_i)}(y^R | x) = p^{do(\tau_j)}(y^R | x)$$

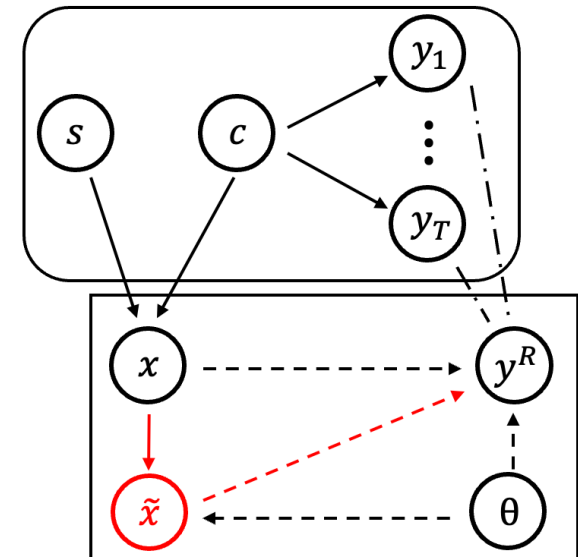


Adversarial contrastive learning

# Adversarial Invariant Regularization (AIR)

- The conditional probability learned via ACL  $p(y^R|x) = p(y^R|\tilde{x})p(\tilde{x}|x)$
- Style-independent criterion: The intervention on the style factor should not affect the conditional probability

$$p^{do(\tau_i)}(y^R|\tilde{x})p^{do(\tau_i)}(\tilde{x}|x) = p^{do(\tau_j)}(y^R|\tilde{x})p^{do(\tau_j)}(\tilde{x}|x) \quad \forall \tau_i, \tau_j \in \mathcal{T}$$



Adversarial contrastive learning

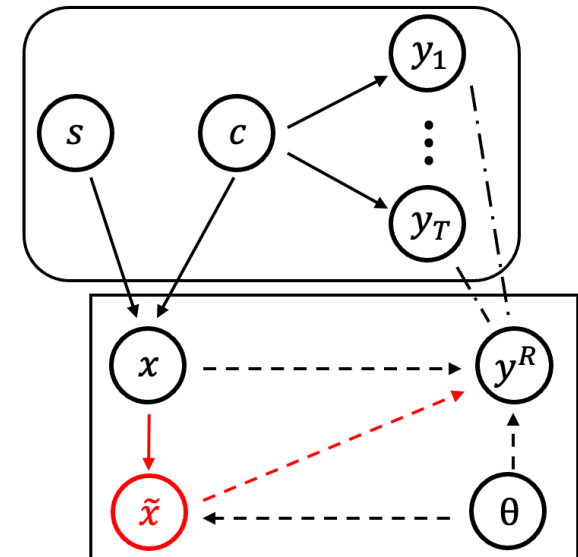
# Adversarial Invariant Regularization (AIR)

- The conditional probability learned via ACL  $p(y^R|x) = p(y^R|\tilde{x})p(\tilde{x}|x)$
- Style-independent criterion: The intervention on the style factor should not affect the conditional probability

$$p^{do(\tau_i)}(y^R|\tilde{x})p^{do(\tau_i)}(\tilde{x}|x) = p^{do(\tau_j)}(y^R|\tilde{x})p^{do(\tau_j)}(\tilde{x}|x) \quad \forall \tau_i, \tau_j \in \mathcal{T}$$

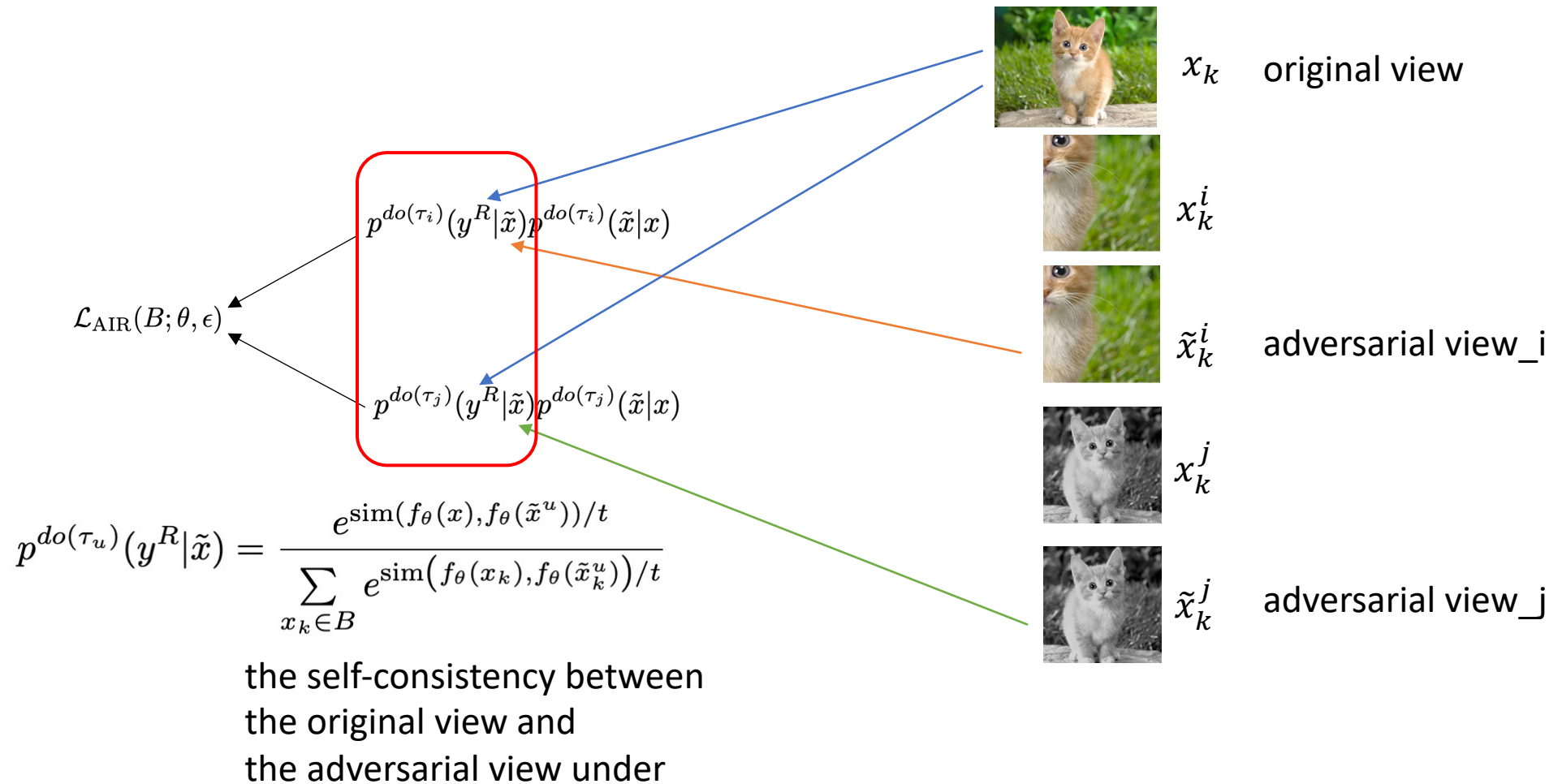
- Loss function of AIR -> to enforce style-independence

$$\mathcal{L}_{\text{AIR}}(B; \theta, \epsilon) = \text{KL} \left( p^{do(\tau_i)}(y^R|\tilde{x})p^{do(\tau_i)}(\tilde{x}|x) \parallel p^{do(\tau_j)}(y^R|\tilde{x})p^{do(\tau_j)}(\tilde{x}|x); B \right)$$

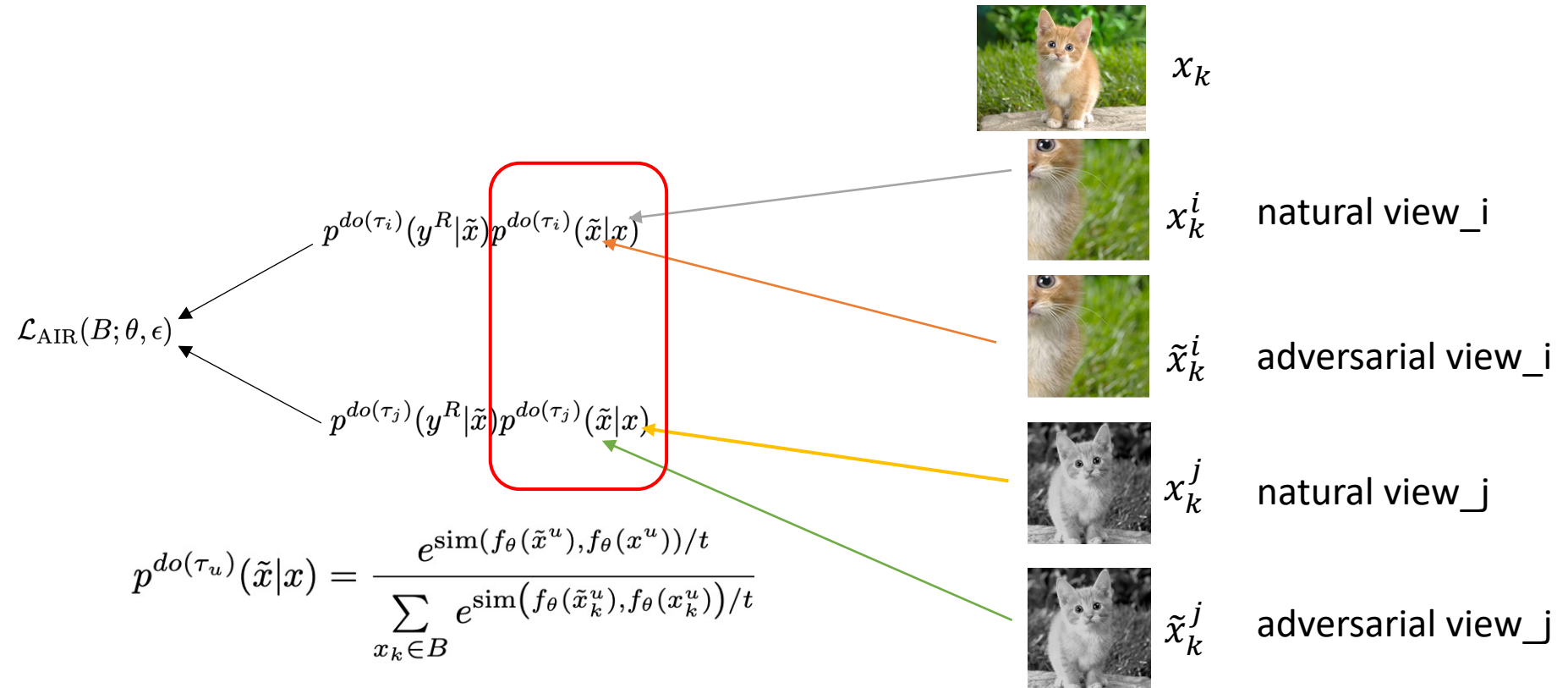


Adversarial contrastive learning

# Understanding of AIR



# Understanding of AIR



the self-consistency between the adversarial view and the natural view

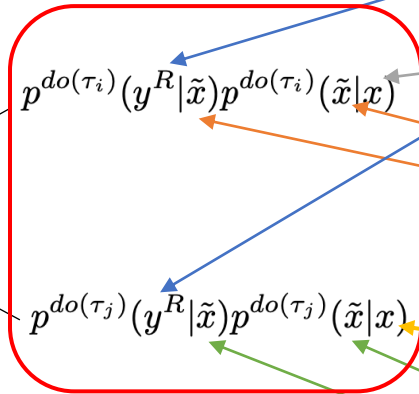


# Understanding of AIR

$$\mathcal{L}_{\text{AIR}}(B; \theta, \epsilon) = \text{KL} \left( p^{\text{do}(\tau_i)}(y^R | \tilde{x}) p^{\text{do}(\tau_i)}(\tilde{x} | x) \parallel p^{\text{do}(\tau_j)}(y^R | \tilde{x}) p^{\text{do}(\tau_j)}(\tilde{x} | x); B \right)$$

AIR enforces self-consistency to be cross-consistent under different augmentations

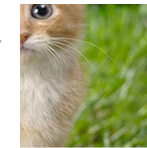
$$\mathcal{L}_{\text{AIR}}(B; \theta, \epsilon)$$



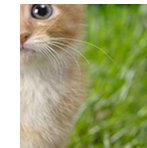
the self-consistency among the different views



$x_k$  original view



$x_k^i$  natural view\_i



$\tilde{x}_k^i$  adversarial view\_i



$x_k^j$  natural view\_j



$\tilde{x}_k^j$  adversarial view\_j

# Algorithm: Enhancing ACL via AIR

---

**Algorithm 1** ACL with Adversarial Invariant Regularization (AIR)

---

- 1: **Input:** Unlabeled training set  $U$ , total training epochs  $E$ , learning rate  $\eta$ , batch size  $\beta$ , adversarial budget  $\epsilon > 0$ , hyperparameters  $\lambda_1$  and  $\lambda_2$
  - 2: **Output:** Pre-trained representation extractor  $h_\theta$
  - 3: Initialize parameters of model  $f_\theta = g \circ h_\theta$
  - 4: **for**  $e = 0$  **to**  $E - 1$  **do**
  - 5:   **for** batch  $m = 1, \dots, \lceil |U|/\beta \rceil$  **do**  $\epsilon > 0$ : Regulate robust representations
  - 6:     Sample a minibatch  $B_m$  from  $U$
  - 7:     Update  $\theta \leftarrow \theta - \eta \cdot \nabla_\theta \sum_{x_k \in B_m} \ell_{\text{ACL}}(x_k; \theta) + \lambda_1 \cdot \mathcal{L}_{\text{AIR}}(B_m; \theta, 0) + \lambda_2 \cdot \mathcal{L}_{\text{AIR}}(B_m; \theta, \epsilon)$   $\epsilon = 0$ : Regulate standard representations
  - 8:   **end for**
  - 9: **end for**
-

# AIR achieves SOTA robustness transferability

Table 3: Cross-task adversarial robustness transferability.  $\mathcal{D}_1 \rightarrow \mathcal{D}_2$  denotes pre-training and finetuning are conducted on the dataset  $\mathcal{D}_1$  and  $\mathcal{D}_2 (\neq \mathcal{D}_1)$ , respectively.

$\mathcal{D}_1 \rightarrow \mathcal{D}_2$	Pre-training	SLF		ALF		AFF	
		AA (%)	SA (%)	AA (%)	SA (%)	AA (%)	SA (%)
CIFAR-10 $\rightarrow$ CIFAR-100	ACL [26]	9.98 $\pm$ 0.02	32.61 $\pm$ 0.04	11.09 $\pm$ 0.06	28.58 $\pm$ 0.06	22.67 $\pm$ 0.16	56.05 $\pm$ 0.19
	ACL-AIR	<b>11.04</b> $\pm$ 0.06	<b>39.45</b> $\pm$ 0.07	<b>13.30</b> $\pm$ 0.02	<b>36.10</b> $\pm$ 0.05	<b>23.45</b> $\pm$ 0.04	<b>56.31</b> $\pm$ 0.06
	DynACL [33]	11.01 $\pm$ 0.02	27.66 $\pm$ 0.03	11.92 $\pm$ 0.05	24.14 $\pm$ 0.09	24.17 $\pm$ 0.10	55.61 $\pm$ 0.17
	DynACL-AIR	<b>12.20</b> $\pm$ 0.04	<b>31.33</b> $\pm$ 0.03	<b>12.70</b> $\pm$ 0.03	<b>28.70</b> $\pm$ 0.05	<b>24.82</b> $\pm$ 0.07	<b>57.00</b> $\pm$ 0.13
CIFAR-10 $\rightarrow$ STL-10	ACL [26]	25.41 $\pm$ 0.08	56.53 $\pm$ 0.10	27.17 $\pm$ 0.09	51.71 $\pm$ 0.17	32.66 $\pm$ 0.07	61.41 $\pm$ 0.13
	ACL-AIR	<b>28.00</b> $\pm$ 0.12	<b>61.91</b> $\pm$ 0.13	<b>30.06</b> $\pm$ 0.10	<b>62.03</b> $\pm$ 0.11	<b>34.26</b> $\pm$ 0.09	<b>62.58</b> $\pm$ 0.10
	DynACL [33]	28.52 $\pm$ 0.09	52.45 $\pm$ 0.10	29.13 $\pm$ 0.13	49.53 $\pm$ 0.17	35.25 $\pm$ 0.15	63.29 $\pm$ 0.18
	DynACL-AIR	<b>29.88</b> $\pm$ 0.04	<b>54.59</b> $\pm$ 0.12	<b>31.24</b> $\pm$ 0.06	<b>57.14</b> $\pm$ 0.09	<b>35.66</b> $\pm$ 0.05	<b>63.74</b> $\pm$ 0.12
CIFAR-100 $\rightarrow$ CIFAR-10	ACL [26]	18.72 $\pm$ 0.07	60.90 $\pm$ 0.02	26.92 $\pm$ 0.11	57.35 $\pm$ 0.07	44.07 $\pm$ 0.11	75.19 $\pm$ 0.10
	ACL-AIR	<b>19.90</b> $\pm$ 0.04	<b>64.89</b> $\pm$ 0.09	<b>27.65</b> $\pm$ 0.06	<b>60.79</b> $\pm$ 0.04	<b>44.84</b> $\pm$ 0.14	<b>75.67</b> $\pm$ 0.13
	DynACL [33]	25.23 $\pm$ 0.12	59.12 $\pm$ 0.10	28.92 $\pm$ 0.10	56.09 $\pm$ 0.14	47.40 $\pm$ 0.23	77.92 $\pm$ 0.18
	DynACL-AIR	<b>25.63</b> $\pm$ 0.07	<b>59.83</b> $\pm$ 0.08	<b>29.32</b> $\pm$ 0.06	<b>56.65</b> $\pm$ 0.06	<b>47.92</b> $\pm$ 0.12	<b>78.44</b> $\pm$ 0.10
CIFAR-100 $\rightarrow$ STL-10	ACL [26]	21.77 $\pm$ 0.07	46.19 $\pm$ 0.05	24.46 $\pm$ 0.09	45.40 $\pm$ 0.12	28.76 $\pm$ 0.07	56.16 $\pm$ 0.13
	ACL-AIR	<b>22.44</b> $\pm$ 0.04	<b>51.52</b> $\pm$ 0.02	<b>26.55</b> $\pm$ 0.06	<b>53.24</b> $\pm$ 0.09	<b>30.40</b> $\pm$ 0.08	<b>58.45</b> $\pm$ 0.11
	DynACL [33]	23.17 $\pm$ 0.09	47.54 $\pm$ 0.14	26.24 $\pm$ 0.13	45.70 $\pm$ 0.14	31.17 $\pm$ 0.14	58.35 $\pm$ 0.18
	DynACL-AIR	<b>23.24</b> $\pm$ 0.07	<b>48.20</b> $\pm$ 0.08	<b>26.60</b> $\pm$ 0.05	<b>48.55</b> $\pm$ 0.12	<b>31.42</b> $\pm$ 0.07	<b>58.59</b> $\pm$ 0.10

# AIR achieves SOTA robustness transferability

- Robustness transferability via automated robust fine-tuning
  - AutoLoRa: an automated and parameter-free robust fine-tuning framework

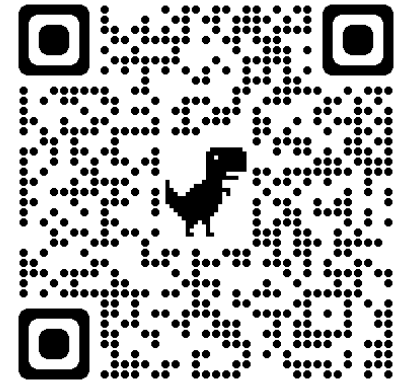
Table 7: Cross-task adversarial robustness transferability evaluated via AutoLoRa [44].  $\mathcal{D}_1 \rightarrow \mathcal{D}_2$  denotes pre-training and finetuning are conducted on the dataset  $\mathcal{D}_1$  and  $\mathcal{D}_2 (\neq \mathcal{D}_1)$ , respectively. “Diff” refers to the gap between the performance achieved by AutoLoRa and that achieved by vanilla finetuning (reported in Table 3).

$\mathcal{D}_1 \rightarrow \mathcal{D}_2$	Finetuning mode	Pre-training	AutoLoRa [44]		Diff	
			AA (%)	SA (%)	AA (%)	SA (%)
CIFAR-10 $\rightarrow$ STL-10	SLF	DynACL [33]	30.18	54.23	+1.01	+1.82
		DynACL-AIR	<b>30.48</b>	<b>56.56</b>	<b>+0.84</b>	<b>+0.72</b>
	ALF	DynACL [33]	31.72	57.30	+2.13	+7.75
		DynACL-AIR	<b>31.81</b>	<b>57.40</b>	<b>+0.57</b>	<b>+0.26</b>
	AFF	DynACL [33]	35.51	64.16	+0.26	+0.63
		DynACL-AIR	<b>35.88</b>	<b>64.25</b>	<b>+0.22</b>	<b>+0.51</b>
CIFAR-100 $\rightarrow$ STL-10	SLF	DynACL [33]	23.27	48.93	+0.10	+1.39
		DynACL-AIR	<b>23.44</b>	<b>50.28</b>	<b>+0.20</b>	<b>+2.08</b>
	ALF	DynACL [33]	26.53	48.56	+0.29	+2.86
		DynACL-AIR	<b>26.89</b>	<b>49.02</b>	<b>+0.29</b>	<b>+0.47</b>
	AFF	DynACL [33]	31.25	58.56	+0.08	+0.06
		DynACL-AIR	<b>31.57</b>	<b>58.65</b>	<b>+0.15</b>	<b>+0.21</b>

# AIR ranks First in RobustSSL Benchmark

Standard Linear Fine-Tuning (SLF)			Vanilla Fine-Tuning		
Rank	Paper	Venue	Robust Accuracy	Corruption Accuracy	Standard Accuracy
1	Enhancing Adversarial Contrastive Learning via Adversarial Invariant Regularization <small>*Using post-processing</small>	NeurIPS 2023	<b>46.99</b>	72.11	81.80
2	Rethinking the Effect of Data Augmentation in Adversarial Contrastive Learning <small>*Using post-processing</small>	ICLR 2023	<b>46.54</b>	71.96	79.82
3	Enhancing Adversarial Contrastive Learning via Adversarial Invariant Regularization	NeurIPS 2023	<b>45.17</b>	70.51	78.08
4	Rethinking the Effect of Data Augmentation in Adversarial Contrastive Learning	ICLR 2023	<b>45.09</b>	68.67	75.41
5	Efficient Adversarial Contrastive Learning via Robustness-Aware Coreset Selection	NeurIPS 2023	<b>44.29</b>	69.56	77.14
6	Decoupled Adversarial Contrastive Learning for Self-supervised Adversarial Robustness	ECCV 2022	<b>43.27</b>	73.06	79.94
7	When Does Contrastive Learning Preserve Adversarial Robustness from Pretraining to Finetuning? <small>*Using ImageNet-1K pre-trained models</small>	NeurIPS 2021	<b>43.18</b>	73.14	82.36
8	Adversarial Contrastive Learning via Asymmetric InfoNCE <small>*Using ImageNet-1K pre-trained models</small>	ECCV 2022	<b>42.72</b>	74.09	83.70
9	Robust Pre-Training by Adversarial Contrastive Learning	NeurIPS 2020	<b>39.17</b>	70.72	78.22
10	Adversarial Self-Supervised Contrastive Learning	NeurIPS 2020	<b>26.12</b>	-	77.90

Adversarial Linear Fine-Tuning (ALF)			Vanilla Fine-Tuning		
Rank	Paper	Venue	Robust Accuracy	Corruption Accuracy	Standard Accuracy
1	Enhancing Adversarial Contrastive Learning via Adversarial Invariant Regularization <small>*Using post-processing</small>	NeurIPS 2023	<b>48.23</b>	71.74	79.56
2	Rethinking the Effect of Data Augmentation in Adversarial Contrastive Learning <small>*Using post-processing</small>	ICLR 2023	<b>47.98</b>	70.89	78.81
3	Enhancing Adversarial Contrastive Learning via Adversarial Invariant Regularization	NeurIPS 2023	<b>46.14</b>	69.97	77.42
4	Efficient Adversarial Contrastive Learning via Robustness-Aware Coreset Selection	NeurIPS 2023	<b>45.75</b>	67.84	74.95
5	Rethinking the Effect of Data Augmentation in Adversarial Contrastive Learning	ICLR 2023	<b>45.67</b>	66.69	72.97
6	When Does Contrastive Learning Preserve Adversarial Robustness from Pretraining to Finetuning? <small>*Using ImageNet-1K pre-trained models</small>	NeurIPS 2021	<b>44.05</b>	71.50	80.04
7	Adversarial Contrastive Learning via Asymmetric InfoNCE <small>*Using ImageNet-1K pre-trained models</small>	ECCV 2022	<b>43.28</b>	71.61	80.30
8	Decoupled Adversarial Contrastive Learning for Self-supervised Adversarial Robustness	ECCV 2022	<b>41.99</b>	71.66	77.71
9	Robust Pre-Training by Adversarial Contrastive Learning	NeurIPS 2020	<b>40.60</b>	68.56	75.53
10	Adversarial Self-Supervised Contrastive Learning	NeurIPS 2020	<b>29.69</b>	-	75.62



Robust Self-Supervised Learning  
(RobustSSL) Benchmark

<https://robustssl.github.io>

# Thank you for your attention!

- Summary
  - With RCS and AIR, we can efficiently build effective robust foundation models!
- Potential future directions
  - Explore the potential applications of ACL in various CV, NLP, and multi-modal tasks.
- References
  1. Xu, Xilie, Jingfeng Zhang, Feng Liu, Masashi Sugiyama, and Mohan Kankanhalli. "Efficient Adversarial Contrastive Learning via Robustness-Aware Coreset Selection." NeurIPS 2023 (spotlight).
  2. Xu, Xilie, Jingfeng Zhang, Feng Liu, Masashi Sugiyama, and Mohan Kankanhalli. "Enhancing Adversarial Contrastive Learning via Adversarial Invariant Regularization." NeurIPS 2023.
  3. Xu, Xilie, Jingfeng Zhang, and Mohan Kankanhalli. "Autolora: A parameter-free automated robust fine-tuning framework." arXiv preprint arXiv:2310.01818 (2023).
- Q&A