# Towards Effective and Efficient Self-Supervised Robust Pre-Training

Xilie Xu

PhD Student at School of Computing, NUS

$26^{th}$ Aug 2023

AIGC'23 Forum 4: "Foundation Model and Fine-Tuning Strategy"

Chaired by Dr. Jingfeng Zhang

# Outline

- Backgrounds
  - Adversarial attack and defense
  - Robust pre-training
- How to make self-supervised robust pre-training
  - More efficient
  - More effective
- Future directions

# Outline

- **Backgrounds**
  - **Adversarial attack and defense**
  - **Robust pre-training**
- How to make self-supervised robust pre-training
  - More efficient
  - More effective
- Future directions

# Gap between AI development and deployment

Develop AI-based applications
in an idealized environment

Deploy AI-base applications
in the wild



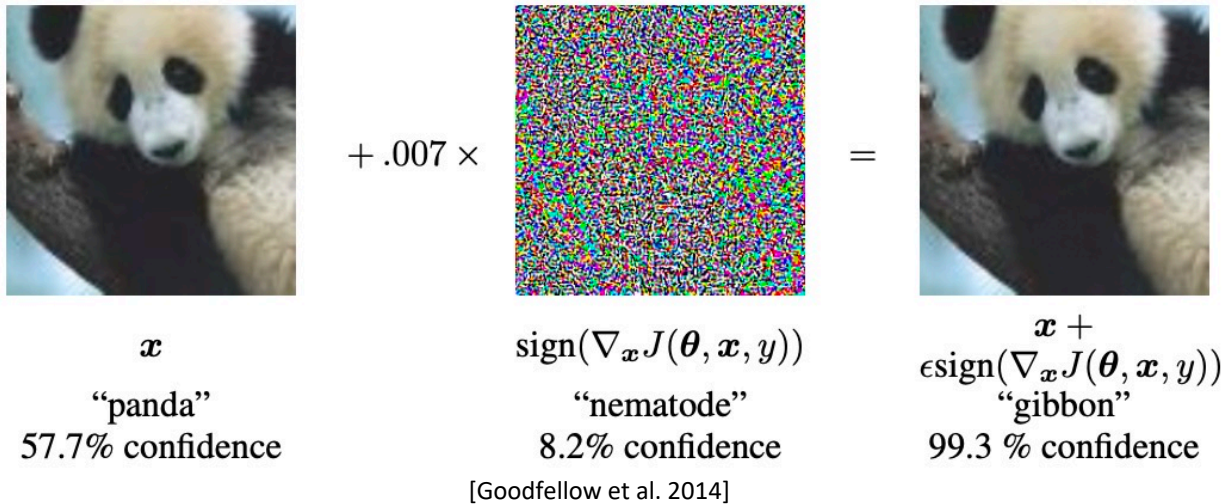Image from https://blog.si-log.net/transport-by-sea-by-land-or-by-air-the-differences-and-similarities

Image from https://www.primeins.com/insurance-news/how-to-protect-your-boat-from-a-tropical-storm-or-hurricane

Threats
("storm")

- Poisoning attack
- Backdoor attack
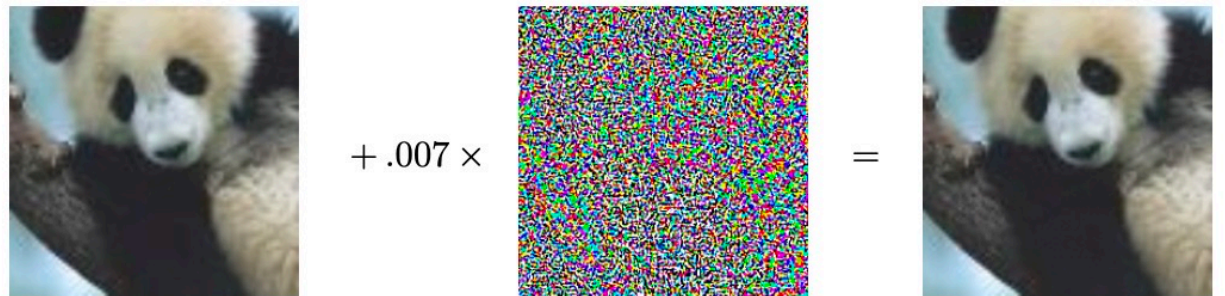- Distribution shift
- **Adversarial attack**

# Adversarial attacks

Objective: Make the model misclassify the adversarial data.



$+.007 \times$

$=$

$x$

"panda"
57.7% confidence

$\text{sign}(\nabla_x J(\boldsymbol{\theta}, \boldsymbol{x}, y))$

"nematode"
8.2% confidence

[Goodfellow et al. 2014]

$\boldsymbol{x} +$
$\epsilon \text{sign}(\nabla_x J(\boldsymbol{\theta}, \boldsymbol{x}, y))$
"gibbon"
99.3 % confidence

Natural data $x$    +    Imperceptible adversarial perturbation    =    Adversarial data $\tilde{x}$

# Adversarial attacks

Objective: Make the model misclassify the adversarial data.

$$\tilde{x} = argmax_{\tilde{x} \in \mathcal{B}_\epsilon[x]} \ell(f(\tilde{x}), y)$$



$+ .007 \times$

$$\text{sign}(\nabla_x J(\boldsymbol{\theta}, \boldsymbol{x}, y))$$

"nematode"
8.2% confidence

$=$

$$\boldsymbol{x} + \epsilon \text{sign}(\nabla_x J(\boldsymbol{\theta}, \boldsymbol{x}, y))$$

"gibbon"
99.3 % confidence

$\boldsymbol{x}$
"panda"
57.7% confidence

[Goodfellow et al. 2014]

Natural data $x$    +    Imperceptible adversarial perturbation    =    Adversarial data $\tilde{x}$



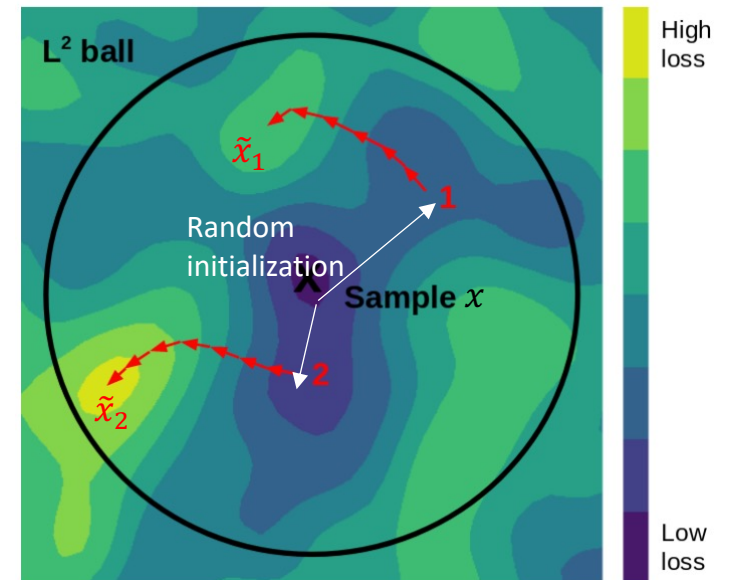Image modified from https://towardsdatascience.com/know-your-enemy-7f7c5038bdf3

Projected gradient descent (PGD)
[Madry et al. ICLR 2018]

# Supervised adversarial training (SAT)

- Minimax formulation of SAT

$$min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^{n} \ell(f(\tilde{x}_i), y_i), \text{ where } \tilde{x}_i = argmax_{\tilde{x}_i \in \mathcal{B}_\epsilon[x_i]} \ell(f(\tilde{x}_i), y_i)$$

outer minimization   [Madry et al. ICLR 2018]   inner maximization

- Realization

Alternatively conduct steps (1) and (2):
(1) generate adversarial data maximizing the loss;
(2) minimize loss on the generated adversarial data w.r.t. model parameters.

# SAT

- Minimax formulation of SAT

$$min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^{n} \ell(f(\tilde{x}_i), y_i), \text{ where } \tilde{x}_i = argmax_{\tilde{x}_i \in \mathcal{B}_\epsilon[x_i]} \ell(f(\tilde{x}_i), y_i)$$

outer minimization     [Madry et al. ICLR 2018]     inner maximization

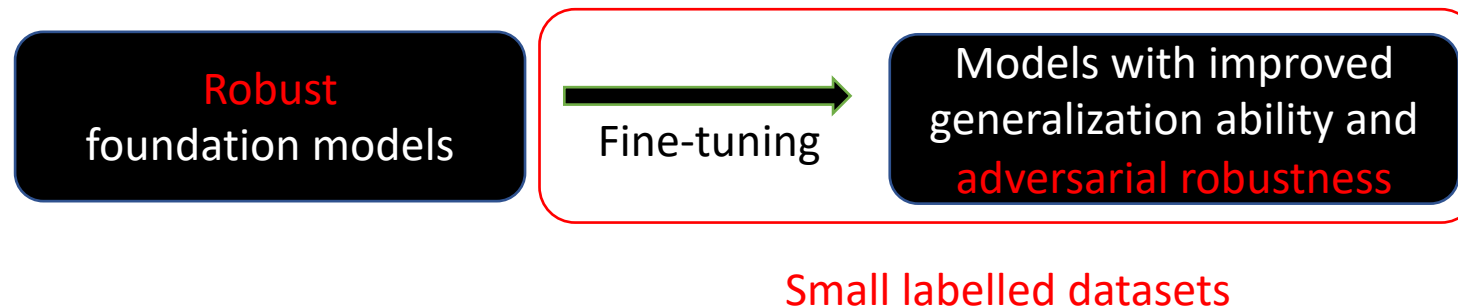- Realization

- Drawback: SAT requires a large amount of labelled data (for each task).

# SAT

- Minimax formulation of SAT

$$min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^{n} \ell(f(\tilde{x}_i), y_i), \text{ where } \tilde{x}_i = argmax_{\tilde{x}_i \in \mathcal{B}_\epsilon[x_i]} \ell(f(\tilde{x}_i), y_i)$$

outer minimization     [Madry et al. ICLR 2018]     inner maximization

- Realization

- Drawback: SAT requires a large amount of labelled data (for each task).



Robust foundation models → Fine-tuning → Models with improved generalization ability and adversarial robustness

Small labelled datasets

9

# Robust pre-training

Large-scale unlabelled dataset

Standard self-supervised pre-training
(e.g., standard contrastive learning)

Foundation models

Fine-tuning

Models with improved
generalization ability

# Robust pre-training

Large-scale unlabelled dataset

Standard self-supervised pre-training (e.g., standard contrastive learning)

Foundation models

Fine-tuning

Models with improved generalization ability

Large-scale unlabelled dataset

Robust self-supervised pre-training (e.g., adversarial contrastive learning)

Robust foundation models

Fine-tuning

Models with improved generalization ability and adversarial robustness

# Adversarial contrastive learning (ACL)



$t_i \in T$

$x_k$

$t_j \in T$

$x_k^i$

$x_k^j$

Standard contrastive learning (e.g., SimCLR)

$$\ell_{\mathrm{CL}}(x_k^i, x_k^j; \theta) = -\sum_{u \in \{i,j\}} \log \frac{e^{\mathrm{sim}\left(f_\theta(x_k^i), f_\theta(x_k^j)\right)/t}}{\sum_{x \in B^i \cup B^j \setminus \{x_k^u\}} e^{\mathrm{sim}\left(f_\theta(x_k^u), f_\theta(x)\right)/t}},$$

# ACL



Standard contrastive learning (e.g., SimCLR)

$$\ell_{\mathrm{CL}}(x_k^i, x_k^j; \theta) = -\sum_{u \in \{i,j\}} \log \frac{e^{\mathrm{sim}\left(f_\theta(x_k^i), f_\theta(x_k^j)\right)/t}}{\sum_{x \in B^i \cup B^j \setminus \{x_k^u\}} e^{\mathrm{sim}\left(f_\theta(x_k^u), f_\theta(x)\right)/t}},$$

PGD

The objective function of ACL

$$\ell_{\mathrm{ACL}}(x_k; \theta) = (1+\omega) \cdot \boxed{\ell_{\mathrm{CL}}(\tilde{x}_k^i, \tilde{x}_k^j; \theta)} + (1-\omega) \cdot \ell_{\mathrm{CL}}(x_k^i, x_k^j; \theta),$$

$$\text{where} \quad \tilde{x}_k^i, \tilde{x}_k^j = \arg\max_{\substack{\tilde{x}_k^i \in \mathcal{B}_\epsilon[x_k^i] \\ \tilde{x}_k^j \in \mathcal{B}_\epsilon[x_k^j]}} \ell_{\mathrm{CL}}(\tilde{x}_k^i, \tilde{x}_k^j; \theta),$$

# Outline

- Backgrounds
  - Adversarial robustness
  - Robust pre-training
- **How to make self-supervised robust pre-training**
  - **More efficient**
  - More effective
- Future directions

# Efficient ACL
# via Robustness-aware Coreset Selection (RCS)

- Why do we need to speed up ACL?
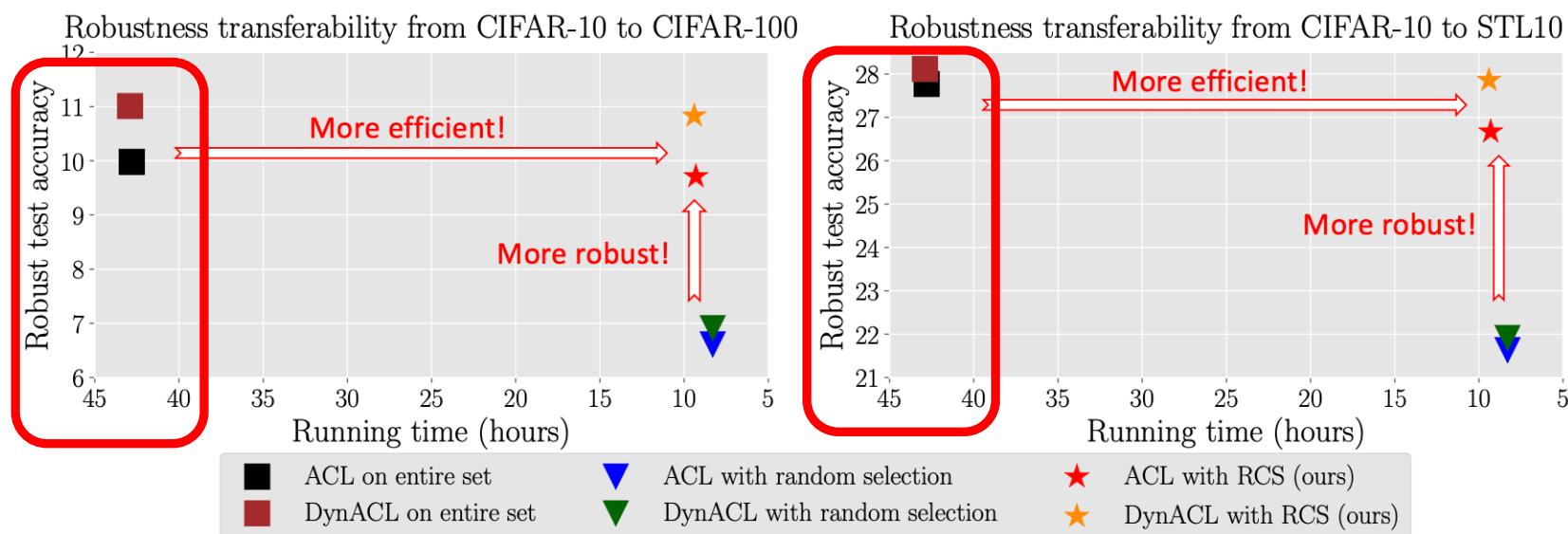  - ACL is extremely time-consuming.



Figure 1: We learn a representation using CIFAR-10 [4] dataset (without requiring labels) via ACL [12] and DynACL [15]. Then, we evaluate the representation's robustness transferability to CIFAR-100 [4] and STL10 [22] (using labels during finetuning) via standard linear finetuning. We demonstrate the running time of robust pre-training w.r.t. different coreset selection (CS) strategies and report the robust test accuracy under AutoAttack [17]. Experimental details are in Appendix B.4.

# Efficient ACL via RCS

- Why do we need to speed up ACL?
  - ACL is extremely time-consuming.
  - ACL has not been applied to ImageNet-1K yet due to computational prohibition.

Figure 4: Robustness evaluations on the CIFAR-10 (left three panels) and CIFAR-100 (right three panels) task. The number after the dash line denotes subset fraction $k \in \{0.05, 0.1, 0.2\}$.

Table 1: 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 | $84.36_{\pm 0.17}$ | $0.01_{\pm 0.01}$ | $10.00_{\pm 0.00}$ | $10.00_{\pm 0.00}$ | $86.63_{\pm 0.12}$ | $49.71_{\pm 0.16}$ |
| 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 | $\mathbf{70.02}_{\pm 0.12}$ | $\mathbf{22.45}_{\pm 0.13}$ | $\mathbf{63.94}_{\pm 0.21}$ | $\mathbf{31.13}_{\pm 0.17}$ | $\mathbf{85.23}_{\pm 0.23}$ | $\mathbf{52.21}_{\pm 0.14}$ |

Table 2: 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 | $57.34_{\pm 0.23}$ | $0.01_{\pm 0.01}$ | $9.32_{\pm 0.01}$ | $0.06_{\pm 0.01}$ | $61.33_{\pm 0.12}$ | $25.11_{\pm 0.15}$ |
| 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 | $\mathbf{40.28}_{\pm 0.17}$ | $\mathbf{14.55}_{\pm 0.10}$ | $\mathbf{33.15}_{\pm 0.26}$ | $\mathbf{14.89}_{\pm 0.16}$ | $\mathbf{60.25}_{\pm 0.18}$ | $\mathbf{28.24}_{\pm 0.13}$ |

16

# Efficient ACL via RCS: Methodology

- Idea: Find an informative training subset
  - Decreasing the number of training samples
  - Preserving the robust representations

# Efficient ACL via RCS: Methodology

- Idea: Find an informative training subset

- Intuitive solution: selects training data from the entire set whose gradients are most beneficial to maintaining adversarial robustness.

# Efficient ACL via RCS: Methodology

- Idea: Find an informative training subset
- Intuitive solution: selects training data from the entire set whose gradients are most beneficial to maintaining adversarial robustness.

- Representational divergence (RD)
  - The smaller the RD is, the representations are of less sensitivity to adversarial perturbations, thus being more robust.

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

# Efficient ACL via RCS: Methodology

- Idea: Find an informative training subset

- Intuitive solution: selects training data from the entire set whose gradients are most beneficial to maintaining adversarial robustness.

- Representational divergence (RD)

- **Objective function of RCS**

Unlabeled validation set

$$S^* = \underset{S \subseteq X, |S|/|X| \leq k}{\arg\min} \mathcal{L}_{\mathrm{RD}}(U; \underset{\theta}{\arg\min} \mathcal{L}_{\mathrm{ACL}}(S; \theta))$$

Coreset

Subset
fraction

Representational
divergence (RD)

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

# Efficient ACL via RCS: Methodology

- Solve the objective function of RCS
  - Transformation of RCS

$$S^* = \underset{S \subseteq X, |S|/|X| \leq k}{\arg\min} \mathcal{L}_{\mathrm{RD}}(U; \underset{\theta}{\arg\min} \mathcal{L}_{\mathrm{ACL}}(S; \theta))$$

One-step gradient approximation

$$S^* = \underset{S \subseteq X, |S|/|X| \leq k}{\arg\min} \mathcal{L}_{\mathrm{RD}}(U; \theta - \eta \nabla_\theta \mathcal{L}_{\mathrm{ACL}}(S; \theta))$$

Transform into a problem of maximizing a set function subject to a cardinality constraint

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

$$G_\theta(S \subseteq X) \triangleq -\mathcal{L}_{\mathrm{RD}}(U; \theta - \eta \nabla_\theta \mathcal{L}_{\mathrm{ACL}}(S; \theta))$$

# Efficient ACL via RCS: Methodology

- Solve the objective function of RCS
  - Transformation of RCS $\quad S^* = \underset{S \subseteq X, |S|/|X|=k}{\arg\max} G_\theta(S) \quad G_\theta(S \subseteq X) \triangleq -\mathcal{L}_{\mathrm{RD}}(U; \theta - \eta \nabla_\theta \mathcal{L}_{\mathrm{ACL}}(S; \theta))$
  - Greedy search for solving a proxy set problem

$$\hat{S}^* = \underset{S \subseteq X, |S|/|X|=k}{\arg\max} \hat{G}_\theta(S)$$

**Theorem 1.** *We define a proxy set function $\hat{G}_\theta(S) \triangleq G_\theta(S) + |S|\sigma$, where $\sigma = 1 + \nu_1 + \nu_2 L_2 + \eta M L_2(L_1 + \eta k N(L_1 L_4 + L_2 L_3))$, $\nu_1 \to 0^+$, and $\nu_2 > 0$ are positive constants. Given Assumption 1, $\hat{G}_\theta(S)$ is monotone and $\gamma$-weakly submodular where $\gamma > \gamma^* = \frac{1}{2\sigma - 1}$.*

# Efficient ACL via RCS: Methodology

- Solve the objective function of RCS
  - Transformation of RCS
  - Greedy search for solving a proxy set problem

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

$$\hat{S}^* = \operatorname*{arg\,max}_{S \subseteq X, |S|/|X|=k} \hat{G}_\theta(S)$$

  - Guaranteed lower bound of the original problem by solving the proxy set problem

**Theorem 2.** *Given a fixed parameter $\theta$, we denote the optimal solution of Eq. (5) as $G_\theta^* = \sup_{S \subseteq X, |S|/|X|=k} G_\theta(S)$. Then, $\hat{S}^*$ in Eq. (6) found via greedy search satisfies*

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

# Efficient ACL via RCS: Methodology

- ## Solve the objective function of RCS

  - Transformation of RCS

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

  - Greedy search for solving a proxy set problem

  $$\hat{S}^* = \underset{S \subseteq X, |S|/|X|=k}{\arg\max} \hat{G}_\theta(S)$$

  - Guaranteed lower bound of the original problem by solving the proxy set problem

  - ## Algorithm

---

**Algorithm 1** Robustness-aware Coreset Selection (RCS)

1: **Input:** Unlabeled training set $X$, unlabeled validation set $U$, batch size $\beta$, model $g \circ f_\theta$, learning rate for RCS $\eta$, subset fraction $k \in (0, 1]$
2: **Output:** Coreset $S$
3: Initialize $S \leftarrow \emptyset$
4: Split entire set into minibatches $X = \{B_m\}_{m=1}^{\lceil |X|/\beta \rceil}$
5: **for** each minibatch $B_m \in X$ **do**
6:     Compute gradient $q_m \leftarrow \nabla_\theta \mathcal{L}_{\mathrm{ACL}}(B_m; \theta)$
7: **end for**
8: // Conduct greedy search via batch-wise selection
9: **for** $1, \ldots, \lfloor k|X|/\beta \rfloor$ **do**
10:     Compute gradient $q_U \leftarrow \nabla_\theta \mathcal{L}_{\mathrm{RD}}(U; \theta)$
11:     Initialize $best\_gain = -\infty$
12:     **for** each minibatch $B_m \in X$ **do**
13:         Compute marginal gain $\hat{G}(B_m|S) \leftarrow \eta q_U^\top q_m$
14:         **if** $\hat{G}(B_m|S) > best\_gain$ **then**
15:             Update $s \leftarrow m, best\_gain \leftarrow \hat{G}(B_m|S)$
16:         **end if**
17:     **end for**
18:     Update $S \leftarrow S \cup B_s, X \leftarrow X \setminus B_s$
19:     Update $\theta \leftarrow \theta - \eta q_s$
20: **end for**

---

**Algorithm 2** Efficient ACL via RCS

1: **Input:** Unlabeled training set $X$, unlabeled validation set $U$, total training epochs $E$, learning rate $\eta'$, batch size $\beta$, warmup epoch $\omega$, epoch interval for executing RCS $\lambda$, subset fraction $k$, learning rate for RCS $\eta$
2: **Output:** Adversarially pre-trained feature extractor $f_\theta$
3: Initialize parameters of model $g \circ f_\theta$
4: Initialize training set $S \leftarrow X$
5: **for** $e = 0$ to $E - 1$ **do**
6:     **if** $e \% \lambda == 0$ **and** $e \geq \omega$ **then**
7:         $S \leftarrow \mathrm{RCS}(X, U, \beta, g \circ f_\theta, \eta, k)$ //by Algorithm 1
8:     **end if**
9:     **for** batch $m = 1, \ldots, \lceil |S|/\beta \rceil$ **do**
10:         Sample a minibatch $B_m$ from $S$
11:         Update $\theta \leftarrow \theta - \eta' \nabla_\theta \mathcal{L}_{\mathrm{ACL}}(B_m; \theta)$
12:     **end for**
13: **end for**

# Efficient ACL via RCS: Empirical results

- Our proposed RCS is
  - more efficient (higher speed-up ratio)
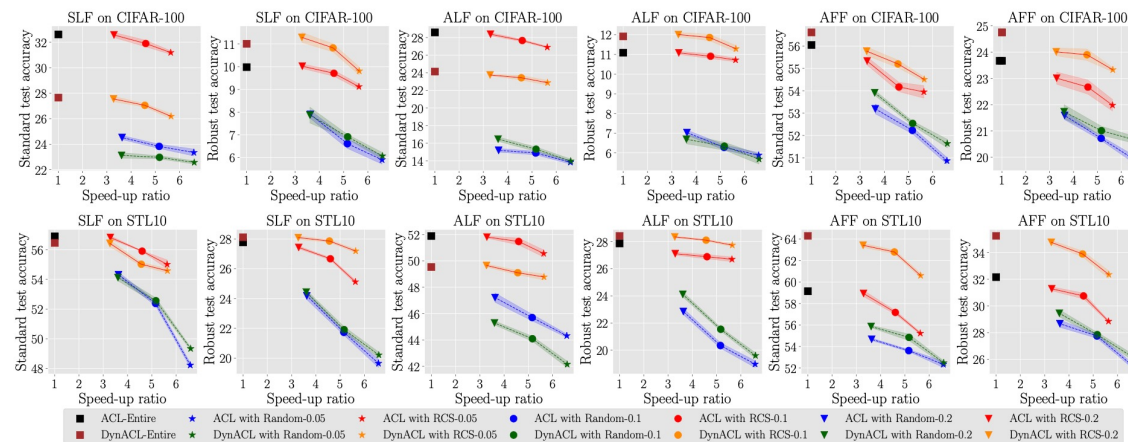  - more effective (higher test accuracy)



Figure 2: Robustness transferability from CIFAR-10 to CIFAR-100 (upper row) and STL10 (bottom row). The number after the dash line denotes subset fraction $k \in \{0.05, 0.1, 0.2\}$.
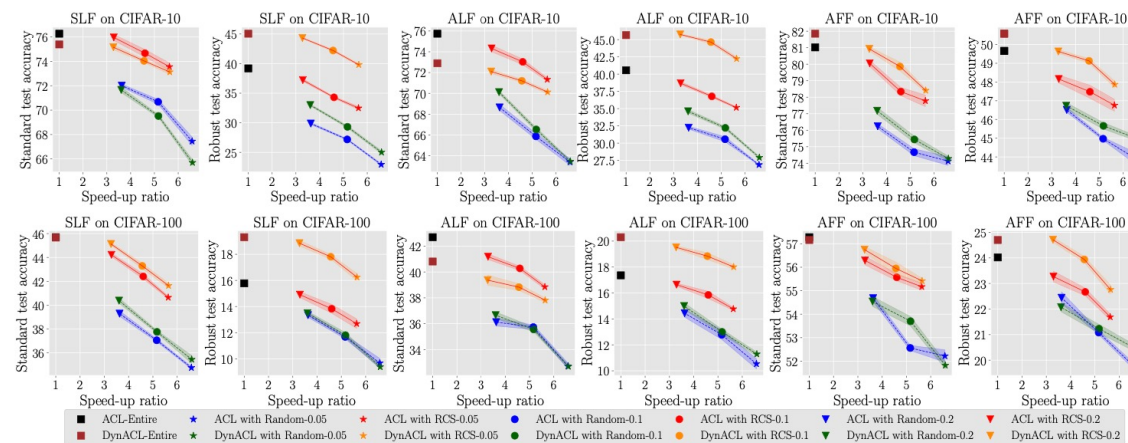


Figure 4: Robustness evaluations on the CIFAR-10 (left three panels) and CIFAR-100 (right three panels) task. The number after the dash line denotes subset fraction $k \in \{0.05, 0.1, 0.2\}$.

The upper-right (ours) is better!

# Efficient ACL via RCS: Empirical results

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

Figure 4: Robustness evaluations on the CIFAR-10 (left three panels) and CIFAR-100 (right three panels) task. The number after the dash line denotes subset fraction $k \in \{0.05, 0.1, 0.2\}$.

Table 1: 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 | $84.36_{\pm 0.17}$ | $0.01_{\pm 0.01}$ | $10.00_{\pm 0.00}$ | $10.00_{\pm 0.00}$ | $86.63_{\pm 0.12}$ | $49.71_{\pm 0.16}$ |
| 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 | $\mathbf{70.02}_{\pm 0.12}$ | $\mathbf{22.45}_{\pm 0.13}$ | $\mathbf{63.94}_{\pm 0.21}$ | $\mathbf{31.13}_{\pm 0.17}$ | $\mathbf{85.23}_{\pm 0.23}$ | $\mathbf{52.21}_{\pm 0.14}$ |

Table 2: 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 | $57.34_{\pm 0.23}$ | $0.01_{\pm 0.01}$ | $9.32_{\pm 0.01}$ | $0.06_{\pm 0.01}$ | $61.33_{\pm 0.12}$ | $25.11_{\pm 0.15}$ |
| 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 | $\mathbf{40.28}_{\pm 0.17}$ | $\mathbf{14.55}_{\pm 0.10}$ | $\mathbf{33.15}_{\pm 0.26}$ | $\mathbf{14.89}_{\pm 0.16}$ | $\mathbf{60.25}_{\pm 0.18}$ | $\mathbf{28.24}_{\pm 0.13}$ |

# Efficient ACL via RCS: Empirical results

- RCS for speeding up SAT on ImageNet-1K (supervised setting)
  - Maintaining standard transferability



  - Maintaining robustness transferability

Table 18: Standard transferability [43] of adversarially pre-trained ResNet-50 from ImageNet-1K to CIFAR-10 and CIFAR-100, respectively. We report the standard test accuracy (%) via standard linear finetuning (SLF) and standard full finetuning (SFF). The number after the dash line denotes subset fraction $k \in \{0.05, 0.1, 0.2\}$.

| Pre-training | Runing time (hours) | CIFAR-10 | | CIFAR-100 | |
|---|---|---|---|---|---|
| | | SLF | SFF | SLF | SFF |
| Standard training [43] on entire set | - | 78.84 | 97.41 | 57.09 | 84.21 |
| SAT [43] on entire set | 286.1 | 93.53 | 98.09 | 77.29 | 86.99 |
| Fast-AT [20] on entire set | 10.4 | 90.91 | 97.54 | 73.35 | 83.33 |
| SAT with Random-0.05 | 38.7 | 85.72 | 95.27 | 69.29 | 82.34 |
| SAT with RCS-0.05 | **48.2** | **92.68** | **97.65** | **75.35** | **84.71** |
| SAT with Random-0.1 | 45.8 | 87.14 | 95.60 | 71.23 | 83.62 |
| SAT with RCS-0.1 | **55.4** | **92.92** | **97.82** | **75.41** | **85.22** |
| SAT with Random-0.2 | 70.3 | 87.69 | 96.10 | 72.05 | 84.14 |
| SAT with RCS-0.2 | **79.8** | **93.48** | **98.06** | **76.39** | **85.44** |

Table 16: 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. [51]. 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 [51] 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 | **68.6** | **79.72** | **44.44** | **87.99** | **56.87** |
| SAT with Random-0.1 | 70.2 | 73.28 | 33.86 | 86.78 | 54.95 |
| SAT with RCS-0.1 | **81.9** | **81.92** | **45.10** | **88.87** | **57.69** |
| SAT with Random-0.2 | 103.4 | 75.46 | 39.62 | 86.64 | 56.46 |
| SAT with RCS-0.2 | **121.9** | **83.94** | **46.88** | **89.54** | **58.13** |

# Efficient ACL via RCS: Conclusions

- We proposed <span style="color:red">robustness-aware coreset selection</span> (RCS) that can
  - <span style="color:red">speed up</span> (supervised and self-supervised) <span style="color:red">robust pre-training</span>
  - <span style="color:red">maintain</span> (standard and robustness) <span style="color:red">transferability</span>

# Outline

- Backgrounds
  - Adversarial attack and defense
  - Robust pre-training
- **How to make self-supervised robust pre-training**
  - More efficient
  - **More effective**
- Future directions

# Effective ACL
# via adversarial invariant regularization (AIR)

- Motivation
  - The style-independence property of learned representations, which eliminates the effects of nuisance style factors in standard contrastive learning (SCL), has been shown to significantly improve the transferability of representations.

| Algorithm | Shorthand | Paper | KNN accuracy |
|---|---|---|---|
| Bootstrap Your Own Latent: A new approach to self-supervised Learning | BYOL | arXiv | 80.09 |
| Representation Learning via Invariant Causal Mechanisms | ReLIC | arXiv | 79.26 |
| A Simple Framework for Contrastive Learning of Visual Representations | SimCLR | arXiv | 77.79 |
| Unsupervised Learning of Visual Features by Contrasting Cluster Assignments | SwAV | arXiv | 72.11 |
| Momentum Contrast for Unsupervised Visual Representation Learning | MoCo | arXiv | 63.14 |
| Barlow Twins: Self-Supervised Learning via Redundancy Reduction | Barlow | arXiv | 56.81 |

Performance evaluated on CIFAR-10

Image from https://github.com/NightShade99/Self-Supervised-Vision

[Mitrovic et al., ICLR 2021]

# Effective ACL via AIR

- Motivation
  - The style-independence property of learned representations, which eliminates the effects of nuisance style factors in standard contrastive learning (SCL), has been shown to improve the transferability of representations.

It is unclear how the style-independence property benefits ACL-learned robust representations.

# Effective ACL via AIR : Methodology

- ACL in the view of causality

$$\ell_{\mathrm{ACL}}(x_k; \theta) = (1 + \omega) \cdot \ell_{\mathrm{CL}}(\tilde{x}_k^i, \tilde{x}_k^j; \theta) + (1 - \omega) \cdot \ell_{\mathrm{CL}}(x_k^i, x_k^j; \theta),$$

$$\text{where} \quad \tilde{x}_k^i, \tilde{x}_k^j = \arg\max_{\substack{\tilde{x}_k^i \in \mathcal{B}_\epsilon[x_k^i] \\ \tilde{x}_k^j \in \mathcal{B}_\epsilon[x_k^j]}} \ell_{\mathrm{CL}}(\tilde{x}_k^i, \tilde{x}_k^j; \theta),$$
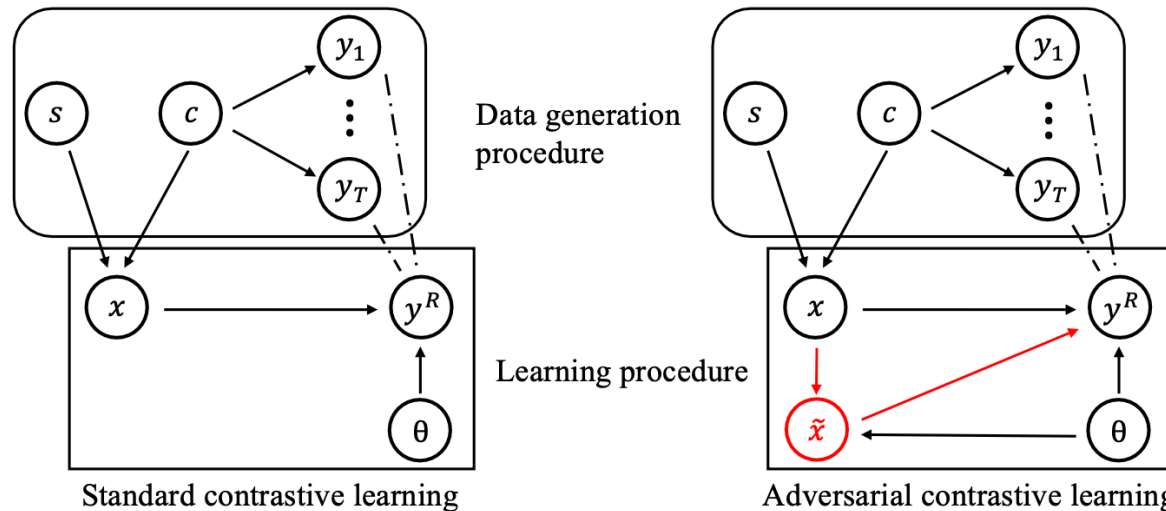


Figure 1: Causal graph of standard contrastive learning [35] (left panel) and adversarial contrastive learning (right panel). $x$ is unlabeled data, $s$ is style variable, $c$ is content variable, $\tilde{x}$ is the generated adversarial data, and $\theta$ is the parameter of representation. The dashdotted lines denote that the proxy label $y^R \in \mathcal{Y}^R$ is a refinement of the target label $y_t \in \mathcal{Y} = \{y_i\}_{i=1}^T$. All other arrows are causal.

# Effective ACL via AIR : Methodology

- ACL in the view of causality

$$\ell_{\mathrm{ACL}}(x_k; \theta) = (1 + \omega) \cdot \ell_{\mathrm{CL}}(\tilde{x}_k^i, \tilde{x}_k^j; \theta) + (1 - \omega) \cdot \ell_{\mathrm{CL}}(x_k^i, x_k^j; \theta),$$

$$\text{where} \quad \tilde{x}_k^i, \tilde{x}_k^j = \arg\max_{\substack{\tilde{x}_k^i \in \mathcal{B}_\epsilon[x_k^i] \\ \tilde{x}_k^j \in \mathcal{B}_\epsilon[x_k^j]}} \ell_{\mathrm{CL}}(\tilde{x}_k^i, \tilde{x}_k^j; \theta),$$



Data generation procedure

Learning procedure

Standard contrastive learning

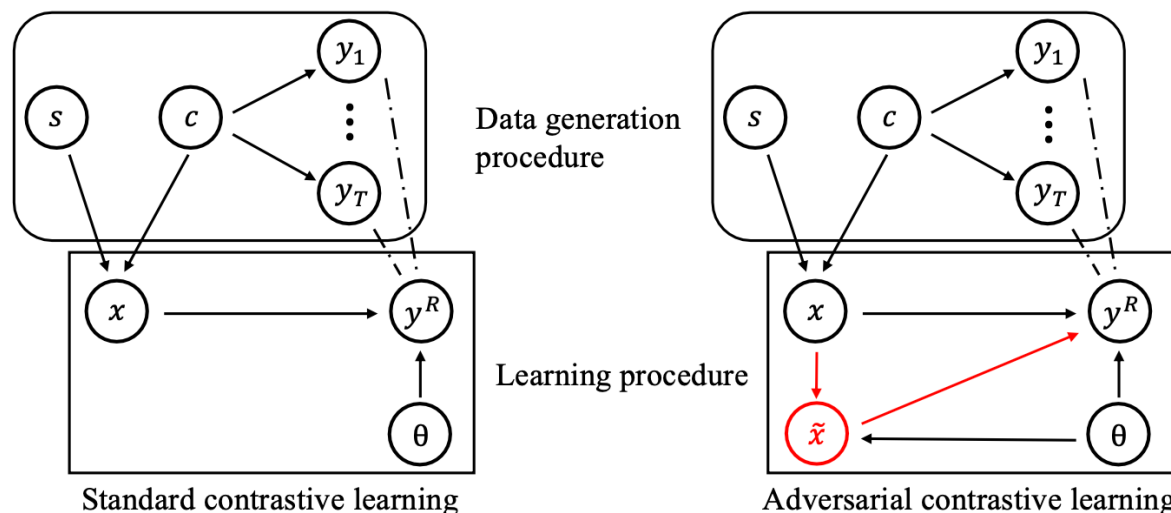Adversarial contrastive learning

Figure 1: Causal graph of standard contrastive learning [35] (left panel) and adversarial contrastive learning (right panel). $x$ is unlabeled data, $s$ is style variable, $c$ is content variable, $\tilde{x}$ is the generated adversarial data, and $\theta$ is the parameter of representation. The dashdotted lines denote that the proxy label $y^R \in \mathcal{Y}^R$ is a refinement of the target label $y_t \in \mathcal{Y} = \{y_i\}_{i=1}^T$. All other arrows are causal.
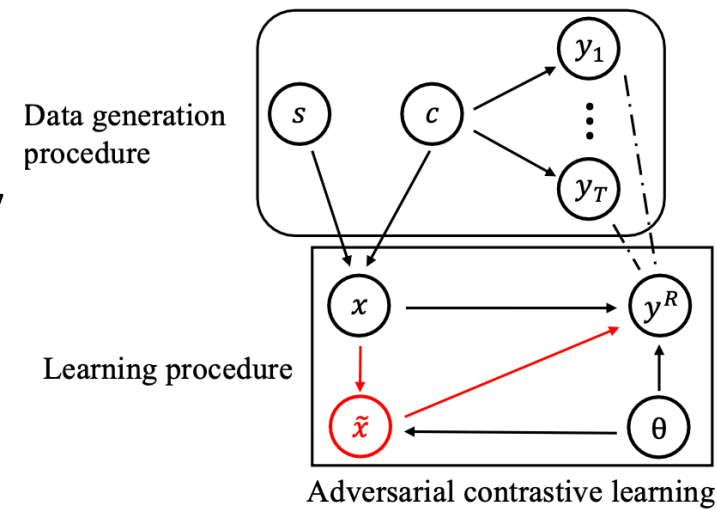
The rationality of the causal graph

**Theorem 1.** *The learning objective of the proxy task used in ACL which is to maximize the conditional probability both $p(y^R|x)$ and $p(y^R|\tilde{x})$ is equivalent to the learning objective of ACL [26] which is to minimize the sum of standard and adversarial contrastive losses.*

# Effective ACL via AIR: Methodology



Data generation procedure

Learning procedure

Adversarial contrastive learning

- Adversarial invariant regularization (AIR)
  - The conditional probability learned via ACL

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

# Effective ACL via AIR: Methodology



Data generation procedure

Learning procedure

Adversarial contrastive learning

- Adversarial invariant regularization (AIR)
  - The conditional probability learned via ACL  $p(y^R|x) = p(y^R|\tilde{x})p(\tilde{x}|x)$
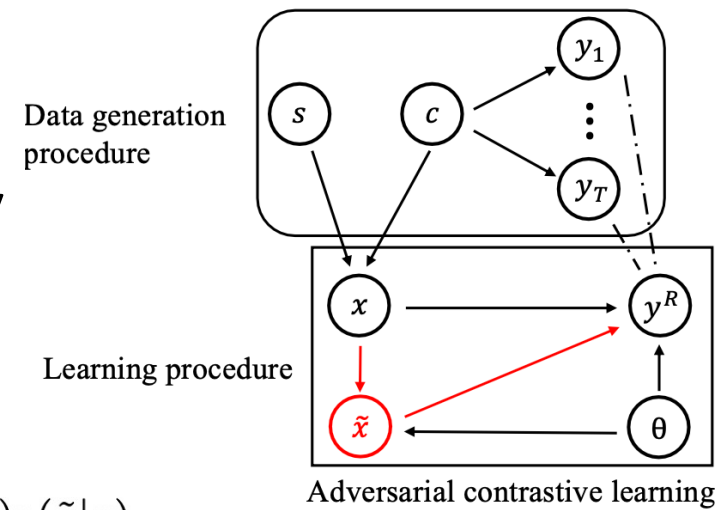  - Style-independent criterion

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

$$p^{do(\tau_u)}(y^R|\tilde{x}) = \frac{e^{\mathrm{sim}(f_\theta(x), f_\theta(\tilde{x}^u))/t}}{\sum\limits_{x_k \in B} e^{\mathrm{sim}(f_\theta(x_k), f_\theta(\tilde{x}_k^u))/t}}, \quad p^{do(\tau_u)}(\tilde{x}|x) = \frac{e^{\mathrm{sim}(f_\theta(\tilde{x}^u), f_\theta(x^u))/t}}{\sum\limits_{x_k \in B} e^{\mathrm{sim}(f_\theta(\tilde{x}_k^u), f_\theta(x_k^u))/t}}$$

# Effective ACL via AIR: Methodology



Data generation procedure

Learning procedure

Adversarial contrastive learning

- ## Adversarial invariant regularization (AIR)

  - The conditional probability learned via ACL $\quad p(y^R|x) = p(y^R|\tilde{x})p(\tilde{x}|x)$

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

  - Loss function of AIR

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

$$p^{do(\tau_u)}(y^R|\tilde{x}) = \frac{e^{\text{sim}(f_\theta(x), f_\theta(\tilde{x}^u))/t}}{\sum_{x_k \in B} e^{\text{sim}(f_\theta(x_k), f_\theta(\tilde{x}_k^u))/t}}, \quad p^{do(\tau_u)}(\tilde{x}|x) = \frac{e^{\text{sim}(f_\theta(\tilde{x}^u), f_\theta(x^u))/t}}{\sum_{x_k \in B} e^{\text{sim}(f_\theta(\tilde{x}_k^u), f_\theta(x_k^u))/t}}$$
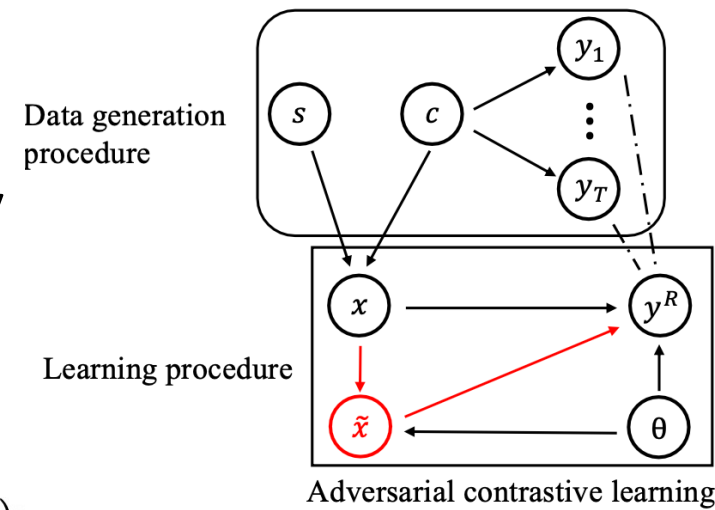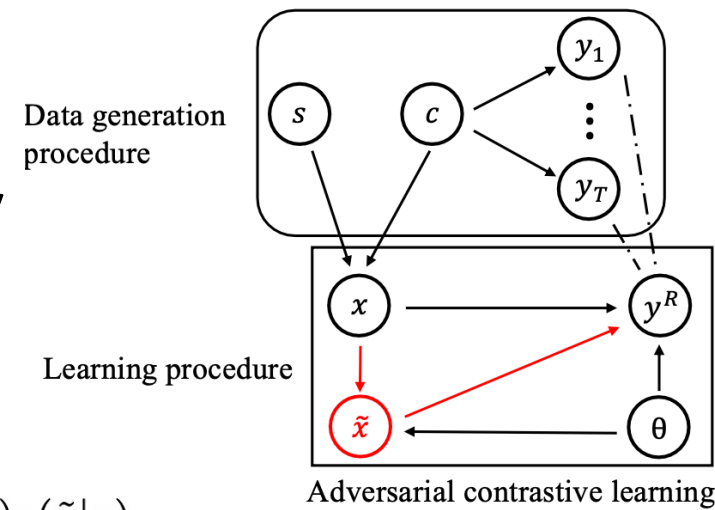
# Effective ACL via AIR: Methodology



Data generation procedure

Learning procedure

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    $p^{do(\tau_i)}(y^R|\tilde{x})p^{do(\tau_i)}(\tilde{x}|x) = p^{do(\tau_j)}(y^R|\tilde{x})p^{do(\tau_j)}(\tilde{x}|x) \quad \forall \tau_i, \tau_j \in \mathcal{T},$
  - Loss function of AIR    $\mathcal{L}_{\mathrm{AIR}}(B;\theta) = \mathrm{KL}\left(p^{do(\tau_i)}(y^R|\tilde{x})p^{do(\tau_i)}(\tilde{x}|x)\|p^{do(\tau_j)}(y^R|\tilde{x})p^{do(\tau_j)}(\tilde{x}|x); B\right)$

  - Standard invariant regularization (SIR): a special case of AIR

$$\mathcal{L}_{\mathrm{SIR}}(B;\theta) = \mathrm{KL}\left(p^{do(\tau_i)}(y^R|x)\|p^{do(\tau_i)}(y^R|x); B\right),$$

$$\text{where} \quad p^{do(\tau_u)}(y^R|x) = \frac{e^{\mathrm{sim}(f_\theta(x), f_\theta(x^u))/t}}{\sum_{x_k \in B} e^{\mathrm{sim}\left(f_\theta(x_k), f_\theta(x_k^u)\right)/t}} \quad \forall u \in \{i, j\}$$

# Effective ACL via AIR: Methodology



Data generation procedure

Learning procedure

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    $p^{do(\tau_i)}(y^R|\tilde{x})p^{do(\tau_i)}(\tilde{x}|x) = p^{do(\tau_j)}(y^R|\tilde{x})p^{do(\tau_j)}(\tilde{x}|x)$    $\forall \tau_i, \tau_j \in \mathcal{T},$
  - Loss function of AIR    $\mathcal{L}_{\mathrm{AIR}}(B;\theta) = \mathrm{KL}\left(p^{do(\tau_i)}(y^R|\tilde{x})p^{do(\tau_i)}(\tilde{x}|x)\|p^{do(\tau_j)}(y^R|\tilde{x})p^{do(\tau_j)}(\tilde{x}|x); B\right)$
  - SIR: a special case of AIR    $\mathcal{L}_{\mathrm{SIR}}(B;\theta) = \mathrm{KL}\left(p^{do(\tau_i)}(y^R|x)\|p^{do(\tau_i)}(y^R|x); B\right)$

- Our proposed invariant regularization (IR)

$$\arg\min_{\theta} \sum_{x \in U} \ell_{\mathrm{ACL}}(x;\theta) + \underbrace{\lambda_1 \cdot \mathcal{L}_{\mathrm{SIR}}(U;\theta) + \lambda_2 \cdot \mathcal{L}_{\mathrm{AIR}}(U;\theta)}_{\textit{invariant regularization}},$$

# Effective ACL via AIR: Theoretical analysis

- Theoretical justification of the effectiveness
  - The style-independence property is generalizable to the downstream tasks

**Proposition 4.** *Let $\mathcal{Y} = \{y_t\}_{t=1}^T$ be a label set of a downstream classification task, $\mathcal{Y}^R$ be a refinement of $\mathcal{Y}$, and $\tilde{x}_t$ be the adversarial data generated on the downstream task. Assuming that $\tilde{x} \in \mathcal{B}_\epsilon[x]$ and $\tilde{x}_t \in \mathcal{B}_\epsilon[x]$, we have the following results:*

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

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

# Effective ACL via AIR: Theoretical analysis

- Theoretical justification of the effectiveness
  - The style-independence property is generalizable to the downstream tasks

**Proposition 4.** *Let $\mathcal{Y} = \{y_t\}_{t=1}^{T}$ be a label set of a downstream classification task, $\mathcal{Y}^R$ be a refinement of $\mathcal{Y}$, and $\tilde{x}_t$ be the adversarial data generated on the downstream task. Assuming that $\tilde{x} \in \mathcal{B}_\epsilon[x]$ and $\tilde{x}_t \in \mathcal{B}_\epsilon[x]$, we have the following results:*

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

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

- We can treat adversarial attacks and common corruptions as style factors
- IR regulates the representations to be invariant of style factors

# Effective ACL via AIR: Experimental results

- Performance evaluated on various tasks

- Performance evaluated via various fine-tuning methods

- Robustness under common corruption

Table 1: Robustness evaluations via SLF across various tasks.

| Pre-training | $\lambda_1$ | $\lambda_2$ | CIFAR-10 | | CIFAR-100 | | STL10 | |
| | | | AA (%) | SA (%) | AA (%) | SA (%) | AA (%) | SA (%) |
|---|---|---|---|---|---|---|---|---|
| ACL [26] | 0.0 | 0.0 | $37.39_{\pm0.06}$ | $78.27_{\pm0.09}$ | $15.78_{\pm0.05}$ | $45.70_{\pm0.09}$ | $35.80_{\pm0.06}$ | $67.90_{\pm0.09}$ |
| ACL with SIR [35] | 0.5 | 0.0 | $37.51_{\pm0.04}$ | $78.97_{\pm0.08}$ | $15.76_{\pm0.06}$ | $47.16_{\pm0.11}$ | $36.36_{\pm0.09}$ | $68.09_{\pm0.13}$ |
| ACL with AIR | 0.0 | 0.5 | $38.70_{\pm0.09}$ | $79.96_{\pm0.05}$ | $16.03_{\pm0.12}$ | $49.60_{\pm0.15}$ | $36.86_{\pm0.08}$ | $68.61_{\pm0.10}$ |
| ACL with IR | 0.5 | 0.5 | $\mathbf{38.89}_{\pm0.06}$ | $\mathbf{80.03}_{\pm0.07}$ | $\mathbf{16.14}_{\pm0.07}$ | $\mathbf{49.75}_{\pm0.10}$ | $\mathbf{36.94}_{\pm0.06}$ | $\mathbf{68.91}_{\pm0.07}$ |
| DynACL [19] | 0.0 | 0.0 | $45.05_{\pm0.04}$ | $75.39_{\pm0.05}$ | $19.31_{\pm0.06}$ | $45.67_{\pm0.09}$ | $46.49_{\pm0.05}$ | $69.59_{\pm0.08}$ |
| DynACL with SIR [35] | 0.5 | 0.0 | $44.70_{\pm0.03}$ | $76.45_{\pm0.06}$ | $19.67_{\pm0.09}$ | $46.13_{\pm0.10}$ | $46.56_{\pm0.08}$ | $70.41_{\pm0.09}$ |
| DynACL with AIR | 0.0 | 0.5 | $45.23_{\pm0.08}$ | $78.01_{\pm0.11}$ | $20.37_{\pm0.08}$ | $46.77_{\pm0.11}$ | $47.62_{\pm0.07}$ | $71.98_{\pm0.12}$ |
| DynACL with IR | 0.5 | 0.5 | $\mathbf{45.27}_{\pm0.04}$ | $\mathbf{78.08}_{\pm0.06}$ | $\mathbf{20.45}_{\pm0.07}$ | $\mathbf{46.84}_{\pm0.12}$ | $\mathbf{47.66}_{\pm0.06}$ | $\mathbf{72.30}_{\pm0.10}$ |

Table 2: Robustness benchmark on the CIFAR-10 task evaluated via SLF, ALF, and AFF.

| Pre-training | $\lambda_1$ | $\lambda_2$ | SLF | | ALF | | AFF | |
| | | | AA (%) | SA (%) | AA (%) | SA (%) | AA (%) | SA (%) |
|---|---|---|---|---|---|---|---|---|
| ACL [26] | 0.0 | 0.0 | $37.39_{\pm0.06}$ | $78.27_{\pm0.09}$ | $40.61_{\pm0.07}$ | $75.56_{\pm0.09}$ | $49.42_{\pm0.07}$ | $82.14_{\pm0.18}$ |
| ACL with SIR [35] | 0.5 | 0.0 | $37.51_{\pm0.04}$ | $78.97_{\pm0.08}$ | $40.30_{\pm0.08}$ | $76.49_{\pm0.05}$ | $50.36_{\pm0.07}$ | $82.62_{\pm0.08}$ |
| ACL with AIR | 0.0 | 0.5 | $38.70_{\pm0.09}$ | $79.96_{\pm0.05}$ | $41.09_{\pm0.06}$ | $77.99_{\pm0.12}$ | $50.32_{\pm0.09}$ | $82.67_{\pm0.09}$ |
| ACL with IR | 0.5 | 0.5 | $\mathbf{38.89}_{\pm0.06}$ | $\mathbf{80.03}_{\pm0.07}$ | $\mathbf{41.39}_{\pm0.08}$ | $\mathbf{78.29}_{\pm0.10}$ | $\mathbf{50.44}_{\pm0.04}$ | $\mathbf{82.71}_{\pm0.06}$ |
| DynACL [33] | 0.0 | 0.0 | $45.05_{\pm0.04}$ | $75.39_{\pm0.05}$ | $45.65_{\pm0.05}$ | $72.90_{\pm0.08}$ | $50.52_{\pm0.05}$ | $81.86_{\pm0.11}$ |
| DynACL with SIR [35] | 0.5 | 0.0 | $44.70_{\pm0.03}$ | $76.45_{\pm0.06}$ | $45.42_{\pm0.10}$ | $74.78_{\pm0.14}$ | $50.58_{\pm0.07}$ | $81.66_{\pm0.18}$ |
| DynACL with AIR | 0.0 | 0.5 | $45.23_{\pm0.08}$ | $78.01_{\pm0.11}$ | $46.12_{\pm0.09}$ | $77.01_{\pm0.12}$ | $50.66_{\pm0.05}$ | $82.62_{\pm0.10}$ |
| DynACL with IR | 0.5 | 0.5 | $\mathbf{45.27}_{\pm0.04}$ | $\mathbf{78.08}_{\pm0.06}$ | $\mathbf{46.14}_{\pm0.07}$ | $\mathbf{77.42}_{\pm0.10}$ | $\mathbf{50.68}_{\pm0.08}$ | $\mathbf{82.74}_{\pm0.11}$ |

Table 3: Test accuracy (%) evaluated on CIFAR-10-C (corruption severity ranges from 1 to 5) of CIFAR-10 pre-trained models after SLF and AFF, respectively. Standard deviation is in Table 20.

| Pre-training | $\lambda_1$ | $\lambda_2$ | SLF | | | | | AFF | | | | |
| | | | CS-1 | CS-2 | CS-3 | CS-4 | CS-5 | CS-1 | CS-2 | CS-3 | CS-4 | CS-5 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ACL [26] | 0.0 | 0.0 | 76.57 | 74.73 | 71.78 | 67.75 | 62.78 | 79.15 | 76.01 | 72.54 | 69.47 | 65.27 |
| ACL with SIR [35] | 0.5 | 0.0 | 77.31 | 75.46 | 72.21 | 68.14 | 63.27 | 79.05 | 76.29 | 72.73 | 69.43 | 65.29 |
| ACL with AIR | 0.0 | 0.5 | 78.30 | 76.34 | 73.27 | 69.10 | 64.24 | 79.24 | 76.54 | 72.81 | 69.64 | 65.32 |
| ACL with IR | 0.5 | 0.5 | **78.55** | **76.67** | **73.33** | **69.12** | **64.28** | **79.49** | **76.86** | **72.95** | **69.73** | **65.37** |
| DynACL [33] | 0.0 | 0.0 | 73.92 | 71.69 | 69.01 | 66.22 | 62.51 | 79.77 | 76.44 | 72.95 | 69.74 | 65.60 |
| DynACL with SIR [35] | 0.5 | 0.0 | 75.81 | 72.88 | 69.31 | 66.24 | 62.20 | 80.59 | 77.31 | 73.67 | 70.39 | 66.05 |
| DynACL with AIR | 0.0 | 0.5 | 76.33 | 73.46 | 69.97 | 67.19 | 63.13 | 80.93 | 77.71 | 74.11 | 70.81 | 66.58 |
| DynACL with IR | 0.5 | 0.5 | **76.62** | **73.62** | **70.16** | **67.37** | **63.29** | **80.98** | **77.87** | **74.31** | **70.96** | **66.75** |

# Effective ACL via AIR: Conclusions

- We proposed an <span style="color:red">invariant regularization</span> that can
  - <span style="color:red">regulate</span> (both standard and robust) <span style="color:red">representations to be style-independent</span>
  - <span style="color:red">improve both generalization ability and robustness transferability</span> against adversarial attacks and common corruptions

# Thank you for your attention!

- Summary
  - More efficient robust pre-training via robustness-aware coreset selection
  - More effective robust pre-training via adversarial invariant regularization

- Future directions
  - The application of robust foundation models in computer vision tasks
    - Segmentation
    - Point cloud classification
    - Human-object interaction detection
    - …
  - The potential of robust self-supervised pre-training in building robust language foundation models