

# Unleashing the Power of Contrastive Self-Supervised Visual Models via Contrast-Regularized Fine-Tuning

Yifan Zhang<sup>1</sup>, Bryan Hooi<sup>1</sup>, Dapeng Hu<sup>1</sup>, Jian Liang<sup>2</sup>, Jiashi Feng<sup>3</sup>

<sup>1</sup>National University of Singapore    <sup>2</sup>Chinese Academy of Sciences    <sup>3</sup>SEA AI Lab

October 16, 2021

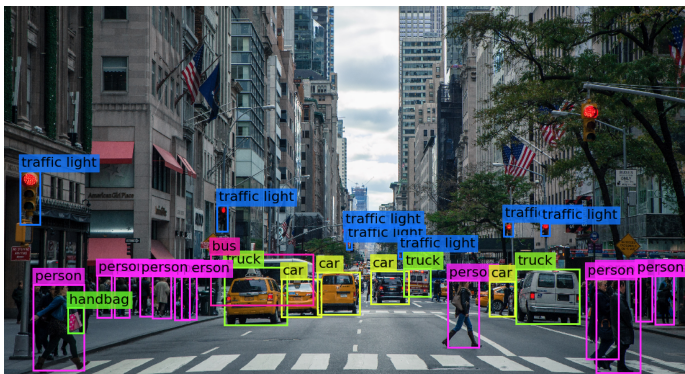
# Contents

- 1 Backgrounds
- 2 Preliminaries
- 3 Proposed Method
- 4 Experiments
- 5 Summary

- 1 Backgrounds
- 2 Preliminaries
- 3 Proposed Method
- 4 Experiments
- 5 Summary

# Deep Learning

Deep learning has achieved great success in many computer vision tasks, e.g., image classification, object detection, semantic segmentation



# Demands for abundant training data

- Deep neural networks (DNNs) have **massive model parameters**
- Hence, they require **a large number of annotated data** for training

Year	Model	#Layers	#Parameter	FLOPs	ImageNet Top-5 Error
2012	AlexNet	5+3	60M	725M	16.4%
2013	Clarifai	5+3	60M	–	11.7%
2014	MSRA	5+3	200M	–	8.06%
2014	VGG-19	16+3	143M	20G	7.32%
2014	GooLeNet	22	6.8M	2G	6.67%
2015	ResNet	152	60.2M	11G	5.79%

# Why self-supervision?

- However, supervised learning with large-scale labeled data is impractical
  - 1 Annotating large-scale data for each new task is **very expensive**
  - 2 Some areas are data-starved, e.g., in the medical area
- Some areas have a vast number of unlabeled image/video/graph data (e.g., Facebook: 1 billion images uploaded per day)
- **Self-supervised learning** is an alternatives to supervised learning

# Self-supervised learning

- Self-supervised learning is based on a pre-training and fine-tuning scheme
  - 1 Pre-train DNNs based on some specific **self-supervised pretext tasks** by using **unlabeled data as self-supervision**
  - 2 Transfer the pre-trained DNNs to solve downstream tasks
  
- Contrastive self-supervised learning has become the most popular one
  - 1 Pre-train DNNs based on an **instance discrimination contrastive task**<sup>1</sup>
  - 2 Transfer the contrastive pre-trained DNNs to solve downstream tasks

---

<sup>1</sup>Wu, Z., et al. Unsupervised feature learning via non-parametric instance discrimination. In CVPR, 2018

- Despite substantial studies on pre-training, few have explored fine-tuning
- The common practice is to directly fine-tune the pre-trained model with the cross-entropy loss on downstream tasks
- Fine-tuning with only cross-entropy may not be an optimal strategy
  - 1 Cross-entropy has limited ability to reduce intra-class feature variations
  - 2 Fine-tuning DNNs with cross-entropy may suffer from overfitting when the data number of downstream task is limited
- Our goal is to investigate how to better fine-tune contrastive self-supervised visual models on downstream tasks



# Contents

1 Backgrounds

2 Preliminaries

3 Proposed Method

4 Experiments

5 Summary

# Contrastive Learning

- Learning from **paired data** instead of single data
- Pull **positive (similar) pairs closer** and push **negative (dissimilar) pairs apart**
- Positive pairs: two transformations of an image or two images from a class

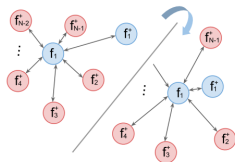
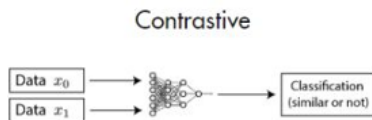


# Contrastive self-supervised learning

- Formulate the contrastive learning as a **multi-class classification problem**
- **Contrastive loss (InfoNCE loss<sup>2</sup>):**

$$L = -\log \frac{\exp(f^\top f^+)}{\exp(f^\top f^+) + \sum_{i=1}^{n-1} \exp(f^\top f_i^-)},$$

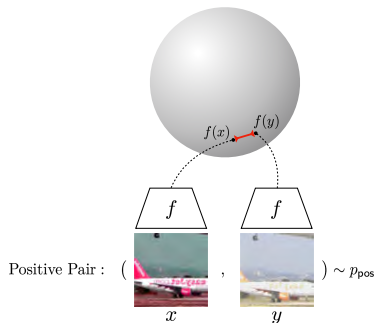
where  $f, f^+, f^-$  denote the features of the anchor image, the positive pair and negative pairs of the anchor



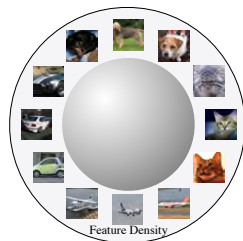
<sup>2</sup>Oord, A. V. D., et al. Representation learning with contrastive predictive coding, arxiv, 2018.

# What contrastive self-supervised learning learns?

- There are two findings in existing work<sup>3</sup>:
  - 1 Instance alignment: similar samples have similar features
  - 2 Instance uniformity: different data are pushed away



**Alignment:** Similar samples have similar features.  
(Figure inspired by [Tian et al. \(2019\)](#).)



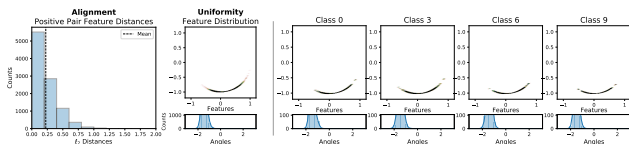
**Uniformity:** Preserve maximal information.

Figure 1: Illustration of alignment and uniformity of feature distributions on the output unit hypersphere. STL-10 ([Coates et al., 2011](#)) images are used for demonstration.

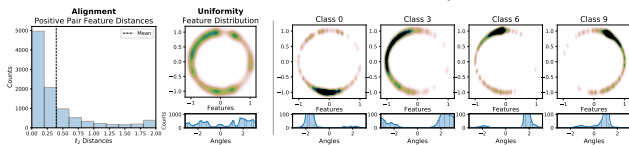
<sup>3</sup>Wang., et al. Understanding Contrastive Representation Learning through Alignment and Uniformity on the Hypersphere, In ICML, 2020.

# What contrastive self-supervised learning learns?

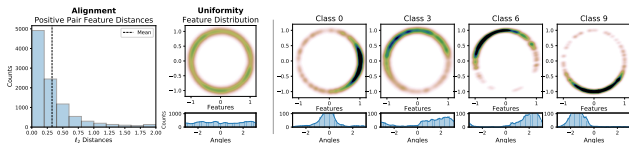
- The learned features are **uniformly distributed in the feature space**<sup>4</sup>



(a) **Random Initialization.** Linear classification validation accuracy: 12.71%.



(b) **Supervised Predictive Learning.** Linear classification validation accuracy: 57.19%.



(c) **Unsupervised Contrastive Learning.** Linear classification validation accuracy: 28.60%.

<sup>4</sup>Wang., et al. Understanding Contrastive Representation Learning through Alignment and Uniformity on the Hypersphere, In ICML, 2020.

# Why not fine-tune with only cross-entropy?

- Although cross-entropy tends to separate inter-class features, the resulted models still have **limited capability of reducing intra-class feature variations**

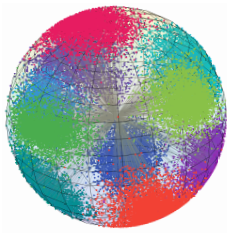


Figure: fine-tuning with only cross-entropy

- How to better fine-tune contrastive self-supervised visual models on downstream tasks?

# Contents

- 1 Backgrounds
- 2 Preliminaries
- 3 Proposed Method**
- 4 Experiments
- 5 Summary

- **Contrastive self-supervised learning:** optimizing unsupervised contrastive loss during pre-training yields models with [instance alignment property](#)
- Whether [applying contrastive learning to fine-tuning](#) would bring benefits?
- To answer this, we begin by analyzing supervised contrastive loss:

$$\mathcal{L}_{con} = -\frac{1}{n|P_i|} \sum_{i=1}^n \sum_{z_j \in P_i} \log \frac{e^{(z_i^\top z_j / \tau)}}{\sum_{z_k \in A_i} e^{(z_i^\top z_k / \tau)}},$$

where features  $z$  are  $\ell_2$ -normalized, and  $\tau$  is a temperature factor, while  $P_i$  and  $A_i$  denote the positive pair set and full pair set of the anchor  $z_i$

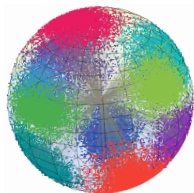


**Theorem 1** *Assuming the features are  $\ell_2$ -normalized and the classes are balanced with equal data number, minimizing the contrastive loss is equivalent to minimizing the class-conditional entropy  $\mathcal{H}(Z|Y)$  and maximizing the feature entropy  $\mathcal{H}(Z)$ :*

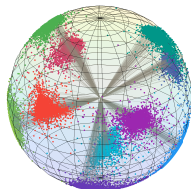
$$\mathcal{L}_{con} \propto \mathcal{H}(Z|Y) - \mathcal{H}(Z)$$

- This theorem shows that  $\mathcal{L}_{con}$  explicitly regularizes representation learning
  - 1 Minimizing  $\mathcal{L}_{con}$  will minimize  $\mathcal{H}(Z|Y)$ , which encourages to learn a low-entropy cluster for each class (i.e., high intra-class compactness)
  - 2 Minimizing  $\mathcal{L}_{con}$  will maximize  $\mathcal{H}(Z)$  and tends to learn a high-entropy feature space (i.e., large inter-class separation degree)

# Regularization effectiveness on representation learning



(a) Training with  $\mathcal{L}_{ce}$



(b) Training with  $\mathcal{L}_{ce} + \mathcal{L}_{con}$

- Based on the results, we confirm:
  - 1 Minimizing  $\mathcal{L}_{con}$  encourages to learn high intra-class compactness
  - 2 Minimizing  $\mathcal{L}_{con}$  tends to learn a feature space with large inter-class separation degree

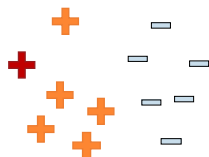
**Theorem 2** *Assuming the features are  $\ell_2$ -normalized and the classes are balanced, the contrastive loss is positive proportional to the infimum of conditional cross-entropy  $\mathcal{H}(Y; \hat{Y}|Z)$ , where the infimum is taken over classifiers:*

$$\mathcal{L}_{con} \propto \underbrace{\inf \mathcal{H}(Y; \hat{Y}|Z)}_{\text{Conditional CE}} - \mathcal{H}(Y)$$

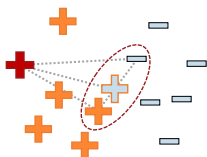
- This theorem shows that  $\mathcal{L}_{con}$  boosts the model optimization
  - The label  $Y$  is given by datasets, so its entropy  $\mathcal{H}(Y)$  is a constant
  - Minimizing  $\mathcal{L}_{con}$  minimizes **the infimum of conditional cross-entropy  $\mathcal{H}(Y; \hat{Y}|Z)$**
  - Pulling positives together and pushing negatives further apart makes the predicted label distribution closer to the ground truth distribution

- Based on the above theoretical analysis, we propose to **use contrastive learning to enhance the fine-tuning** of contrastive pre-trained models
- Instead of simply adding contrastive loss to the objective function, we further consider **two practical challenges**:
  - 1 How to **mine hard sample pairs** for better contrastive fine-tuning
  - 2 How to **improve the generalizability of the model**

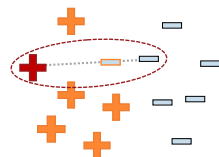
# Challenges I



(a) Sample features



(b) Hard positive mixup

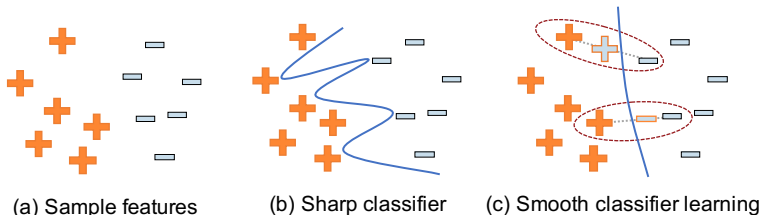


(c) Hard negative mixup

- The majority of sample are easy-to-contrast in training<sup>5</sup>, and may produce negligible contrastive loss gradients without performance contributions
- Our solution: to generate both hard positive and hard negative sample pairs based on feature mixup

<sup>5</sup>Hardwood, etc. Smart mining for deep metric learning. In ICCV, 2017.

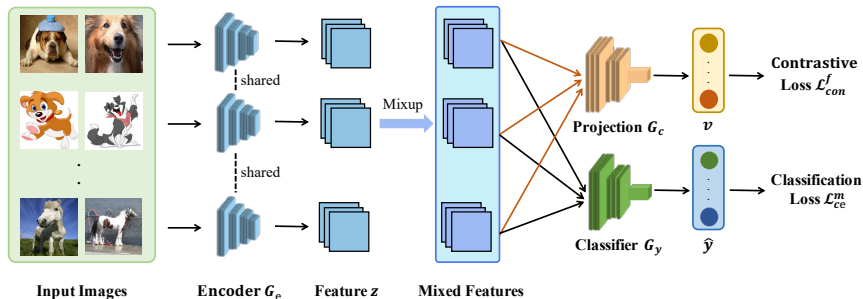
# Challenges II



- The decision boundary trained with cross-entropy is **often sharp and close to training data**<sup>6</sup>, especially when the data number is limited
- This may lead to **incorrect yet confident predictions** when evaluated on **slightly different test samples**
- Our solution: to **smooth classifiers** by using the mixed features

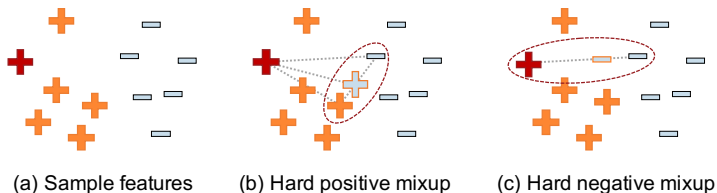
<sup>6</sup>Bengio, Y., et al. Manifold mixup: Better rep-representations by interpolating hidden states. In ICML, 2019.

# Overall scheme



- We propose a contrast-regularized tuning (Core-tuning) method
- Core-tuning consists of two main components:
  - 1 Hard sample pair mining
  - 2 Smooth classifier learning

# Hard positive sample pair generation



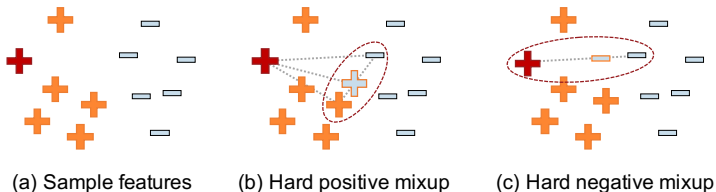
- Given a feature anchor  $z_i$ , we first find out its **hardest positive sample**  $z_i^{hp}$  and **hardest negative sample**  $z_i^{hn}$  based on cosine similarity
- We generate hard positive pair as a convex combination of  $z_i^{hp}$  and  $z_i^{hn}$ :

$$z_i^+ = \lambda z_i^{hp} + (1 - \lambda) z_i^{hn}$$

- $\lambda \sim \text{Beta}(\alpha, \alpha) \in [0, 1]$ , and  $\alpha$  is a parameter to decide the beta distribution



# Hard negative sample pair generation



- Given a feature anchor  $z_i$ , we randomly select a negative sample  $(z_i^n, y_i^n)$  to synthesize the hard negative pair:

$$z_i^- = (1 - \lambda)z_i + \lambda z_i^n; \quad y_i^- = (1 - \lambda)y_i + \lambda y_i^n$$

- We clip  $\lambda \geq 0.8$  to make the generated negative closer to negatives
- Hardness: the generated pairs are located between positives/negatives

# Hard pair reweighting

- We first use a two-layer MLP head  $G_c$  to obtain  $\ell_2$ -normalized contrastive features  $v_i = G_c(z_i) / \|G_c(z_i)\|_2$
- Since hard positives are more informative for contrastive learning, we propose to assign higher importance weights to them
- Inspired by the focal loss<sup>7</sup>, we find hard positive pairs generally lead to a low prediction probability:  $p_{ij} = \frac{\exp(v_i^\top v_j / \tau)}{\sum_{v_k \in A_i} \exp(v_i^\top v_k / \tau)}$
- Focal contrastive loss:

$$\mathcal{L}_{con}^f = -\frac{1}{n|P_i|} \sum_{i=1}^n \sum_{v_j \in P_i} (1-p_{ij}) \log \frac{e^{(v_i^\top v_j / \tau)}}{\sum_{v_k \in A_i} e^{(v_i^\top v_k / \tau)}},$$

where  $P_i$  and  $A_i$  denote the anchor's positive pair set and full pair set, both of which contain the generated hard pairs

<sup>7</sup>Lin, T-T, et al. Focal loss for dense object detection. In ICCV, 2017

# Smooth classifier learning

- The classifier trained with cross-entropy is often sharp and close to data, leading to limited generalization performance
- Inspired by that [mixup helps learn a smoother classifier](#)<sup>8</sup>, we further use the mixed data from the generated hard pair set  $\mathcal{B}$  for training

$$\mathcal{L}_{ce}^m = -\frac{1}{n} \sum_{i=1}^n y_i \log(\hat{y}_i) - \frac{1}{|\mathcal{B}|} \sum_{(z_j, y_j) \in \mathcal{B}} y_j \log(G_y(z_j))$$

---

<sup>8</sup>Bengio, Y., etc. Manifold mixup: Better representations by interpolating hidden states. In ICML, 2019.

# Contents

- 1 Backgrounds
- 2 Preliminaries
- 3 Proposed Method
- 4 Experiments**
- 5 Summary

# Experiment objectives

- 1 Verify Core-tuning on **image classification** and **semantic segmentation**
- 2 Evaluate how Core-tuning affects **model generalization to new domains** and **model robustness to adversarial examples**

# Ablation studies on Image Classification

- 1 Each component in Core-tuning improves the performance
- 2 Mixup (row 3) is expected to outperform mixup-hard w.r.t. classification, but Core-tuning still shows obvious improvement [via contrast regularizer](#)

$\mathcal{L}_{ce}$	$\mathcal{L}_{con}$	$\mathcal{L}_{con}^f$	mixup	mixup-hard	ImageNet-20	CIFAR10	CIFAR100	Caltech101	DTD	Aircraft	Cars	Pets	Flowers	Avg.
✓					88.28	94.70	80.27	91.87	71.68	86.87	88.61	89.05	98.49	87.76
✓	✓				89.29	95.33	81.49	92.84	72.73	87.44	89.37	89.71	98.65	88.54
✓			✓		90.67	95.43	81.03	92.68	73.31	88.37	89.06	91.37	98.74	88.96
✓	✓			✓	92.20	97.01	83.89	93.22	74.78	88.88	89.79	91.95	98.94	90.07
✓		✓		✓	<b>92.59</b>	<b>97.31</b>	<b>84.13</b>	<b>93.46</b>	<b>75.37</b>	<b>89.48</b>	<b>90.17</b>	<b>92.36</b>	<b>99.18</b>	<b>90.45</b>

# Results on Image Classification

- 1 MoCo-v2 pre-trained models perform worse than supervised ones
- 2 Fine-tuning methods for supervised pre-trained models cannot perform well
- 3 Using contrastive loss for fine-tuning may also performs unsatisfactory
- 4 Core-tuning achieves state-of-the-art fine-tuning performance

Method	ImageNet-20	CIFAR10	CIFAR100	Caltech101	DTD	Aircraft	Cars	Pets	Flowers	Avg.
SL-CE-tuning	91.01	94.23	83.40	93.39	74.40	87.03	89.77	92.17	98.78	89.35
CE-tuning	88.28	94.70	80.27	91.87	71.68	86.87	88.61	89.05	98.49	87.76
L2SP (Li et al., 2018)	88.49	95.14	81.43	91.98	72.18	86.55	89.00	89.43	98.66	88.10
M&M (Zhan et al., 2018)	88.53	95.02	80.58	92.91	72.43	87.45	88.90	89.60	98.57	88.22
DELTA (Li et al., 2019)	88.35	94.76	80.39	92.19	72.23	87.05	88.73	89.54	98.65	87.99
BSS (Chen et al., 2019)	88.34	94.84	80.40	91.95	72.22	87.18	88.50	89.50	98.57	87.94
RIFLE (Li et al., 2020)	89.06	94.71	80.36	91.94	72.45	87.60	89.72	90.05	98.70	88.29
SCL (Gunel et al., 2021)	89.29	95.33	81.49	92.84	72.73	87.44	89.37	89.71	98.65	88.54
Bi-tuning (Zhong et al., 2021)	89.06	95.12	81.42	92.83	73.53	87.39	89.41	89.90	98.57	88.58
Core-tuning (ours)	<b>92.59</b>	<b>97.31</b>	<b>84.13</b>	<b>93.46</b>	<b>75.37</b>	<b>89.48</b>	<b>90.17</b>	<b>92.36</b>	<b>99.18</b>	<b>90.45</b>

# Results on semantic segmentation

- Core-tuning contributes to fine-tuning performance of all pre-trained models
- Self-supervised models have already outperformed supervised ones
- This inspires us to explore unsupervised contrastive regularizers in the future

Pre-training	Fine-tuning	MPA	FWIoU	MIoU
Supervised	CE	87.10+/-0.20	89.12+/-0.17	76.52+/-0.34
InsDis	CE	83.64+/-0.12	88.23+/-0.08	74.14+/-0.21
	ours	<b>84.53+/-0.31</b>	<b>88.67+/-0.07</b>	<b>74.81+/-0.13</b>
PIRL	CE	83.16+/-0.26	88.22+/-0.24	73.99+/-0.42
	ours	<b>85.30+/-0.24</b>	<b>88.95+/-0.08</b>	<b>75.49+/-0.36</b>
MoCo-v2	CE	87.31+/-0.31	90.26+/-0.12	78.42+/-0.28
	ours	<b>88.76+/-0.34</b>	<b>90.75+/-0.04</b>	<b>79.62+/-0.12</b>
SimCLR-v2	CE	87.37+/-0.48	90.27+/-0.12	78.16+/-0.19
	ours	<b>87.95+/-0.20</b>	<b>90.71+/-0.13</b>	<b>79.15+/-0.33</b>
InfoMin	CE	87.17+/-0.20	89.84+/-0.09	77.84+/-0.24
	ours	<b>88.92+/-0.36</b>	<b>90.65+/-0.09</b>	<b>79.48+/-0.30</b>



# Effectiveness on cross-domain generalization

- Enforcing contrastive regularizer improves generalization performance
- Core-tuning further improves the generalization performance

Pre-training	Fine-tuning	PACS					VLCS				
		A	C	P	S	Avg.	C	L	V	S	Avg.
Supervised	CE	83.65	79.21	96.11	81.46	85.11	98.41	63.81	68.55	75.45	76.56
MoCo-v2	CE	78.71	76.92	90.87	75.67	80.54	94.96	66.87	68.96	64.98	73.94
	CE-Con	85.11	81.77	95.58	80.12	85.65	95.94	67.76	69.31	73.57	77.67
	ours	<b>87.31</b>	<b>84.06</b>	<b>97.53</b>	<b>83.43</b>	<b>88.08</b>	<b>98.50</b>	<b>68.19</b>	<b>73.15</b>	<b>81.53</b>	<b>80.34</b>

# Robustness to adversarial samples

- We generate adversarial samples via projected gradient descent (PGD)<sup>9</sup>
- Core-tuning is **beneficial to model robustness** on various image datasets
- We hope that Core-tuning can motivate people to **rethink the accuracy-robustness trade-off in adversarial training**

Method	$\ell_2$ -attack						$\ell_\infty$ -attack					
	$\epsilon=0.5$		$\epsilon=1.5$		$\epsilon=2.5$		$\epsilon=2/255$		$\epsilon=4/255$		$\epsilon=8/255$	
	Robust	Clean	Robust	Clean	Robust	Clean	Robust	Clean	Robust	Clean	Robust	Clean
CE	50.25	94.70	48.29	94.70	46.82	94.70	25.13	94.70	12.28	94.70	4.57	94.70
AT-CE	86.59	92.00	89.60	94.28	89.16	94.15	83.20	93.05	75.82	91.99	69.27	92.79
AT-CE-Con	90.74	94.71	90.29	94.80	89.70	94.27	85.07	94.56	79.75	93.79	70.70	93.38
AT-ours	<b>92.97</b>	<b>96.82</b>	<b>92.32</b>	<b>96.90</b>	<b>92.05</b>	<b>96.87</b>	<b>86.92</b>	<b>96.29</b>	<b>82.01</b>	<b>95.95</b>	<b>74.83</b>	<b>95.90</b>

<sup>9</sup>Madry, A., et al. Towards deep learning models resistant to adversarial attacks. In ICLR, 2018

# Contents

- 1 Backgrounds
- 2 Preliminaries
- 3 Proposed Method
- 4 Experiments
- 5 Summary**

- This work studied how to fine-tune contrastive self-supervised visual models
- There are three main contributions:
  - 1 We propose a simple yet effective contrast-regularized tuning method
  - 2 We analyze [the benefits of supervised contrastive loss to fine-tuning](#)
  - 3 We empirically show that Core-tuning effectively improve the fine-tuning performance of contrastive self-supervised models
- We call for more attention to the fine-tuning of contrastive self-supervised models [on understanding its underlying theories and better approaches](#)
- Source code is available: <https://github.com/Vanint/Core-tuning>.

# Thanks!