

# Improving Transferability of Representations via Augmentation-Aware Self-Supervision

**NeurIPS 2021**

Hankook Lee<sup>1</sup> Kibok Lee<sup>23</sup> Kimin Lee<sup>4</sup> Honglak Lee<sup>25</sup> Jinwoo Shin<sup>1</sup>

<sup>1</sup>Korea Advanced Institute of Science and Technology

<sup>2</sup>University of Michigan

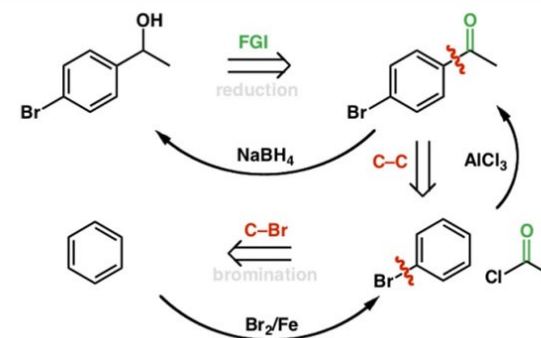
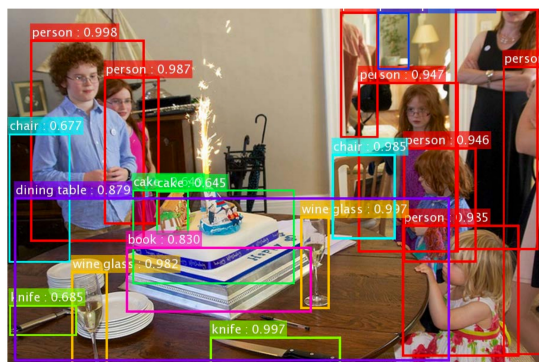
<sup>3</sup>Amazon Web Services

<sup>4</sup>University of California, Berkeley

<sup>5</sup>LG AI Research

# Unsupervised Representation Learning

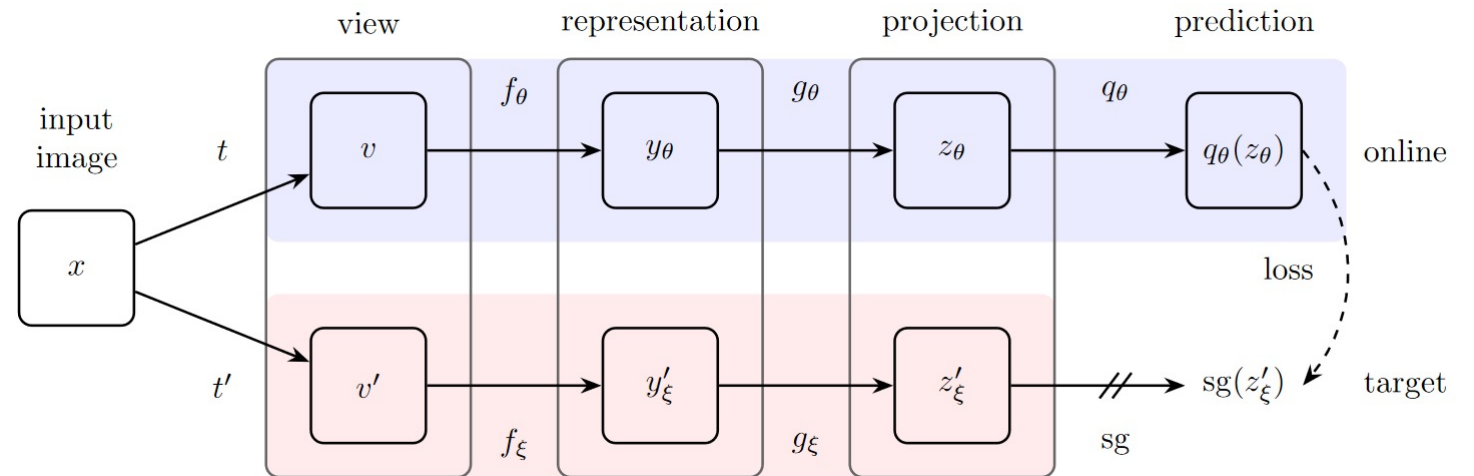
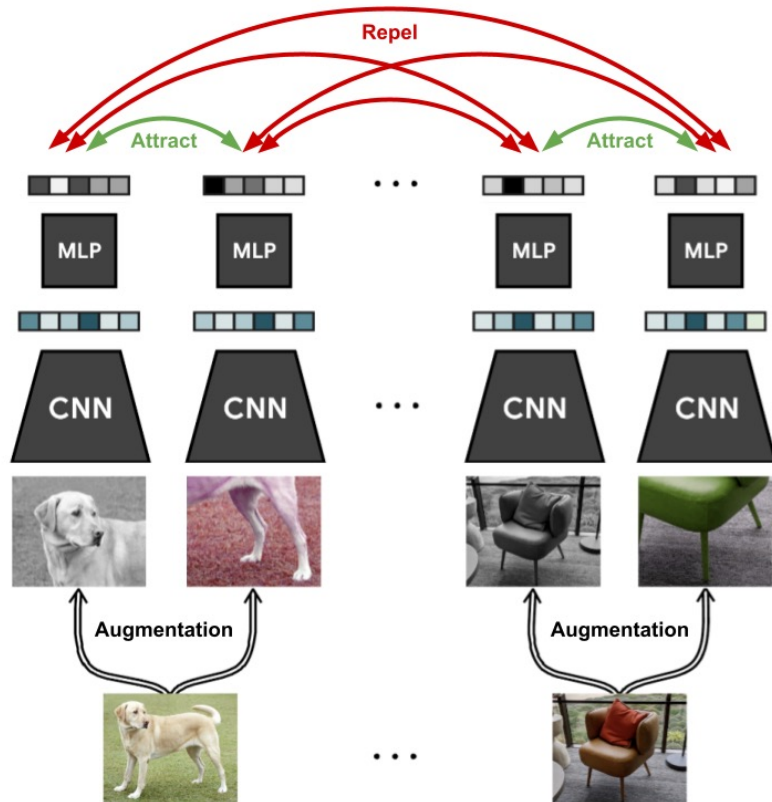
- DNNs have achieved a **remarkable success** in various applications
  - They often **require a massive amount of manually labeled data**
  - The **annotation cost is often expensive** because
    - It is **time-consuming**: e.g., annotating bounding boxes
    - It requires **expert knowledge**: e.g., medical diagnosis and retrosynthesis



- Hence, **collecting unlabeled samples is easier** than doing labeled samples
- **Question:** How to utilize the unlabeled samples for representation learning?

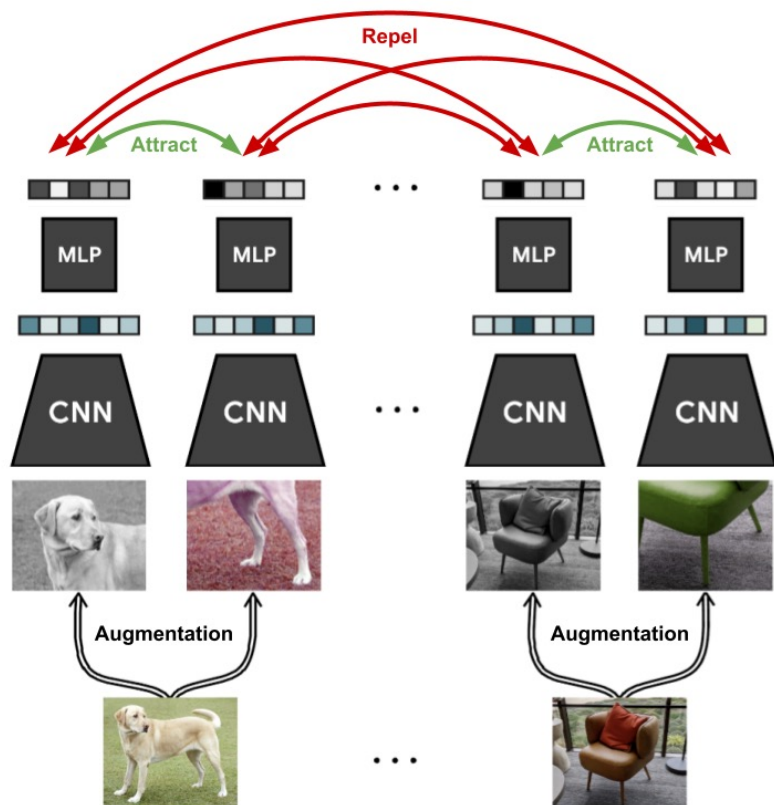
# Recent Advances in Self-supervised Learning

- State-of-the-art self-supervised learning methods have shown promising results
  - The SSL methods remarkably reduce the gap to supervised learning
  - They commonly **learn augmentation-invariant representations**



# Recent Advances in Self-supervised Learning

- State-of-the-art self-supervised learning methods have shown promising results
  - The SSL methods remarkably reduce the gap to supervised learning
  - They commonly **learn augmentation-invariant representations**



Contrastive methods (e.g., SimCLR [1] and MoCo [2])

$$\mathcal{L} = -\log \frac{\exp(\text{sim}(\mathbf{z}_1, \mathbf{z}_2)/\tau)}{\exp(\text{sim}(\mathbf{z}_1, \mathbf{z}_2)/\tau) + \sum_{\mathbf{z}'} \exp(\text{sim}(\mathbf{z}_1, \mathbf{z}')/\tau)}$$

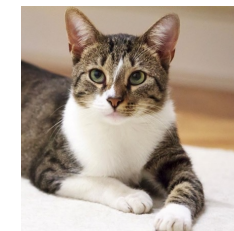
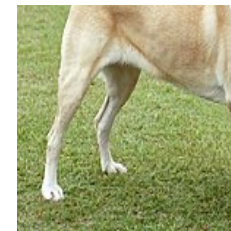
**Maximize**  $\text{sim}(\mathbf{z}_1, \mathbf{z}_2)$

**Minimize**  $\text{sim}(\mathbf{z}_1, \mathbf{z}')$

$$\mathbf{z}_1 = f(\mathbf{x}_1)$$

$$\mathbf{z}_2 = f(\mathbf{x}_2)$$

$$\mathbf{z}' = f(\mathbf{x}')$$



$\mathbf{x}_1$

$\mathbf{x}_2$

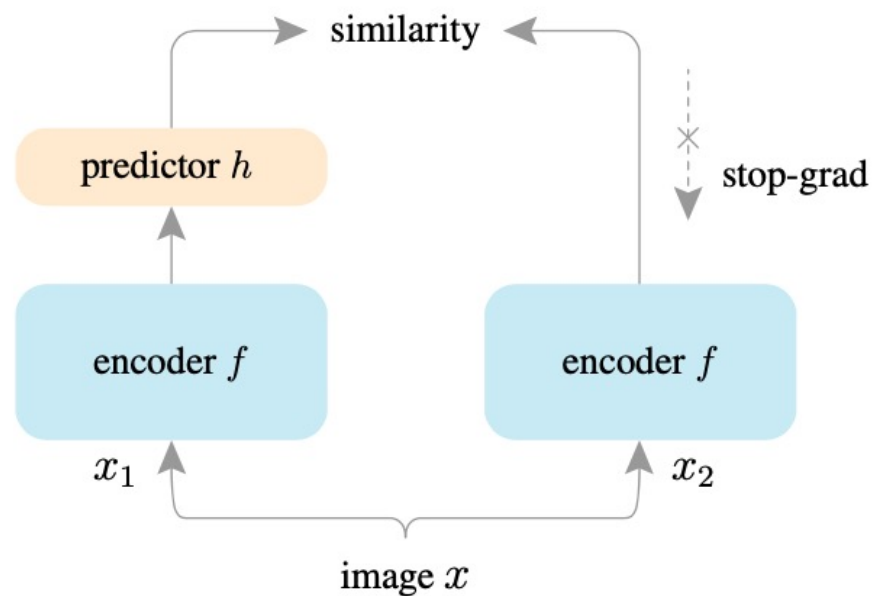
$\mathbf{x}'$

[1] Chen et al., A Simple Framework for Contrastive Learning of Visual Representations, ICML 2020

[2] He et al., Momentum Contrast for Unsupervised Visual Representation Learning, CVPR 2020

# Recent Advances in Self-supervised Learning

- State-of-the-art self-supervised learning methods have shown promising results
  - The SSL methods remarkably reduce the gap to supervised learning
  - They commonly **learn augmentation-invariant representations**

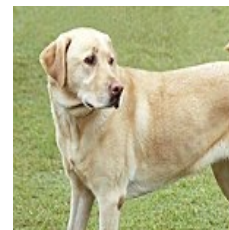


Non-contrastive methods (e.g., BYOL [3] and SimSiam [4])

$$\mathcal{L} = \|h(\mathbf{z}_1) - \text{stop\_grad}(\mathbf{z}_2)\|_2^2$$

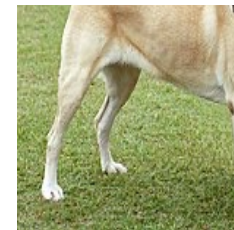
$$\begin{aligned}\nabla \mathbb{E}[\mathcal{L}] &= \nabla \mathbb{E}[\|h^*(\mathbf{z}_1) - \mathbf{z}_2\|_2^2] \\ &= \nabla \mathbb{E}[\|\mathbb{E}[\mathbf{z}_2|\mathbf{z}_1] - \mathbf{z}_2\|_2^2] \\ &= \nabla \mathbb{E}\left[\sum_i \text{Var}(\mathbf{z}_2^{(i)}|\mathbf{z}_1)\right]\end{aligned}$$

$$\mathbf{z}_1 = f(\mathbf{x}_1)$$



$\mathbf{x}_1$

$$\mathbf{z}_2 = f(\mathbf{x}_2)$$



$\mathbf{x}_2$

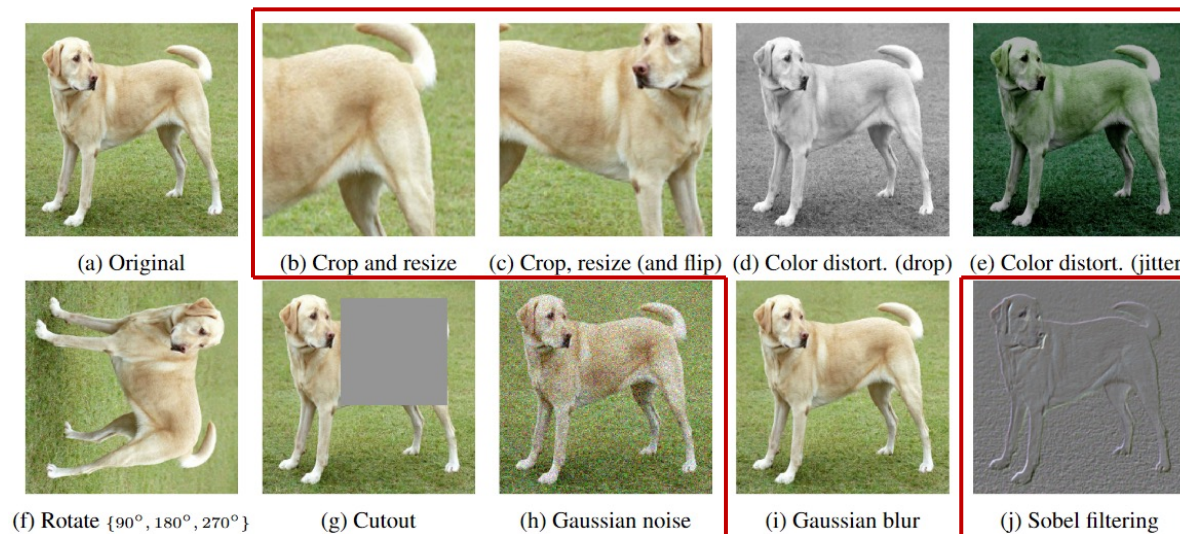
[3] Grill et al., Bootstrap Your Own Latent: A New Approach to Self-supervised Learning, 2020

[4] Chen & He, Exploring Simple Siamese Representation Learning, 2020

# Recent Advances in Self-supervised Learning

- State-of-the-art self-supervised learning methods have shown promising results
  - The SSL methods remarkably reduce the gap to supervised learning
  - They commonly **learn augmentation-invariant representations**
  - Augmentations:
    - **Geometric augmentations:** Cropping, Resizing, Flipping
    - **Color augmentations:** Color Jittering, Color Dropping, Gaussian Blurring

## Commonly used augmentations for invariant representation learning



# Motivation

- Total = (a) augmentation-invariant information + (b) augmentation-aware information

Original image



**Yellow** flower

Augmented image



**Flower** of unknown color

Augmentation  
Color Dropping

- (a) augmentation-invariant information = **Flower**
  - (b) augmentation-aware information = **Yellow**
- **Q) Is augmentation-aware information not or less important?**

# Motivation

- **Q) Is augmentation-aware information not or less important?**
- Learning augmentation-invariance may hurt performance in certain downstream tasks
  - Learning invariance to color augmentations (e.g., color dropping) forces the representations of color-modified and original images to be same as much as possible

$$f \left( \text{img}_1 \right) \approx f \left( \text{img}_2 \right)$$
$$f \left( \text{img}_3 \right) \approx f \left( \text{img}_4 \right)$$

Which flower is yellow? 🙄

- It degrades the representation qualities for color-sensitive downstream tasks such as flower classification



# Motivation

- **Q) Is augmentation-aware information not or less important?**
- Learning augmentation-invariance may hurt performance in certain downstream tasks
  - Learning invariance to color augmentations (e.g., color dropping) forces the representations of color-modified and original images to be same as much as possible

$$f\left(\text{img}_1\right) \approx f\left(\text{img}_2\right) \quad \text{v.s.} \quad f\left(\text{img}_3\right) \approx f\left(\text{img}_4\right)$$


- It degrades the representation qualities for color-sensitive downstream tasks such as flower classification
- **Q) How to learn more generalizable and transferable representations?**
- **Our goal** is to prevent information loss from learning augmentation-invariance, i.e., to **learn both augmentation-invariant and augmentation-aware representations**

# AugSelf: Auxiliary Augmentation-aware Self-supervision

- **Notations**

- Original image  $\mathbf{x}$
- Augmentation function  $t_\omega$  where  $\omega \sim \Omega$  is augmentation-specific parameter
- Augmented view  $\mathbf{v} = t_\omega(\mathbf{x})$
- Examples:



Original image



Random cropping

$$\begin{aligned}\omega^{\text{crop}} &= (y_{\text{center}}, x_{\text{center}}, H, W) \\ &= (0.4, 0.3, 0.6, 0.4)\end{aligned}$$



Horizontal flipping

$$\begin{aligned}\omega^{\text{flip}} &= \mathbb{1}[\mathbf{v} \text{ is flipped}] \\ &= 1\end{aligned}$$



Color jittering

$$\begin{aligned}\omega^{\text{color}} &= (\lambda_{\text{bright}}, \lambda_{\text{contrast}}, \lambda_{\text{sat}}, \lambda_{\text{hue}}) \\ &= (0.3, 1.0, 0.8, 1.0)\end{aligned}$$

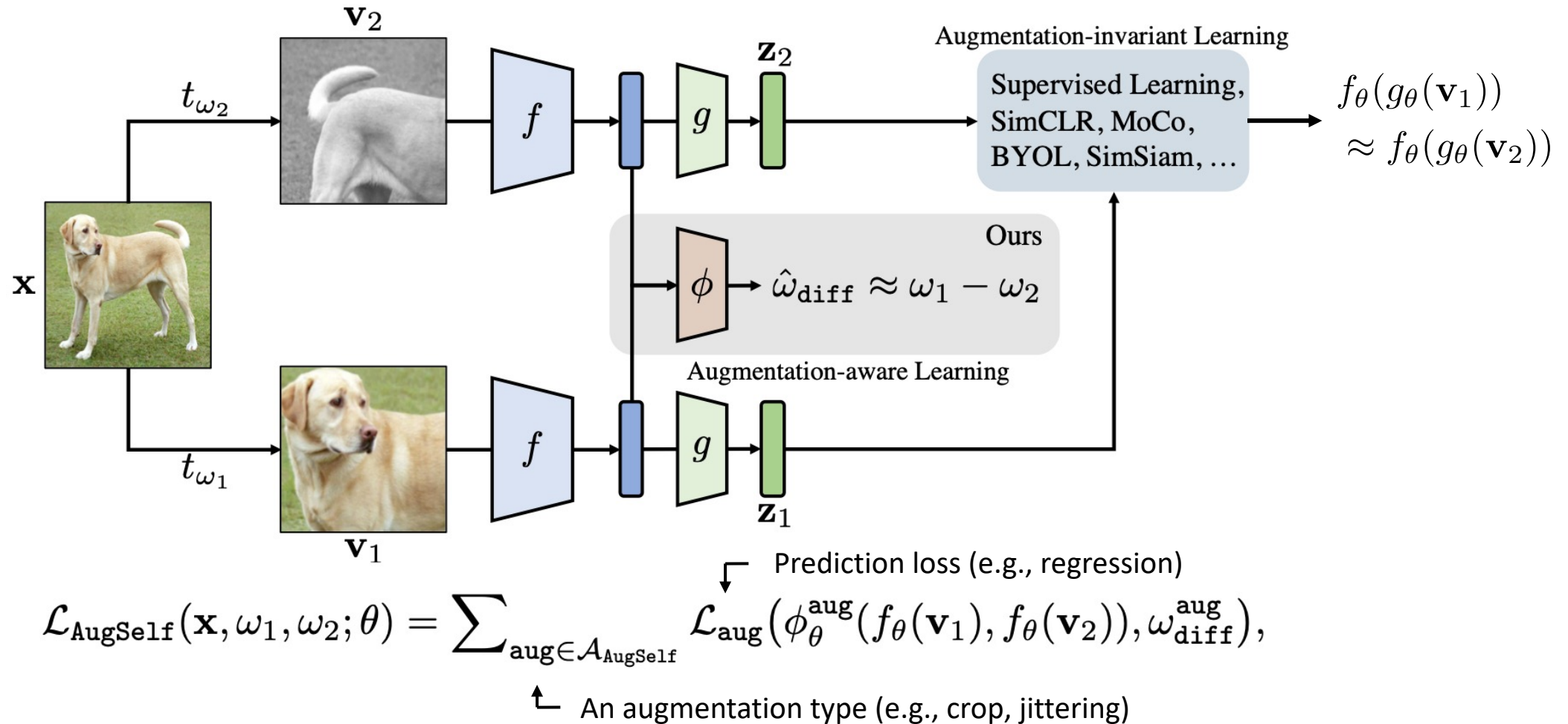


Gaussian blurring

$$\begin{aligned}\omega^{\text{blur}} &= \text{std. dev. of Gaussian kernel} \\ &= 1.0\end{aligned}$$

- Augmentation parameters  $\omega$  explain how the image is modified
- Main idea is to **predict the augmentation parameters from augmented views**

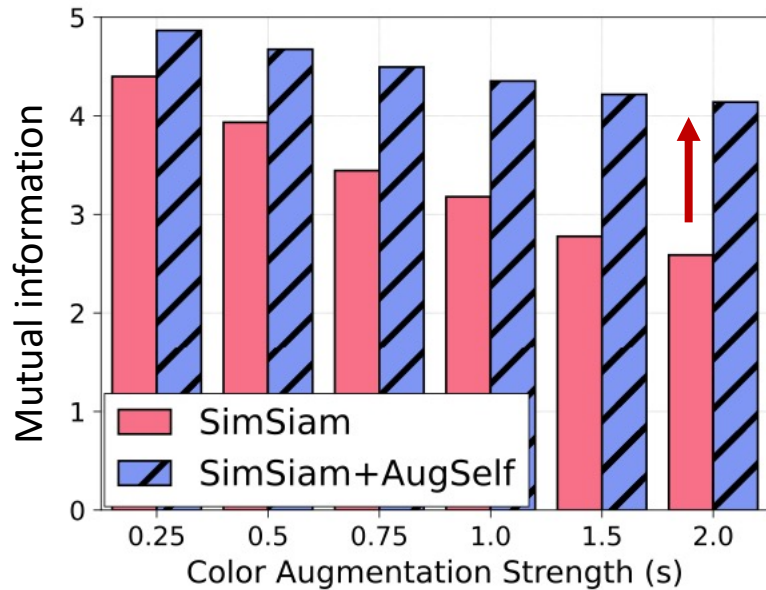
# AugSelf: Auxiliary Augmentation-aware Self-supervision



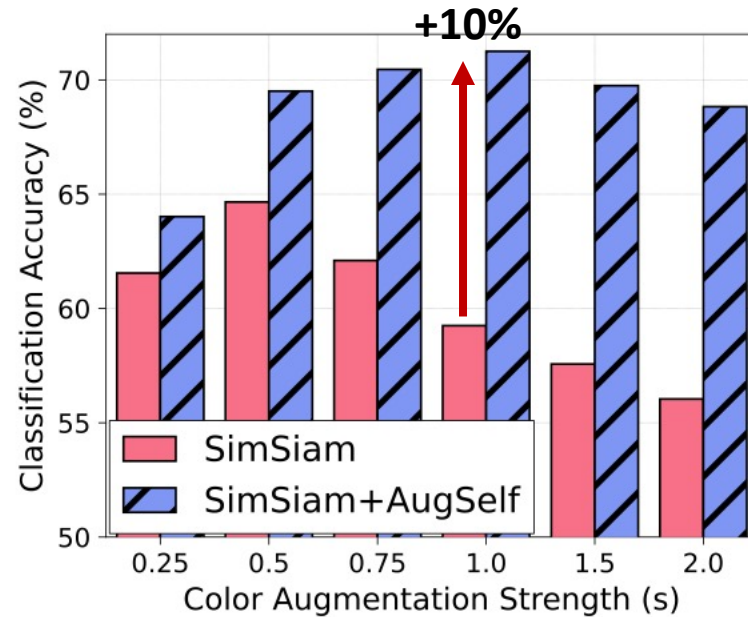
- **AugSelf** learns to predict difference between augmentation parameters of two views
  - This prediction task encourages  $f(x)$  to learn augmentation-aware information
  - This design allows to incorporate AugSelf into existing frameworks **without additional training costs**

# Analysis: Mutual Information

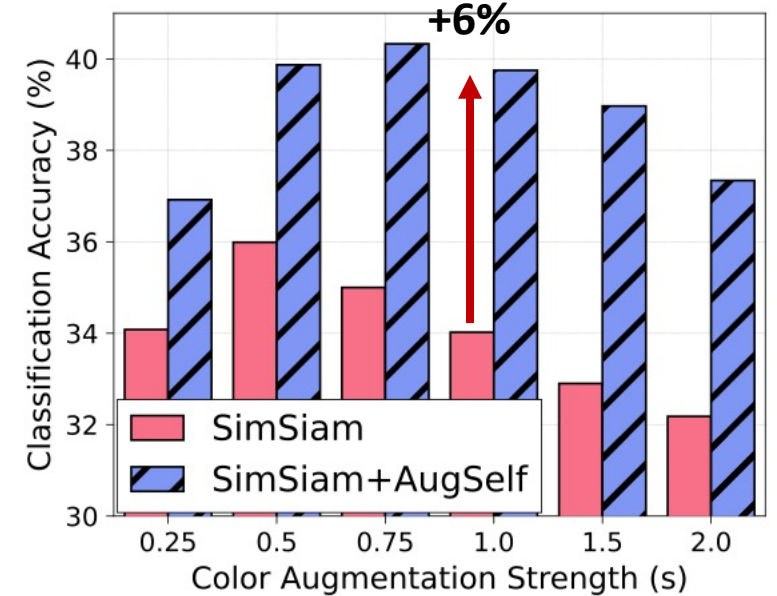
- AugSelf preserves the augmentation-aware information
  - $I_{NCE}(C; Z)$  = the mutual information between color histogram (i.e.,  $C$ ) and representation (i.e.,  $Z = f(x)$ )
  - AugSelf significantly improves the linear evaluation accuracy in the color-sensitive downstream tasks



(a) Mutual information



(b) STL10→Flowers



(c) STL10→Food

# Ablation Study: All Information Is Useful

- Both color/geometric information is useful in various downstream tasks
  - Learn **color** information by predicting **Color Jittering** parameters
  - Learn **geometric** information by predicting **Random Cropping** parameters

	$\mathcal{A}_{\text{AugSelf}}$	STL10	CIFAR10	CIFAR100	Food	MIT67	Pets	Flowers
Aug. parameters we predict	$\emptyset$	85.19	82.35	54.90	33.99	39.15	44.90	59.19
	{crop}	<b>85.98</b>	82.82	55.78	35.68	43.21	47.10	62.05
	{color}	85.55	<b>82.90</b>	58.11	40.32	43.56	47.85	71.08
	{crop, color}	85.70	82.76	<b>58.65</b>	<b>41.58</b>	<b>45.67</b>	<b>48.42</b>	<b>72.18</b>

- The improvement depends on the characteristic of the downstream tasks
- Learning all information achieves best performance in most downstream tasks

# Experimental Results: Fine-grained Classification Tasks

- AugSelf consistently improves Supervised Learning, SimSiam, MoCo in various settings
  - 11 fine-grained classification benchmarks

Method	CIFAR10	CIFAR100	Food	MIT67	Pets	Flowers	Caltech101	Cars	Aircraft	DTD	SUN397
<i>ImageNet100-pretrained ResNet-50</i>											
SimSiam	86.89	66.33	61.48	65.75	74.69	88.06	84.13	<b>48.20</b>	48.63	65.11	50.60
+ AugSelf	<b>88.80</b>	<b>70.27</b>	<b>65.63</b>	<b>67.76</b>	<b>76.34</b>	<b>90.70</b>	<b>85.30</b>	47.52	<b>49.76</b>	<b>67.29</b>	<b>52.28</b>
MoCo v2	84.60	61.60	59.37	61.64	70.08	82.43	77.25	33.86	<b>41.21</b>	64.47	46.50
+ AugSelf	<b>85.26</b>	<b>63.90</b>	<b>60.78</b>	<b>63.36</b>	<b>73.46</b>	<b>85.70</b>	<b>78.93</b>	<b>37.35</b>	39.47	<b>66.22</b>	<b>48.52</b>
Supervised	<b>86.16</b>	62.70	53.89	52.91	73.50	76.09	<b>77.53</b>	30.61	36.78	61.91	40.59
+ AugSelf	86.06	<b>63.77</b>	<b>55.84</b>	<b>54.63</b>	<b>74.81</b>	<b>78.22</b>	77.47	<b>31.26</b>	<b>38.02</b>	<b>62.07</b>	<b>41.49</b>
<i>STL10-pretrained ResNet-18</i>											
SimSiam	82.35	54.90	33.99	39.15	44.90	59.19	66.33	16.85	26.06	42.57	29.05
+ AugSelf	<b>82.76</b>	<b>58.65</b>	<b>41.58</b>	<b>45.67</b>	<b>48.42</b>	<b>72.18</b>	<b>72.75</b>	<b>21.17</b>	<b>33.17</b>	<b>47.02</b>	<b>34.14</b>
MoCo v2	81.18	53.75	33.69	39.01	42.34	61.01	64.15	16.09	26.63	41.20	28.50
+ AugSelf	<b>82.45</b>	<b>57.17</b>	<b>36.91</b>	<b>41.67</b>	<b>43.80</b>	<b>66.96</b>	<b>66.02</b>	<b>17.53</b>	<b>28.02</b>	<b>45.21</b>	<b>30.93</b>

# Experimental Results: Fine-grained Classification Tasks

- AugSelf consistently improves Supervised Learning, SimSiam, MoCo in various settings
  - 11 fine-grained classification benchmarks

Method	CIFAR10	CIFAR100	Food	MIT67	Pets	Flowers	Caltech101	Cars	Aircraft	DTD	SUN397
<i>ImageNet100-pretrained ResNet-50</i>											
SimSiam											20.60
+ AugSelf											22.28
MoCo											26.50
+ AugSelf											28.52
Supervised											20.59
+ AugSelf											21.49
SimSiam											29.05
+ AugSelf											31.14
MoCo											28.50
+ AugSelf	<b>82.45</b>	<b>57.17</b>	<b>36.91</b>	<b>41.67</b>	<b>43.80</b>	<b>66.96</b>	<b>66.02</b>	<b>17.53</b>	<b>28.02</b>	<b>45.21</b>	<b>30.93</b>

# Experimental Results: Few-shot Classification Tasks

- AugSelf consistently improves Supervised Learning, SimSiam, MoCo in various settings
  - 11 fine-grained classification benchmarks
  - 3 few-shot classification benchmarks

Method	FC100		CUB200		Plant Disease		→ 5-way 5-shot task
	(5, 1)	(5, 5)	(5, 1)	(5, 5)	(5, 1)	(5, 5)	
<i>ImageNet100-pretrained ResNet-50</i>							
SimSiam	36.19 $\pm$ 0.36	50.36 $\pm$ 0.38	45.56 $\pm$ 0.47	62.48 $\pm$ 0.48	75.72 $\pm$ 0.46	89.94 $\pm$ 0.31	
+ AugSelf (ours)	<b>39.37<math>\pm</math>0.40</b>	<b>55.27<math>\pm</math>0.38</b>	<b>48.08<math>\pm</math>0.47</b>	<b>66.27<math>\pm</math>0.46</b>	<b>77.93<math>\pm</math>0.46</b>	<b>91.52<math>\pm</math>0.29</b>	
MoCo v2	31.67 $\pm$ 0.33	43.88 $\pm$ 0.38	41.67 $\pm$ 0.47	56.92 $\pm$ 0.47	65.73 $\pm$ 0.49	84.98 $\pm$ 0.36	
+ AugSelf (ours)	<b>35.02<math>\pm</math>0.36</b>	<b>48.77<math>\pm</math>0.39</b>	<b>44.17<math>\pm</math>0.48</b>	<b>57.35<math>\pm</math>0.48</b>	<b>71.80<math>\pm</math>0.47</b>	<b>87.81<math>\pm</math>0.33</b>	
Supervised	33.15 $\pm$ 0.33	46.59 $\pm$ 0.37	46.57 $\pm$ 0.48	63.69 $\pm$ 0.46	68.95 $\pm$ 0.47	88.77 $\pm$ 0.30	
+ AugSelf (ours)	<b>34.70<math>\pm</math>0.35</b>	<b>48.89<math>\pm</math>0.38</b>	<b>47.58<math>\pm</math>0.48</b>	<b>65.31<math>\pm</math>0.45</b>	<b>70.82<math>\pm</math>0.46</b>	<b>89.77<math>\pm</math>0.29</b>	
<i>STL10-pretrained ResNet-18</i>							
SimSiam	36.72 $\pm$ 0.35	51.49 $\pm$ 0.36	37.97 $\pm$ 0.43	50.61 $\pm$ 0.45	58.13 $\pm$ 0.50	75.98 $\pm$ 0.40	
+ AugSelf (ours)	<b>40.68<math>\pm</math>0.39</b>	<b>56.26<math>\pm</math>0.38</b>	<b>41.60<math>\pm</math>0.42</b>	<b>56.33<math>\pm</math>0.44</b>	<b>62.85<math>\pm</math>0.49</b>	<b>81.14<math>\pm</math>0.37</b>	
MoCo v2	35.69 $\pm$ 0.34	49.26 $\pm$ 0.36	37.62 $\pm$ 0.42	50.71 $\pm$ 0.44	57.87 $\pm$ 0.48	75.98 $\pm$ 0.40	
+ AugSelf (ours)	<b>39.66<math>\pm</math>0.39</b>	<b>55.58<math>\pm</math>0.39</b>	<b>38.33<math>\pm</math>0.41</b>	<b>51.93<math>\pm</math>0.44</b>	<b>60.78<math>\pm</math>0.50</b>	<b>78.76<math>\pm</math>0.38</b>	



# Experimental Results: Object Localization

- AugSelf consistently improves Supervised Learning, SimSiam, MoCo in various settings
  - 11 fine-grained classification benchmarks
  - 3 few-shot classification benchmarks
  - Object localization on CUB200 benchmark

Method	Error
SimSiam	0.00462
+ AugSelf	<b>0.00335</b>
MoCo	0.00487
+ AugSelf	<b>0.00429</b>
Supervised	0.00520
+ AugSelf	<b>0.00473</b>

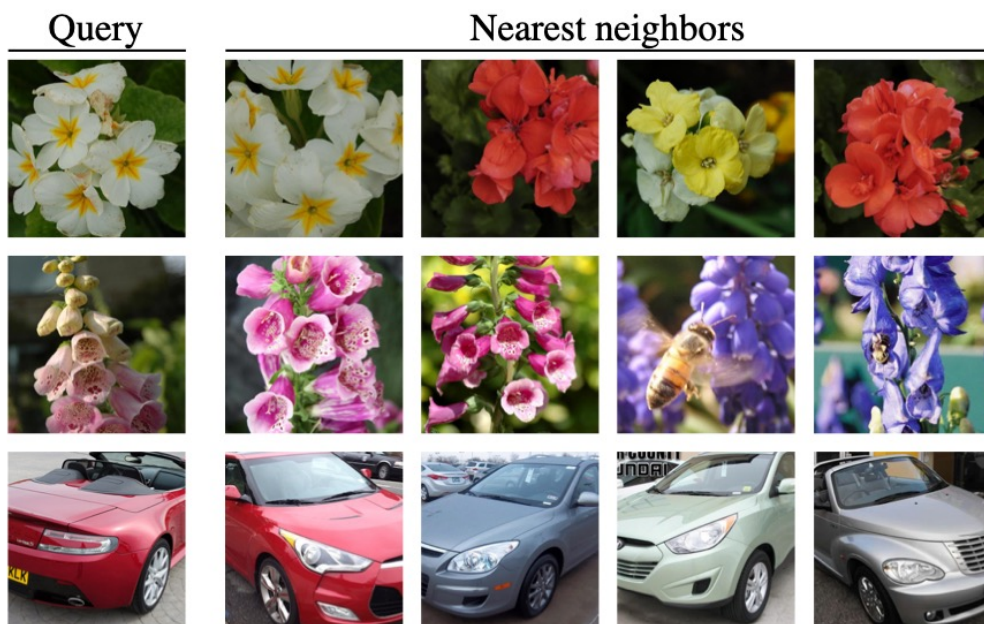
Table 4:  $\ell_2$  errors of bounding box predictions on CUB200.



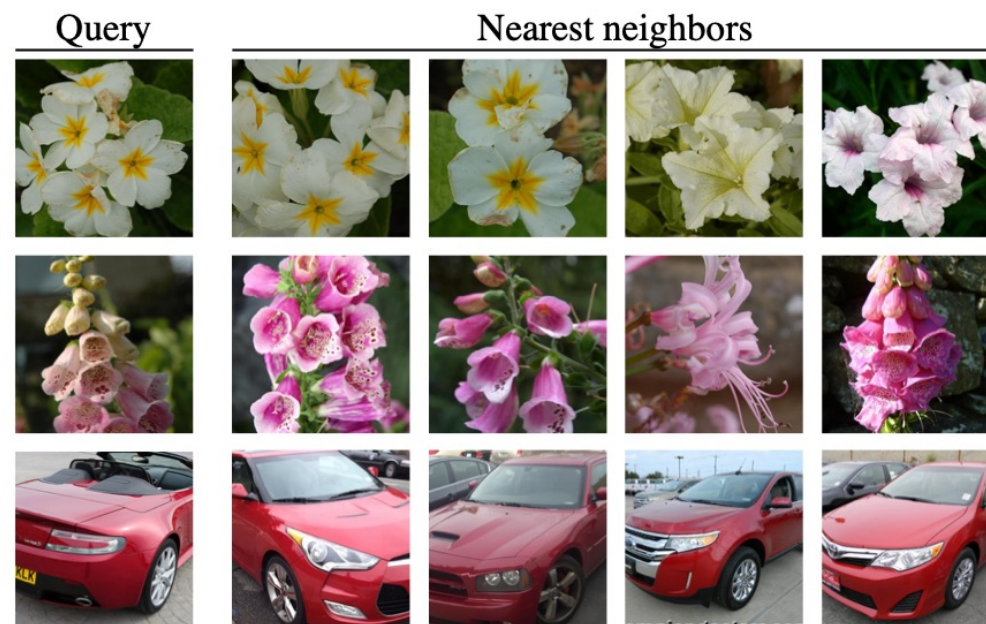
Figure 4: Examples of bounding box predictions on CUB200. Blue and red boxes are ground-truth and model prediction, respectively.

# Experimental Results: Retrieval

- AugSelf consistently improves Supervised Learning, SimSiam, MoCo in various settings
  - 11 fine-grained classification benchmarks
  - 3 few-shot classification benchmarks
  - Object localization on CUB200 benchmark
- Quantitative analysis (based on retrieval)



(a) SimSiam



(b) SimSiam + AugSelf

# Conclusion

- We propose AugSelf for learning more transferable and generalizable representations
  - AugSelf encourages to preserve augmentation-aware information by learning the difference of augmentation parameters between two randomly augmented samples
  - AugSelf can easily be incorporated into recent state-of-the-art self-supervised learning methods with a negligible additional training cost
  - Extensive experiments demonstrate that AugSelf consistently improves the transferability of representations learned by supervised and unsupervised methods in various transfer learning scenarios

Thank you for your attention!