



# Augmentation-aware Self-supervised Learning with Conditioned Projector

**Marcin Przewięźlikowski**, Mateusz Pyla, Bartosz Zieliński,  
Bartłomiej Twardowski, Jacek Tabor, Marek Śmieja

ML in PL 2024

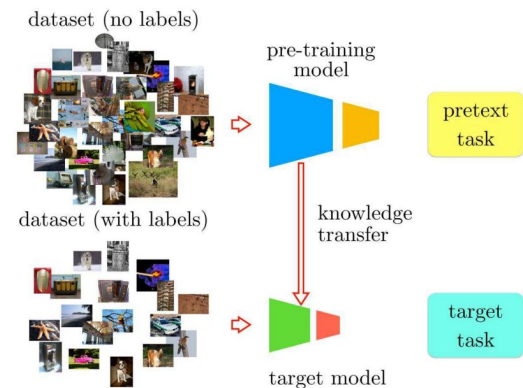
# Three pillars of AI success



# Problems with data



# Self-supervised learning

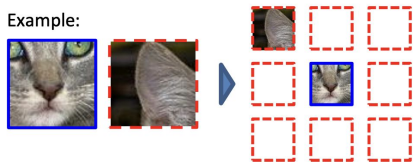


*“the dark matter of intelligence”* - Yann LeCun

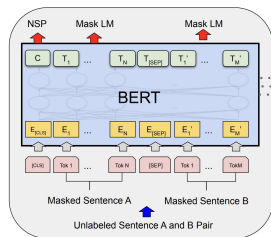
# What can serve as a pretext task?

And how they did it before 2020s

Example:



Context prediction



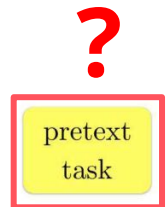
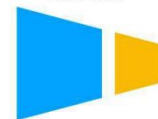
Masked modeling  
(i.e. BERT)

dataset (no labels)



dataset (with labels)

pre-training  
model



90° rotation

270° rotation

Rotation prediction

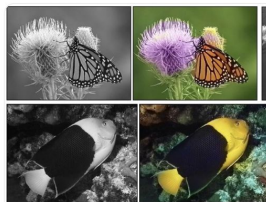


Image colorization

# Modern pretext tasks

Contrastive Siamese Joint-Embedding models

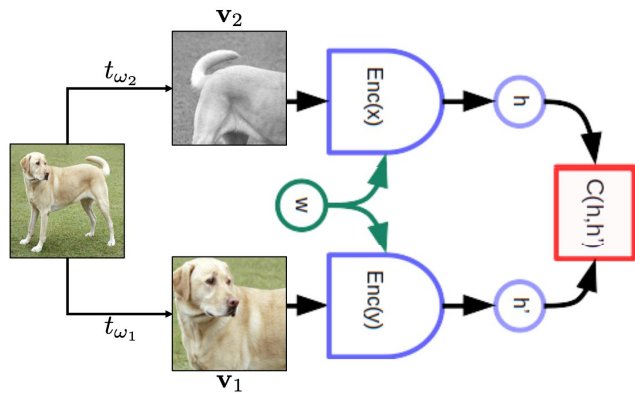


# Modern pretext tasks

## Contrastive Siamese Joint-Embedding models

### Intuition:

- augment an image in two different ways
- obtain network representations of two augmented images
- optimize the (pairwise) similarity of image representations **and** their diversity



$$\mathcal{L}_{SSL} = \sum_{(\mathbf{x}, \mathbf{x}') \in \text{PositivePairs}} \text{Distance}(f_{\theta}(\mathbf{x}), f_{\theta}(\mathbf{x}')) - \text{Diversity}(\{f_{\theta}(\mathbf{x}), \mathbf{x} \in \mathbb{X}\}),$$

# What are the problems with Joint-Embedding SSL?

Joint-embedding SSL methods are inherently bound to augmentations

Augmentations need to be carefully selected for pretraining datasets

- solved for ImageNet, but what about other datasets?

Invariance to augmentations can be detrimental for downstream tasks

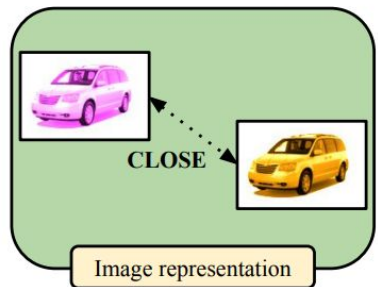
- invariance to color shifts may not transfer well to flower classification



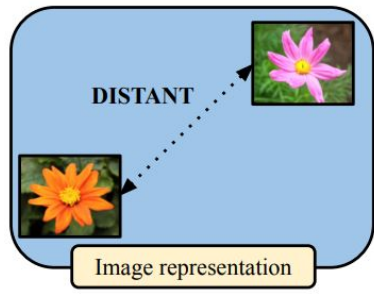
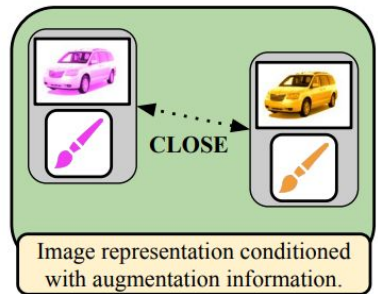
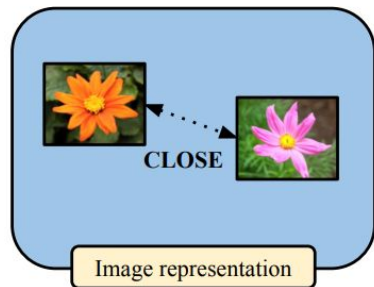


# Conditional Augmentation-aware Self-Supervised Learning (CASSLE)

Projector representation used  
in the contrastive objective



Feature extractor representation  
used in downstream tasks



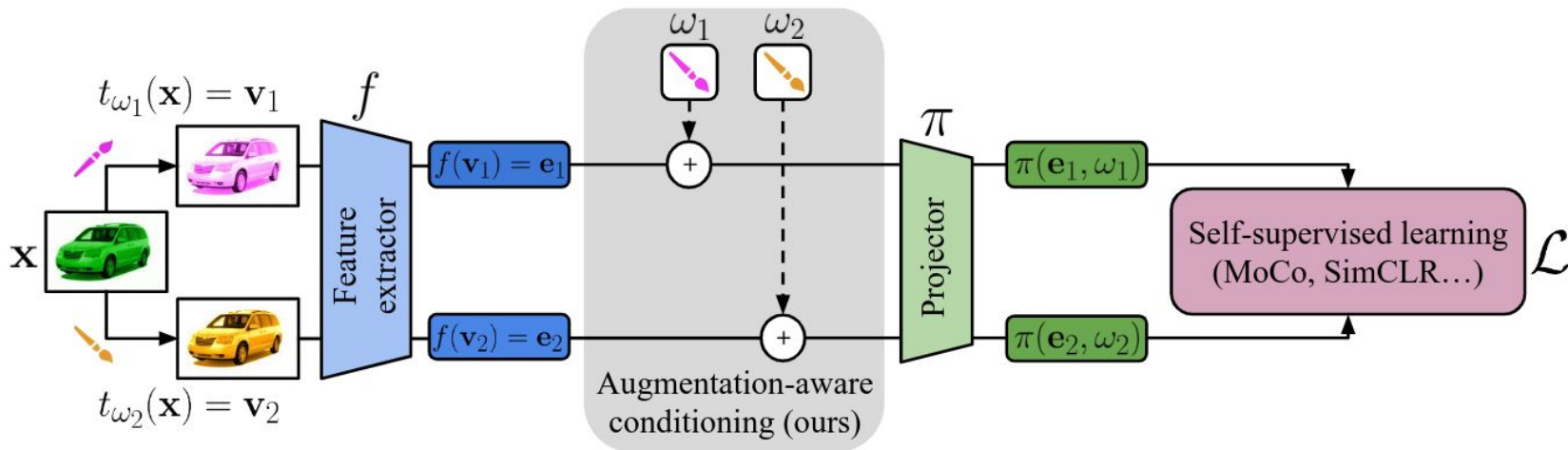
**Previously:**

*Join two image embeddings together*

**Now:**

*Join two image embeddings together **on condition of knowing how they were augmented***

# Architecture of CASSLE



- Image and augmentation information are joined together before the projector
- Feature extractor (e.g. a ResNet) remains unaware of augmentation information
- In order for projector to act upon the knowledge of augmentations, feature extractor must learn to preserve information about features modified by them

**CASSLE is applicable to all J-E architectures, regardless of their loss function.**

# What defines the augmentations?



Original image



Random cropping

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



Color jittering

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



Horizontal flipping

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



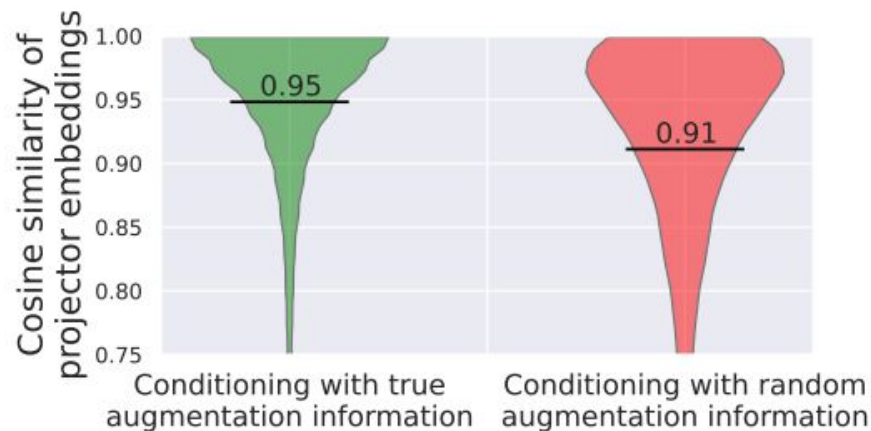
Gaussian blurring

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

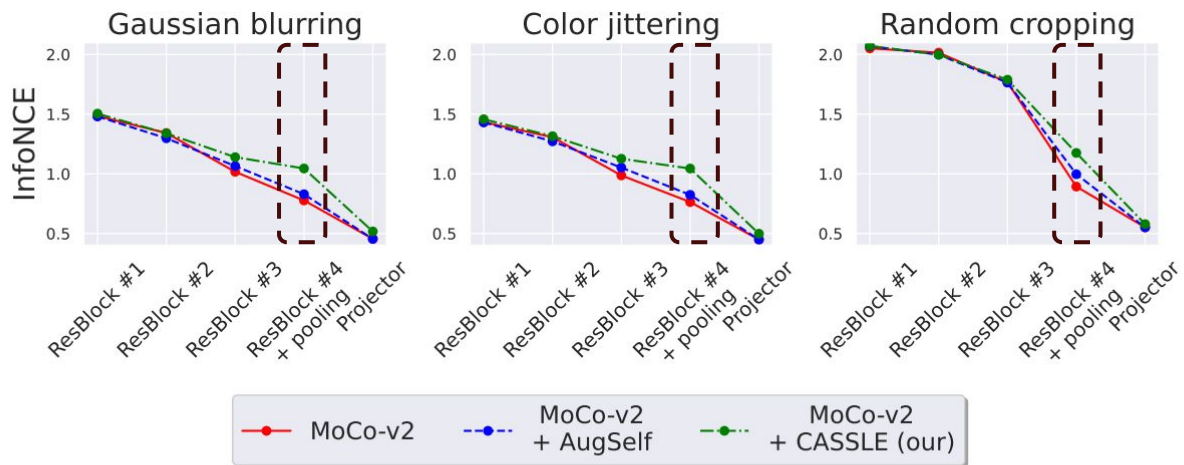
Figure 3: Examples of the commonly-used augmentations and their parameters  $\omega^{\text{aug}}$ .

# Is the knowledge of augmentations useful for SSL tasks?

- We **do not** explicitly force the model to utilize the augmentation information.
- Conditioning the CASSLE projector with wrong augmentation information decreases its ability to draw image pairs together.
- CASSLE projector indeed relies on augmentation information to perform its task.



# CASSLE learns an augmentation-aware data representation



- we measure the error of matching embeddings of augmented image pairs
- embeddings generated by CASSLE are **hardest** to match together (the highest InfoNCE value)
- CASSLE preserves the largest amount of augmentation-induced noise

# CASSLE improves the transferability of SSL models

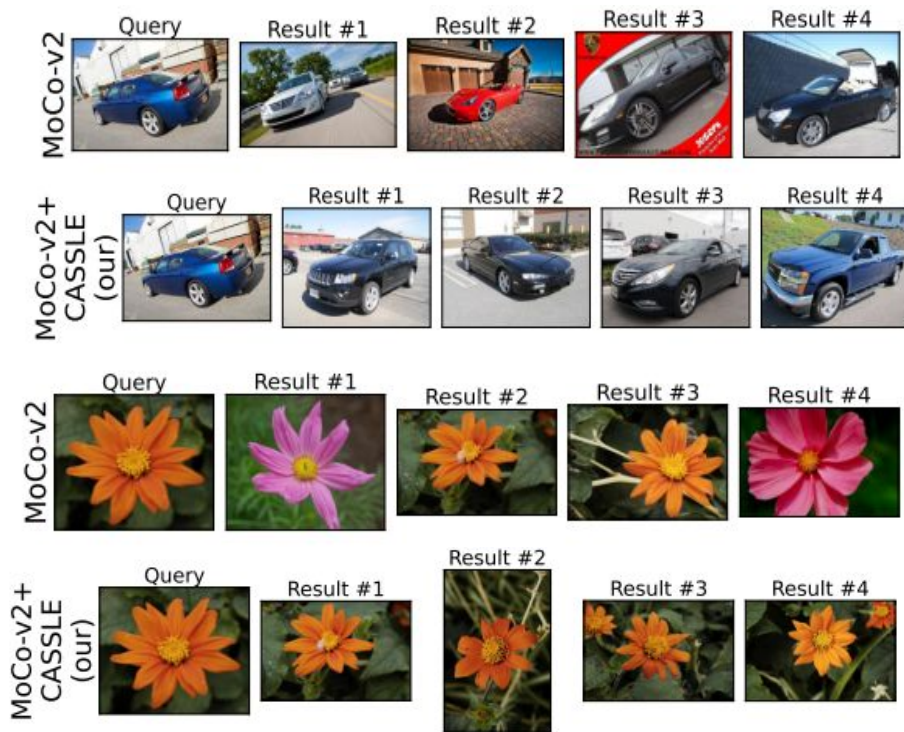
Table 1: Linear evaluation on downstream classification and regression tasks. CASSLE consistently improves representations formed by vanilla SSL approaches and performs better or comparably to other techniques of increasing sensitivity to augmentations [69, 40, 14].

Method	C10	C100	Food	MIT	Pets	Flowers	Caltech	Cars	FGVCA	DTD	SUN	CUB	300W
<i>SimCLR</i> [15]													
Vanilla	84.41 <sup>†</sup>	61.40	57.48 <sup>†</sup>	63.10 <sup>†</sup>	71.60 <sup>†</sup>	83.37 <sup>†</sup>	<b>79.67<sup>†</sup></b>	35.14 <sup>†</sup>	40.03 <sup>†</sup>	64.90	46.92 <sup>†</sup>	30.98 <sup>†</sup>	88.59 <sup>†</sup>
AugSelf [40] <sup>†</sup>	84.45	62.67	59.96	63.21	70.61	<b>85.77</b>	77.78	37.38	42.86	65.53	<b>49.18</b>	34.24	88.27
AI [14]	83.90	63.10	–	–	69.50	68.30	74.20	–	–	53.70	–	<b>38.60</b>	88.00
<b>CASSLE</b>	<b>86.31</b>	<b>64.36</b>	<b>60.67</b>	<b>63.96</b>	<b>72.33</b>	85.22	79.62	<b>39.86</b>	<b>43.10</b>	<b>65.96</b>	48.91	33.21	<b>88.88</b>
<i>MoCo-v2</i> [32, 17]													
Vanilla	84.60	61.60	59.67	61.64	70.08	82.43	77.25	33.86	41.21	64.47	46.50	32.20	88.77 <sup>†</sup>
AugSelf [40]	85.26	63.90	60.78	63.36	73.46	85.70	78.93	37.35	39.47	66.22	48.52	37.00	89.49 <sup>†</sup>
AI [14]	81.30	64.60	–	–	<b>74.00</b>	81.30	78.90	–	–	<b>68.80</b>	–	<b>41.40</b>	<b>90.00</b>
LooC [69]	–	–	–	–	–	–	–	–	–	–	–	39.60	–
IFM [57] <sup>†</sup>	83.36	60.22	59.86	60.60	72.99	85.73	78.77	36.54	41.05	62.34	47.48	35.90	88.92
<b>CASSLE</b>	<b>86.32</b>	<b>65.29</b>	<b>61.93</b>	<b>63.86</b>	72.86	<b>86.51</b>	<b>79.63</b>	<b>38.82</b>	<b>42.03</b>	66.54	<b>49.25</b>	36.22	88.93
<i>MoCo-v3</i> [19] with ViT-Small [23] pretrained on the full ImageNet dataset.													
Vanilla <sup>†</sup>	83.17	62.40	56.15	53.28	62.29	81.48	69.63	28.63	32.84	57.18	42.16	35.00	87.42
AugSelf [40] <sup>†</sup>	84.25	64.12	<b>58.28</b>	<b>56.12</b>	<b>63.93</b>	<b>83.13</b>	72.45	29.64	32.54	<b>60.27</b>	43.22	<b>37.16</b>	87.85
<b>CASSLE</b>	<b>85.13</b>	<b>64.67</b>	57.30	55.90	63.88	82.42	<b>73.53</b>	<b>30.92</b>	<b>35.91</b>	58.24	<b>43.37</b>	36.09	<b>88.53</b>

Table 2: Average Precision of object detection on VOC dataset [25, 42]. CASSLE extension of MoCo-v2 and SimCLR outperforms the vanilla approaches and AugSelf extension by a small margin.

Method	<i>MoCo-v2</i>	<i>SimCLR</i>
Vanilla	45.12	44.74
AugSelf [40]	45.20	44.50
<b>CASSLE</b>	<b>45.90</b>	<b>45.60</b>

# What data is similar for SSL?



# Summary

- augmentation invariance is a key component of modern Self-Supervised Learning
- it can lead to learning representations that are suboptimal for downstream tasks which rely on features of data modified by augmentations
- we propose to increase augmentation-awareness of SSL methods by conditioning them with information about used augmentations



paper



# Thank you for your attention!

group of machine  
**gmum**  
learning research



[gmum.net](http://gmum.net)

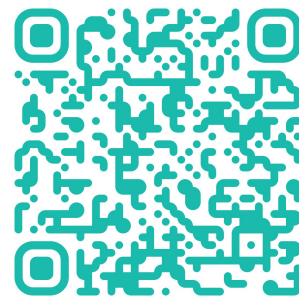


JAGIELLONIAN UNIVERSITY  
IN KRAKÓW



[paper](#)

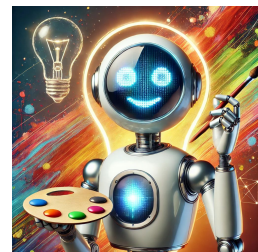
**IDEAS**  
NCBR ○ ○ ○



[ideas-ncbr.pl](http://ideas-ncbr.pl)



Some  
pics by  
DALLE-3



Friday:

Session 2 / Lecture Hall B / 10:35

**Deep learning for effective  
analysis  
of high content screening**

Adriana Borowa

Session 4 / Lecture Hall A / 14:30

**Efficient fine-tuning of LLMs:  
exploring PEFT methods and  
LORA-XS insights**

Klaudia Bałazy

Session 5 / Lecture Hall B / 14:30

**Current trends in intrinsically  
interpretable Deep Learning**

Dawid Rymarczyk

**Neural rendering: the future of  
3D modeling**

Przemysław Spurek

**Check out  
our talks during  
ML in PL 2024!**



Saturday:

Session 7 / Lecture Hall A / 12:00

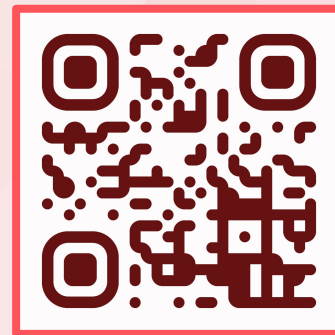
**AdaGlimpse: Active Visual Exploration  
with Arbitrary Glimpse Position and Scale**

Adam Pardyl

Session 8 / Lecture Hall B / 12:00

**Augmentation-aware Self-supervised  
Learning with Conditioned Projector**

Marcin Przewięźlikowski



gmum.net