# "A STUDY ON THE DISTRIBUTION OF SOCIAL BIASES IN SELF-SUPERVISED LEARNING VISUAL MODELS"   Kirill Sirotkin, Pablo Carballeira, Marcos Escudero-Viñolo

## MOTIVATION

- Self-Supervised Learning (SSL) models can learn biases, **contrary** to common belief.

- It is unclear if some **types** of SSL learn **more** biases than others.

## STUDIED SOCIAL BIASES INCLUDE

| Concepts (X, Y) | Attributes (A, B) |
|---|---|
| Male and Female | Career and Family |
| Black and White | Weapon and Tool |
| Skinny and Overweight | Pleasant and Un-pleasant |

## MAIN FINDINGS

- **Contrastive** SSL models learn more biases than geometric or clustering ones.

- Most biases are learned in the **deep layers** of a network.

- SSL-learned biases **transfer** to the downstream task.

- There is a **trade-off** between the accuracy and a number of learned biases.

UAM Universidad Autónoma de Madrid

VPU Video Processing and Understanding Lab



ResNet-50 pre-trains on unlabelled data

Goldfish
Shark
Rock python

*Unlabeled data*

*Pre-training*

ResNet blocks — 1 2 3 4 5 GAP — Contrastive head

*Biases are transferred to the downstream task!*



Family — Male — Female — Career

*Correct predictions are shown in green*



Number of biases

Geometric | Cluster-ing | **Contrastive**

Jigsaw, Rotation, RL, ClusterFit, ODC, SwAV, NPID, MoCo v1, MoCo v2, SimCLR, BYOL, Supervised, Random



Female (X) — Male (Y)
Family (A) — Career (B)

$x_1$, $x_2$, $x_n$

*Images in the feature space*

Bias-detection test

$$\sum_{x \in X} s(x, A, B) - \sum_{y \in Y} s(y, A, B)$$

$$\operatorname*{mean}_{a \in A} \cos(x, a) - \operatorname*{mean}_{b \in B} \cos(x, b)$$

$$\operatorname*{mean}_{a \in A} \cos(y, a) - \operatorname*{mean}_{b \in B} \cos(y, b)$$



Number of biases

BYOL, SimCLR, MoCo v2, MoCo v1, NPID — **Contrastive**

SwAV

ClusterFit, ODC — **Clustering**

Relative Location, Rotation, Jigsaw — **Geometric**

Supervised, Random

*Uncertainty in the bias-detection*