

# SELF-SUPERVISED LEARNING OF CONTEXTUALIZED LOCAL VISUAL EMBEDDINGS

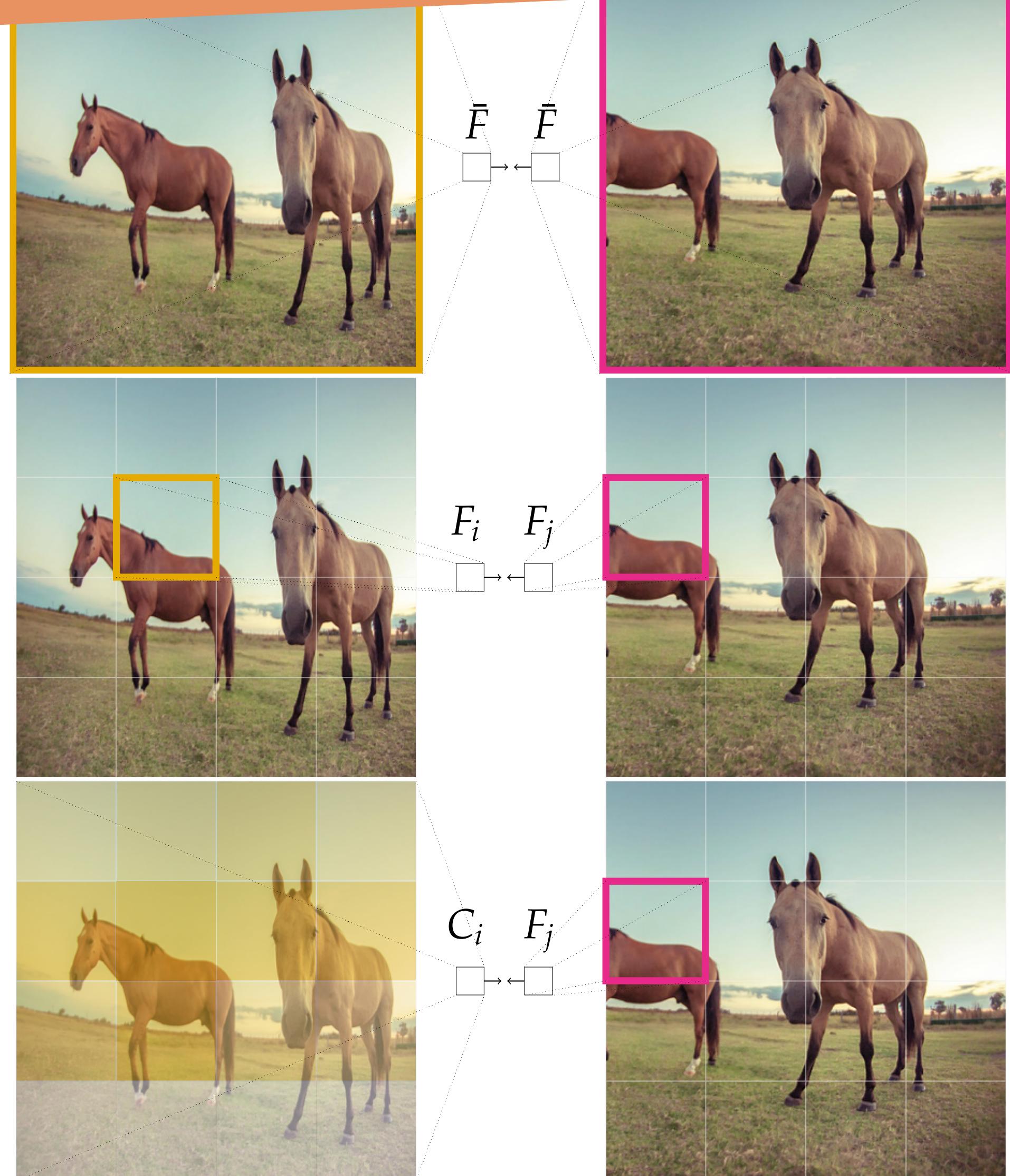


Thalles Silva Helio Pedrini Adín Ramírez Rivera

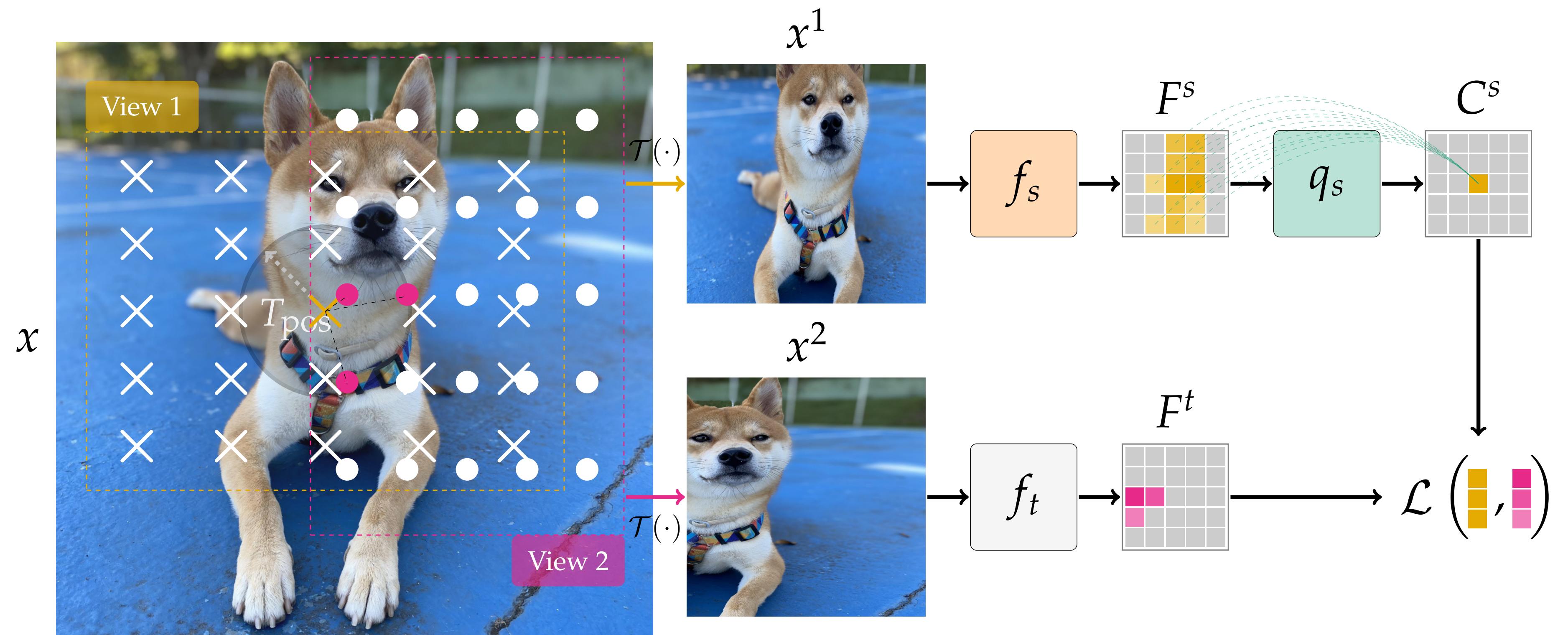
University of Campinas University of Oslo



**Goal** Learn representations that **preserve local information** from the input by finding correlations among similar regions of a view to predict local parts of a different view



## Methodology



- We create views from unlabeled images using random transformations such as flip, color distortions, and cropping.
- Each view is forwarded through a student encoder  $f_s$  and a teacher encoder  $f_t$ .
- Encoders are composed of a feature extractor, e.g., a CNN, and a MLP projection head.
- For each view, we obtain a tensor of projected local feature maps  $F$  taken from the last layers of the CNN encoder (before average pooling).

## Obtaining self-supervised dense targets

### Strategy

match local features based on the pixel's spatial localities

- $I^1$  and  $I^2$  are lists of 2D points in the pixel space for each view
- For each point  $I_i^1$ , we look for pixel correspondences in  $I^2$  to create  $M$
- $M$  is the set of pairs  $(I_i^1, I_j^2)$  such that the euclidean distance between points  $I_i^1$  and  $I_j^2$  is smaller than a threshold  $T_{\text{pos}}$

$$M = \left\{ \left( I_i^1, I_j^2 \right) : d \left( I_i^1, I_j^2 \right) < T_{\text{pos}} \right\}, \quad (1)$$

- We map points in  $M$  from the pixel space to the feature space, to obtain a pair of indices representing matching features from the two views.

## Results

### Obj. detection and segmentation on COCO (R50-FPN)

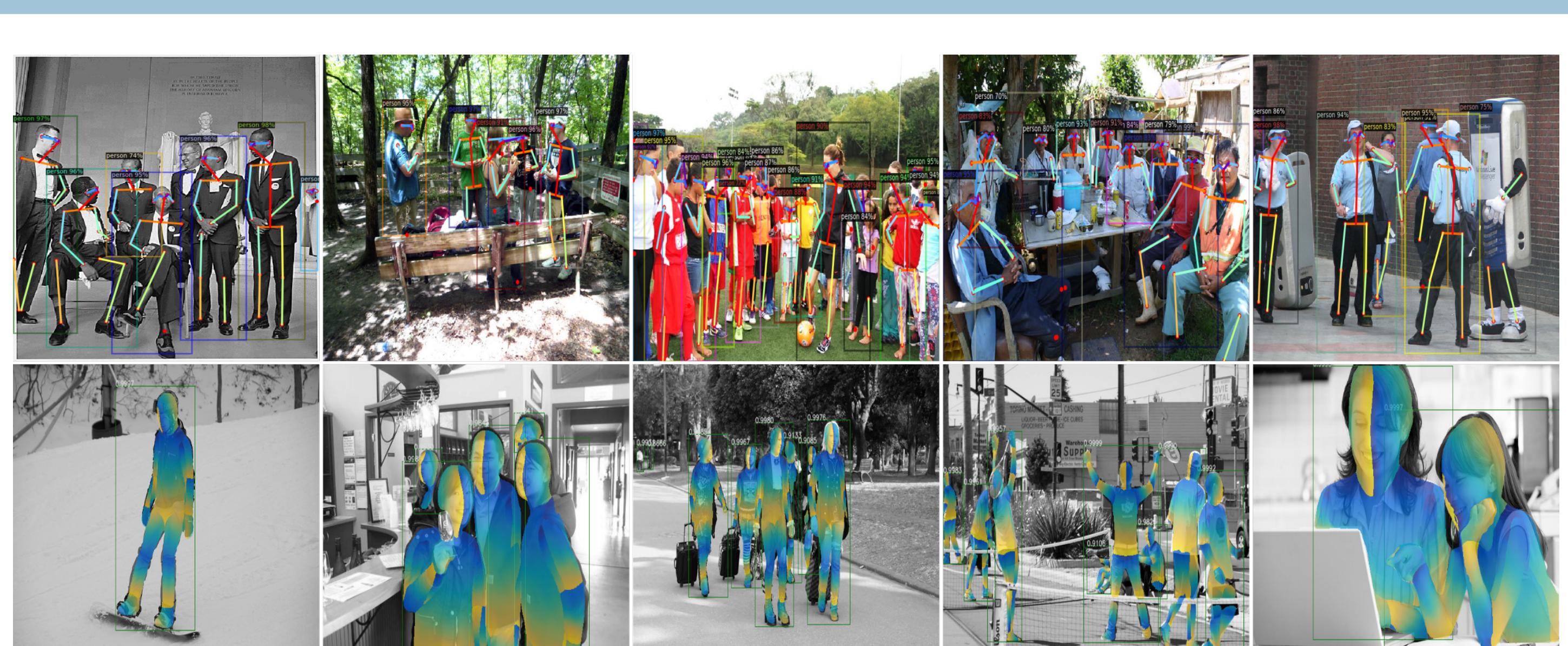
Method	ep	AP <sup>bb</sup>	AP <sub>50</sub> <sup>bb</sup>	AP <sub>75</sub> <sup>bb</sup>	AP <sup>mb</sup>	AP <sub>50</sub> <sup>mb</sup>	AP <sub>75</sub> <sup>mb</sup>
Supervised	100	38.9	59.6	42.7	35.4	56.5	38.1
Rand init	-	32.8	51	35.3	28.5	46.8	30.4
DenseCL	200	39.4	59.9	42.7	35.6	56.7	38.2
ReSim	200	39.3	59.7	43.1	35.7	56.7	38.1
PixPro	400	39.8	59.5	43.7	36.1	56.5	38.9
SetSim	200	40.2	60.7	43.9	36.4	57.7	39
VICRegL	300	37.3	57.6	40.7	34.1	54.7	36.5
CLoVE	200	40.8	60.5	45.0	36.8	57.6	39.8
	400	<b>41.2</b>	<b>61.1</b>	<b>45</b>	<b>37.1</b>	<b>58.1</b>	<b>40.1</b>

### Instance segmentation on Cityscapes (R50-FPN)

Method	ep	AP	AP <sub>50</sub>
Supervised	100	26.5	52.9
Rand init	-	19.9	40.7
DenseCL	200	33.1	61.7
PixPro	400	<b>35.8</b>	63.7
VICRegL	300	29.8	58.5
SlotCon	200	35.2	63.8
CLoVE	200	35.7	<b>64.1</b>
	400	<b>37.2</b>	<b>65.3</b>

### Keypoint detection on COCO (R50-FPN)

Method	ep	AP <sup>kp</sup>	AP <sub>50</sub> <sup>kp</sup>	AP <sub>75</sub> <sup>kp</sup>
Supervised	100	65.3	87	71.3
Rand init	-	63	85.1	68.4
DenseCL	200	66.3	87.1	71.9
PixPro	400	66.6	87.2	73.0
ReSim	200	66.3	87.2	72.4
SetSim	200	66.7	<b>87.8</b>	72.4
SlotCon	200	66.5	87.5	72.5
CLoVE	200	66.9	87.5	73.2
	400	<b>67.0</b>	87.4	<b>73.3</b>



## Acknowledgements

We thank Sigma2 (the National Infrastructure for High Performance Computing and Data Storage in Norway), Project NN8104K; the RCN Centre for Research-based Innovation funding scheme (grant no. 309439); and Coordenação de Aperfeiçoamento de Pessoal de Nível Superior—Brasil (CAPES)—Finance Code 001.