

DeepMind

# The virtuous cycle of object discovery and representation learning

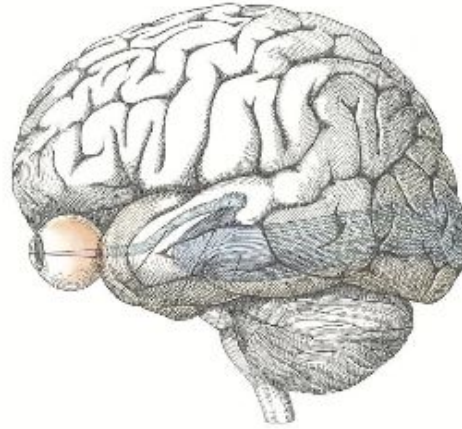
Olivier Hénaff

ECCV 2022 Workshop on Self-Supervised  
Representation Learning in Computer Vision



# Principles of biological and artificial intelligence

Efficient generalization



Neural representations

- Enable intelligent behavior

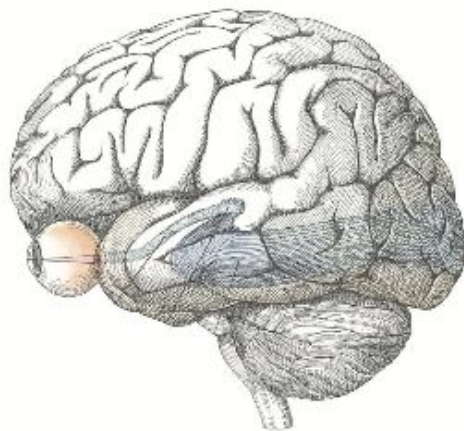


# Principles of biological and artificial intelligence

## Self-supervised pretraining



*But thy eternal summer shall not fade,  
Nor lose possession of that fair thou ow'st,  
Nor shall death brag thou wander'st in his shade,  
When in eternal lines to time thou grow'st,  
So long as men can breathe, or eyes can see,  
So long lives this, and this gives life to thee.*



## Neural representations

- Enable intelligent behavior
- Require minimal supervision
- Are generally applicable

## Efficient generalization

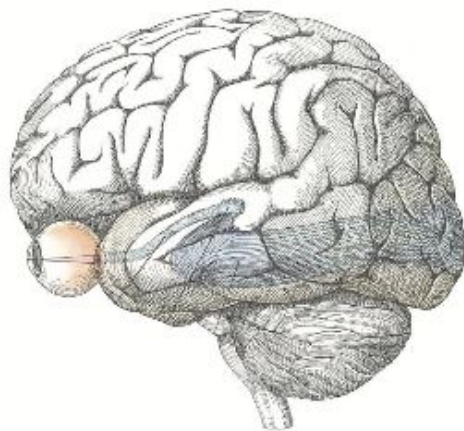


# Principles of biological and artificial intelligence

## Self-supervised pretraining



*But thy eternal summer shall not fade,  
Nor lose possession of that fair thou ow'st,  
Nor shall death brag thou wander'st in his shade,  
When in eternal lines to time thou grow'st,  
So long as men can breathe, or eyes can see,  
So long lives this, and this gives life to thee.*



## Neural representations

- Enable intelligent behavior
- Require minimal supervision
- Are generally applicable

## Efficient generalization



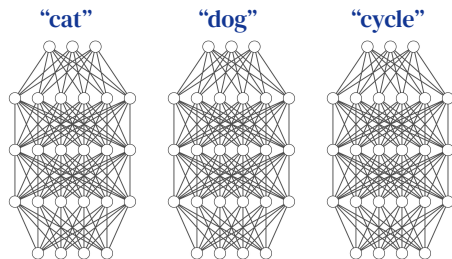
- I'm looking for a place to eat
- I know lots of restaurants!  
Anything specific?
- I love chatbots
- Everybody does!
- So, which cuisine?
- A cheap one
- Ok, and what city?
- Actually, something fancy
- Ok, expensive it is.
- So which city?
- Can you show me some  
restaurants, yet?
- I need some more info first.  
Where would you like to eat?

# Principles of biological and artificial intelligence

## Self-supervised pretraining



*But thy eternal summer shall not fade,  
Nor lose possession of that fair thou ow'st,  
Nor shall death brag thou wander'st in his shade,  
When in eternal lines to time thou grow'st,  
So long as men can breathe, or eyes can see,  
So long lives this, and this gives life to thee.*



Supervised ConvNet

## Neural representations

- Enable intelligent behavior
- Require minimal supervision
- Are generally applicable

## Efficient generalization



- I'm looking for a place to eat
- I know lots of restaurants!  
Anything specific?
- I love chatbots
- Everybody does!
- So, which cuisine?
- A cheap one
- Ok, and what city?
- Actually, something fancy
- Ok, expensive it is.
- So which city?
- Can you show me some restaurants yet?
- I need some more info first.  
Where would you like to eat?

# Principles of biological and artificial intelligence

## Self-supervised pretraining



*But thy eternal summer shall not fade,  
Nor lose possession of that fair thou ow'st,  
Nor shall death brag thou wander'st in his shade,  
When in eternal lines to time thou grow'st,  
So long as men can breathe, or eyes can see,  
So long lives this, and this gives life to thee.*



image colorization  
(Zhang, 2016)

## Neural representations

- Enable intelligent behavior
- Require minimal supervision
- Are generally applicable

## Efficient generalization



- I'm looking for a place to eat
- I know lots of restaurants!  
Anything specific?
- I love chatbots
- Everybody does!
- So, which cuisine?
- A cheap one
- Ok, and what city?
- Actually, something fancy
- Ok, expensive it is.
- So which city?
- Can you show me some  
restaurants yet?
- I need some more info first.  
Where would you like to eat?

# Principles of biological and artificial intelligence

## Self-supervised pretraining



*But thy eternal summer shall not fade,  
Nor lose possession of that fair thou ow'st,  
Nor shall death brag thou wander'st in his shade,  
When in eternal lines to time thou grow'st,  
So long as men can breathe, or eyes can see,  
So long lives this, and this gives life to thee.*



## Efficient generalization



I'm looking for a place to eat

I know lots of restaurants!  
Anything specific?

I love chatbots

Everybody does!

So, which cuisine?

A cheap one

Ok, and what city?

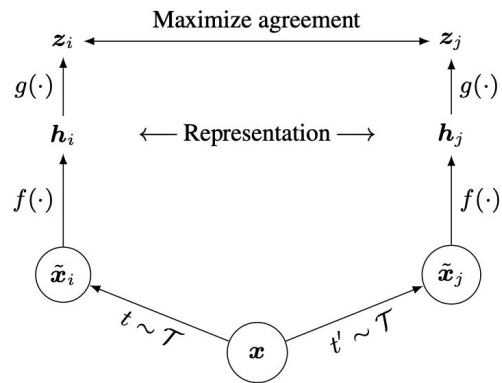
Actually, something fancy

Ok, expensive it is.

So which city?

Can you show me some restaurants yet?

I need some more info first.  
Where would you like to eat?

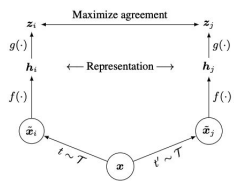
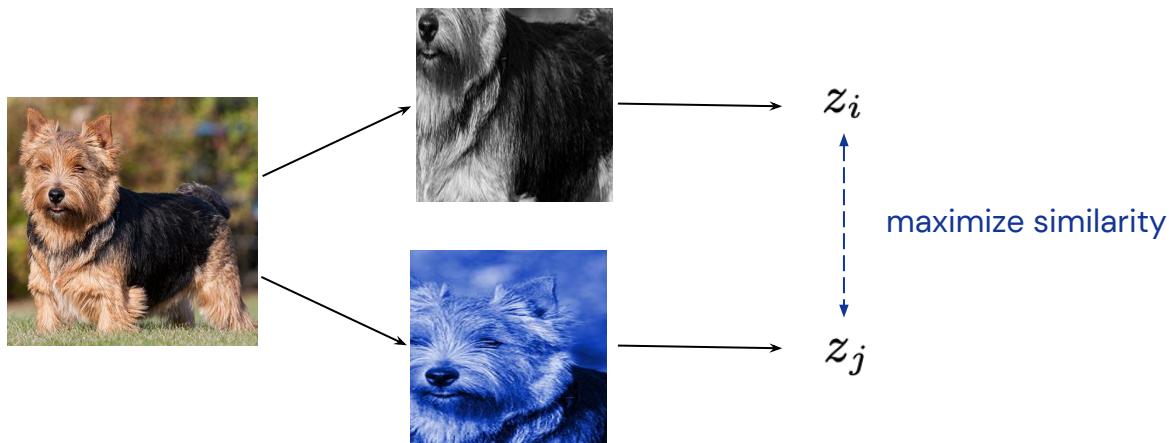


Contrastive learning (Chen, 2020)

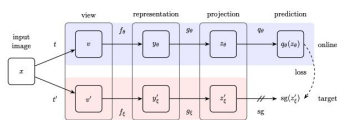
### Neural representations

- Enable intelligent behavior
- Require minimal supervision
- Are generally applicable

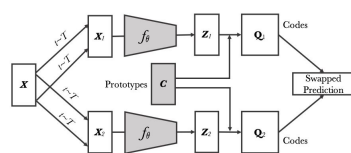
# Is the current self-supervised paradigm too simple?



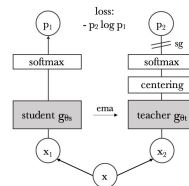
SimCLR



BYOL



SwAV



DINO





# Is the current self-supervised paradigm too simple?

Real-world data is complex

- multiple objects in natural scenes
- multiple speakers in natural speech
- multiple scenes in natural videos

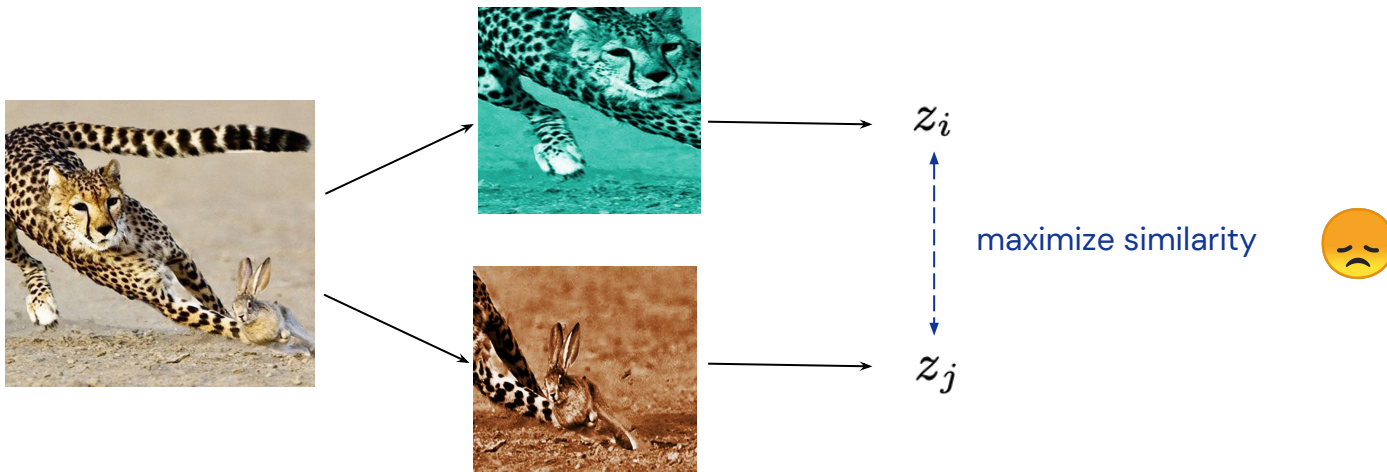


# Is the current self-supervised paradigm too simple?

Real-world data is complex

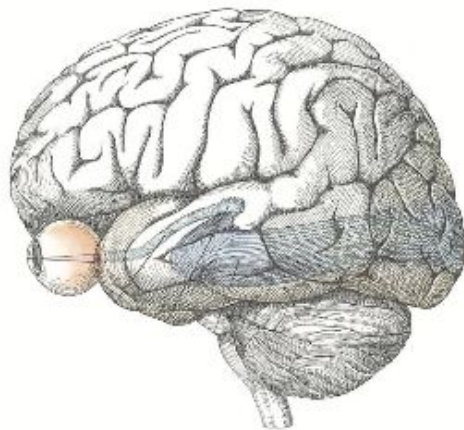
- multiple objects in natural scenes
- multiple speakers in natural speech
- multiple scenes in natural videos

→ invariance across views dampens instance selectivity



# Is the current self-supervised paradigm too simple?

Pretrain a ResNet-50 on ImageNet



Neural representations

- Enable intelligent behavior
- Require minimal supervision
- Are domain-agnostic

Fine-tune for object detection, segmentation

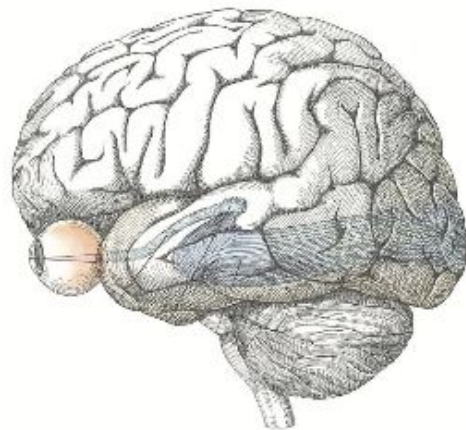


Fine-tune on

- Semantic segmentation (PASCAL or ADE20K)
- Object detection (COCO or LVIS)

# Is the current self-supervised paradigm too simple?

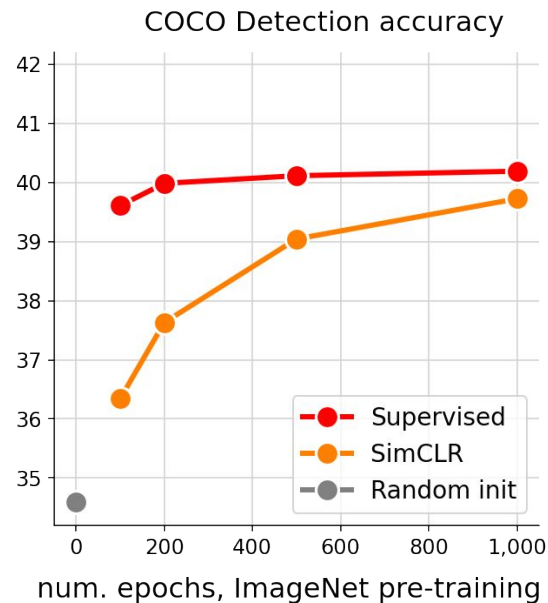
Pretrain a ResNet-50 on ImageNet



**Neural representations**

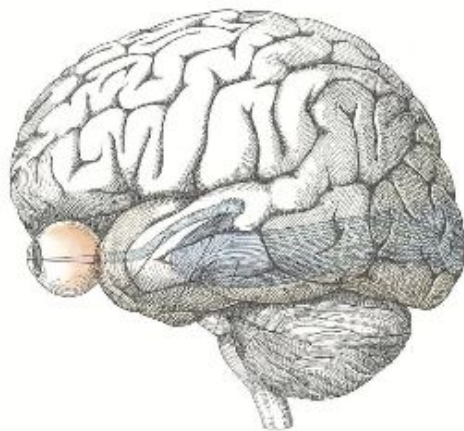
- Enable intelligent behavior
- Require minimal supervision
- Are domain-agnostic

Fine-tune for object detection, segmentation



# Is the current self-supervised paradigm too simple?

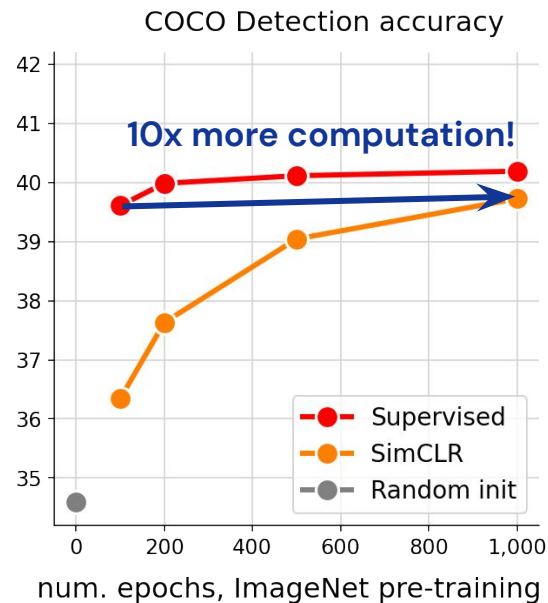
Pretrain a ResNet-50 on ImageNet



**Neural representations**

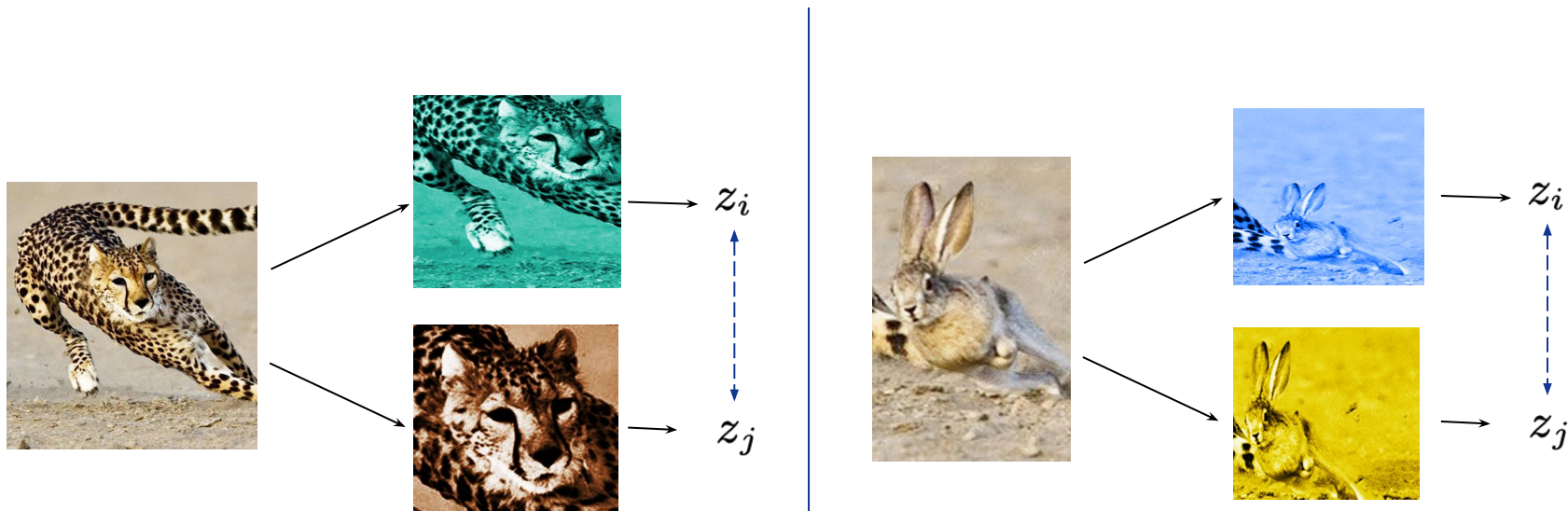
- Enable intelligent behavior
- Require minimal supervision
- Are domain-agnostic

Fine-tune for object detection, segmentation



# Hypothesis for handling real-world data

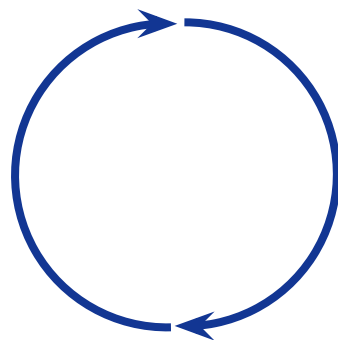
- Break images into their constituent objects
- Apply contrastive learning to each **object** rather than each **image**



# Outline

1. **Knowledge of objects accelerates and improves representation learning**  
→ DetCon objective (ICCV, 2021)
2. **Knowledge of objects can be extracted from learned representations**  
→ Odin framework (ECCV, 2022)
3. **Videos can be used to learn strong image representations**  
→ VITO framework (arXiv, 2022)

Better object knowledge



Better representations

DeepMind

# Efficient visual pretraining with contrastive detection

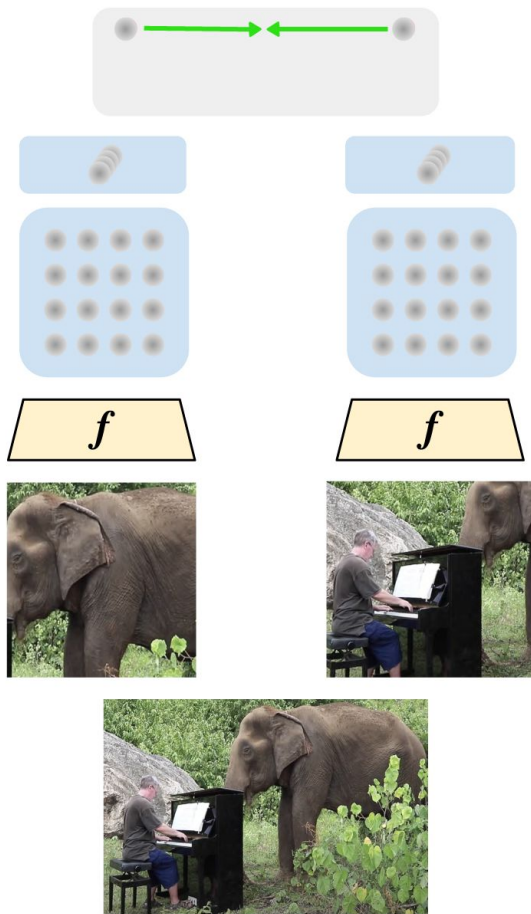
Olivier Hénaff, Skanda Koppula, Jean-Baptiste Alayrac,  
Aaron van den Oord, Oriol Vinyals, João Carreira

ICCV 2021





# Contrastive learning



Contrastive objective

$$\mathcal{L} = -\log \frac{\exp(\mathbf{v} \cdot \mathbf{v}')}{\exp(\mathbf{v} \cdot \mathbf{v}') + \sum_n \exp(\mathbf{v} \cdot \mathbf{v}_n)}$$

Global pooling

Convolutional features

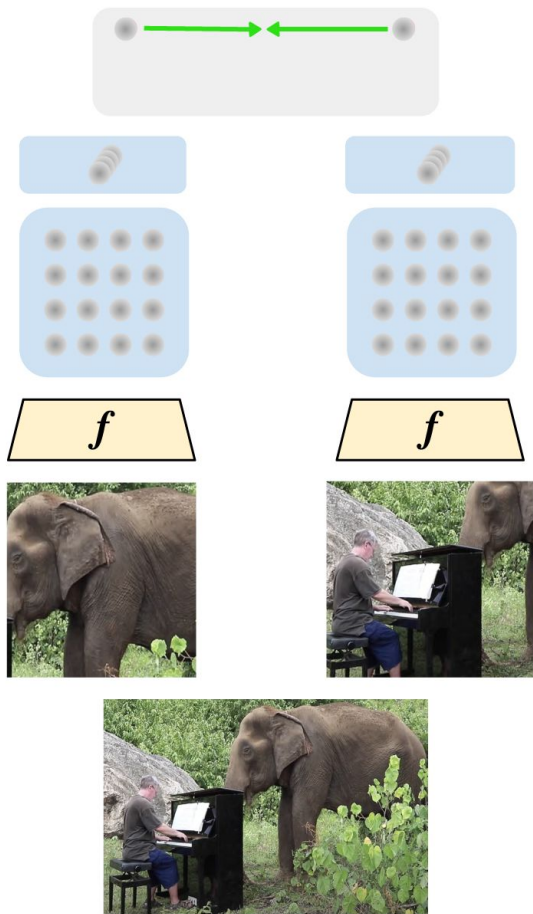
Encoder

Augmented views

Training image



# Contrastive learning



Contrastive objective

Global pooling

Convolutional features

Encoder

Augmented views

Training image

1 positive pair per image

$$\mathcal{L} = -\log \frac{\exp(\mathbf{v} \cdot \mathbf{v}')}{\exp(\mathbf{v} \cdot \mathbf{v}') + \sum_n \exp(\mathbf{v} \cdot \mathbf{v}_n)}$$

1 negative sample per image

- Each image contributes a single positive pair and negative sample
- Positive pairs can be semantically different



# Contrastive detection



Training image and  
heuristic masks



# Contrastive detection



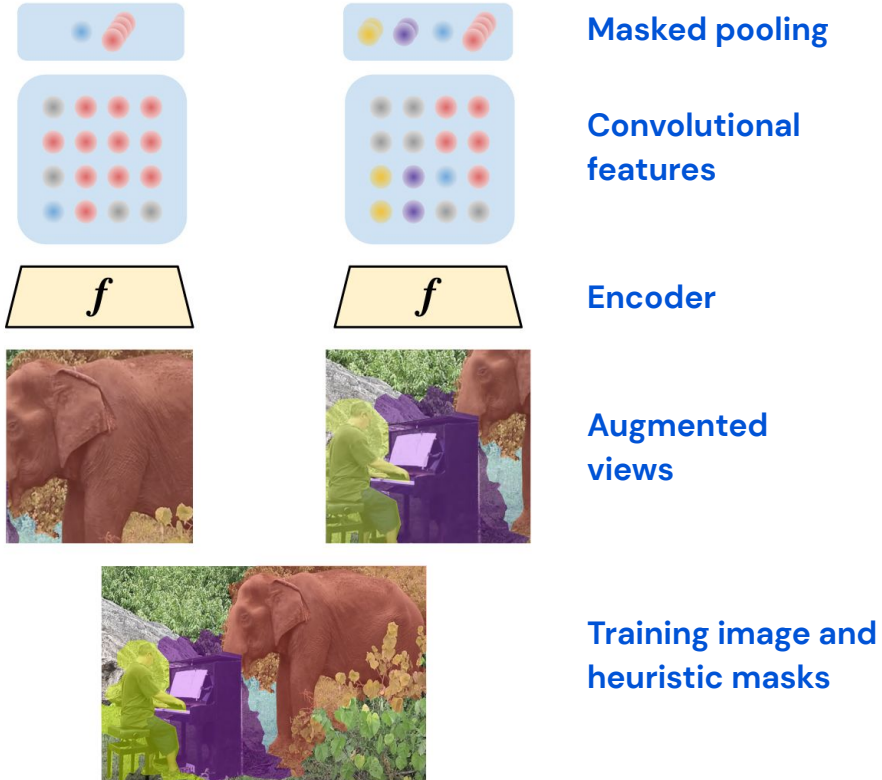
Augmented views



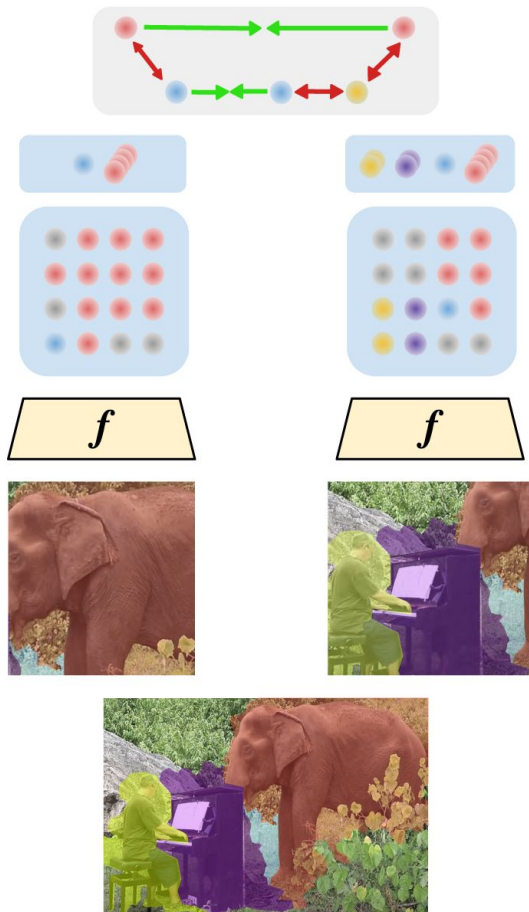
Training image and heuristic masks



# Contrastive detection



# Contrastive detection



DetCon objective

Masked pooling

Convolutional features

Encoder

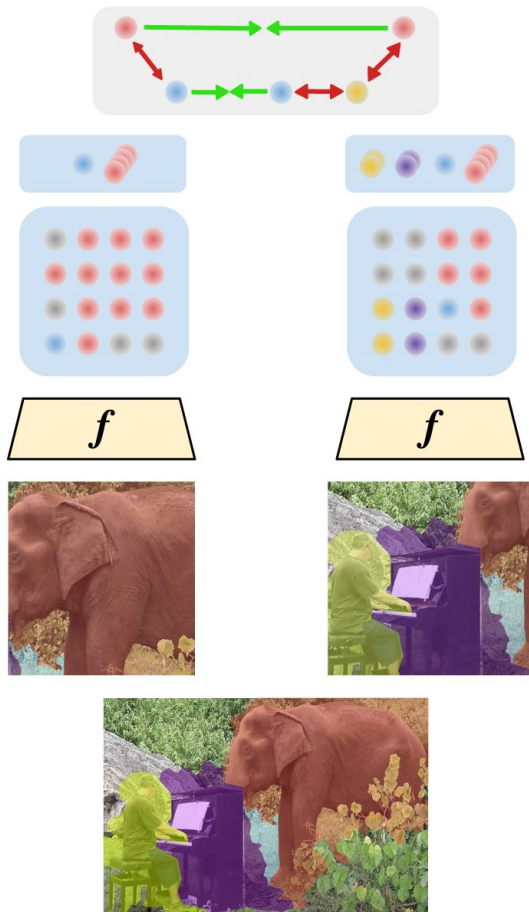
Augmented views

Training image and heuristic masks

$$\mathcal{L} = - \sum_m \log \frac{\exp(\mathbf{v}_m \cdot \mathbf{v}'_m)}{\exp(\mathbf{v}_m \cdot \mathbf{v}'_m) + \sum_n \exp(\mathbf{v}_m \cdot \mathbf{v}_n)}$$



# Contrastive detection



DetCon objective

Masked pooling

Convolutional features

Encoder

Augmented views

Training image and heuristic masks

M positive pairs per image

$$\mathcal{L} = - \sum_m \log \frac{\exp(\mathbf{v}_m \cdot \mathbf{v}'_m)}{\exp(\mathbf{v}_m \cdot \mathbf{v}'_m) + \sum_n \exp(\mathbf{v}_m \cdot \mathbf{v}_n)}$$

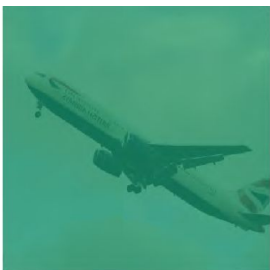
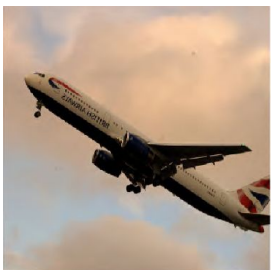
M negative samples per image

- Each image contributes multiple positive pairs and negative samples
- Positive pairs are spatially aligned



# Unsupervised segmentation

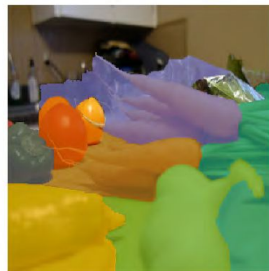
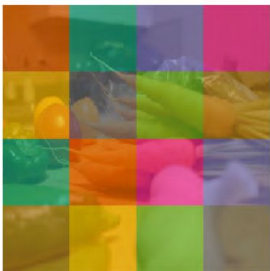
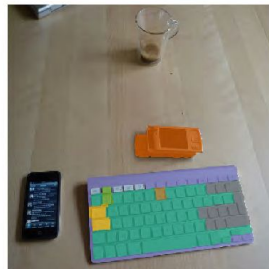
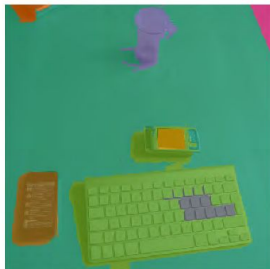
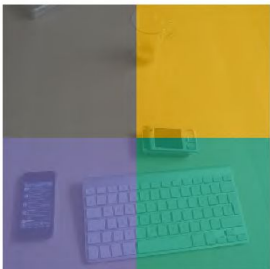
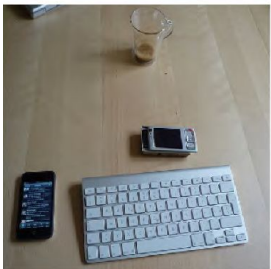
Spatial heuristic



Heuristic: FH



Heuristic: MCG



*Felzenszwalb & Huttenlocher 2004*

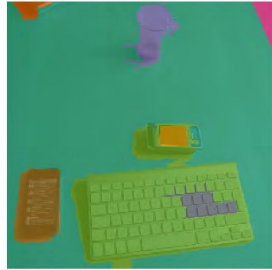
*Arbeláez et al. 2014*





# Unsupervised segmentation

Heuristic: FH



*Felzenszwalb & Huttenlocher 2004*

[skimage.segmentation.felzenszwalb](http://skimage.segmentation.felzenszwalb)

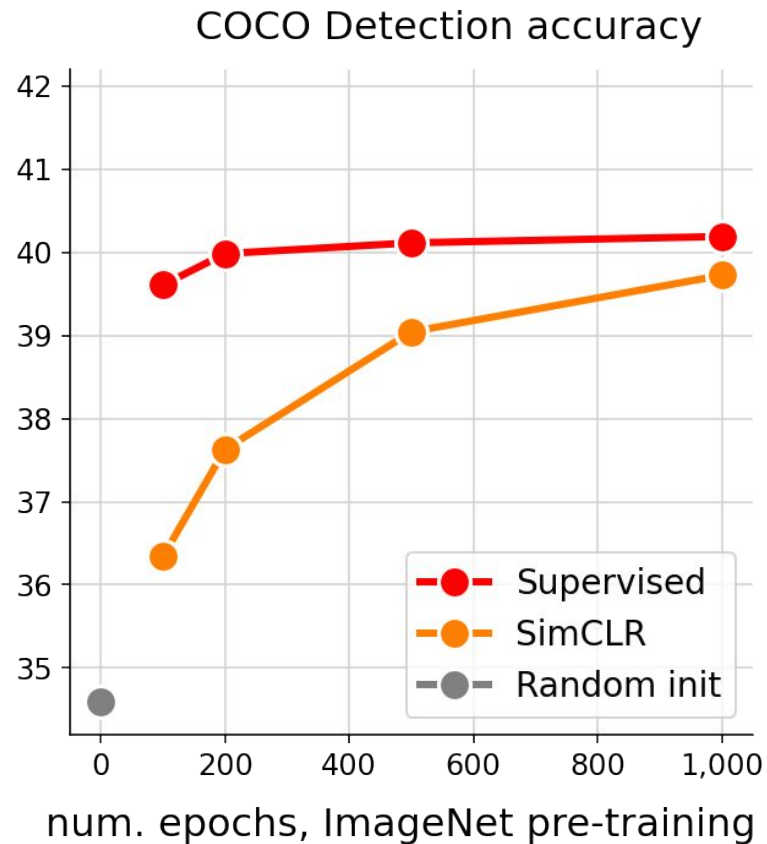


# Experiments

**Pretrain** ResNet-50 on ImageNet

**Objective:** Supervised, SimCLR, or DetCon

**Transfer** to COCO detection and instance segmentation using Mask-RCNN

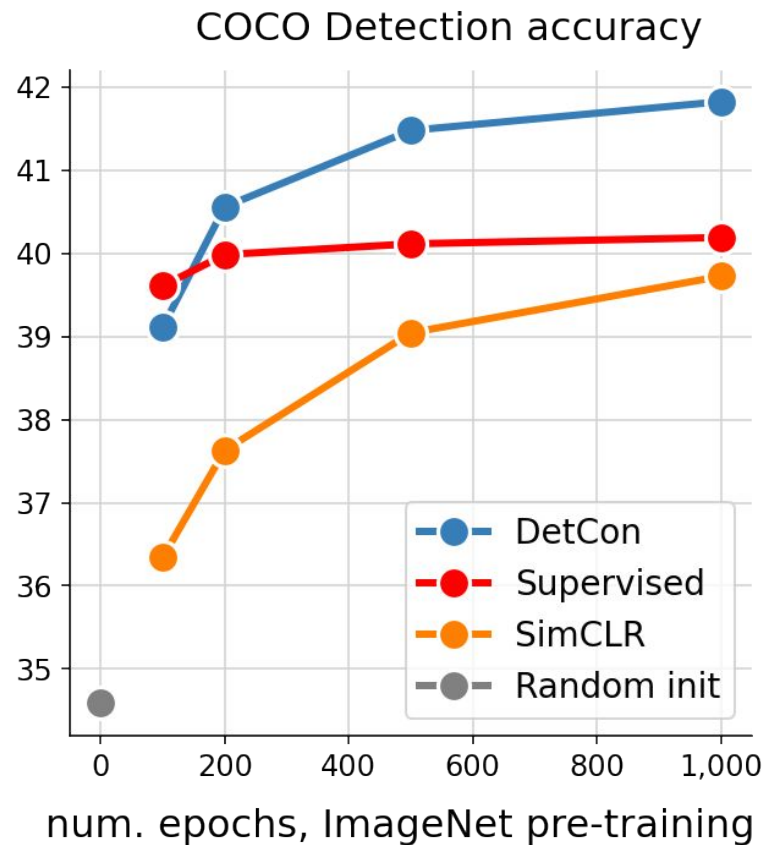


# Experiments

**Pretrain** ResNet-50 on ImageNet

**Objective:** Supervised, SimCLR, or DetCon

**Transfer** to COCO detection and instance segmentation using Mask-RCNN

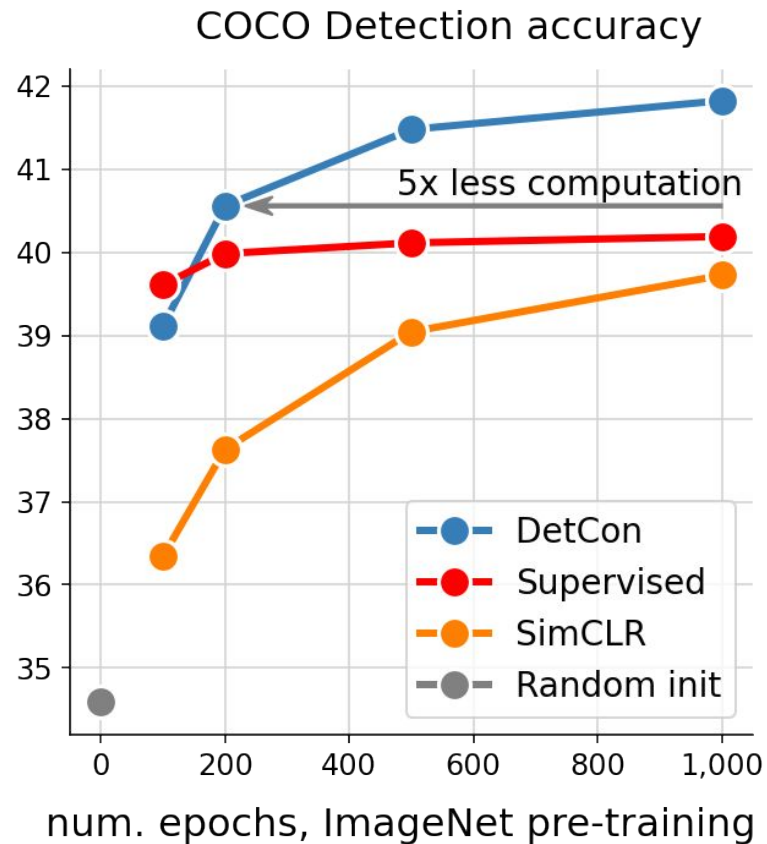


# Experiments

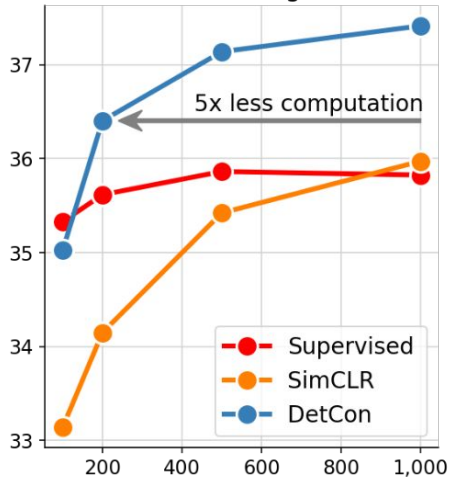
**Pretrain** ResNet-50 on ImageNet

**Objective:** Supervised, SimCLR, or DetCon

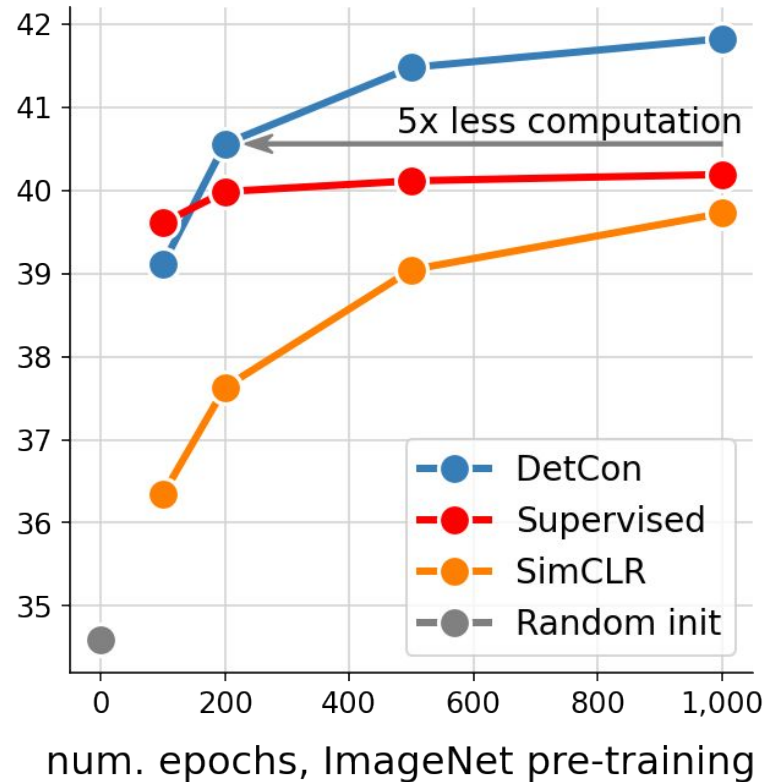
**Transfer** to COCO detection and instance segmentation using Mask-RCNN



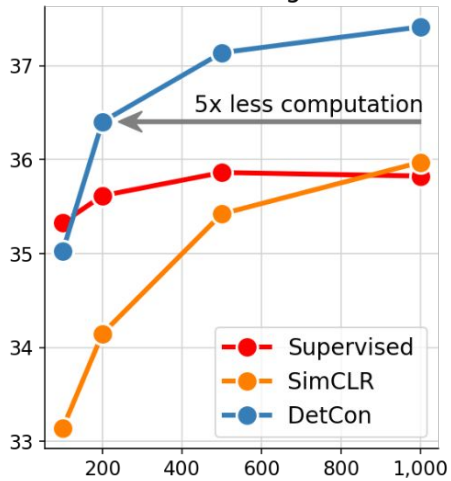
### COCO Instance Segmentation



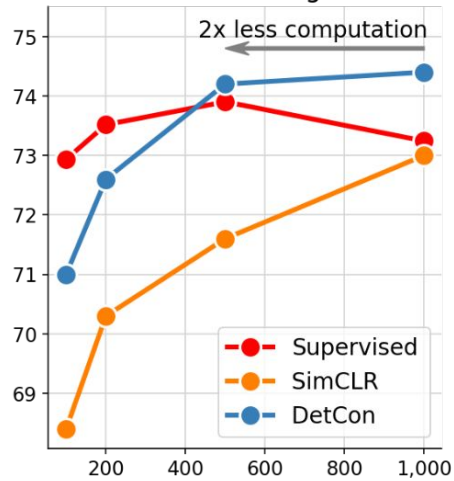
### COCO Detection accuracy



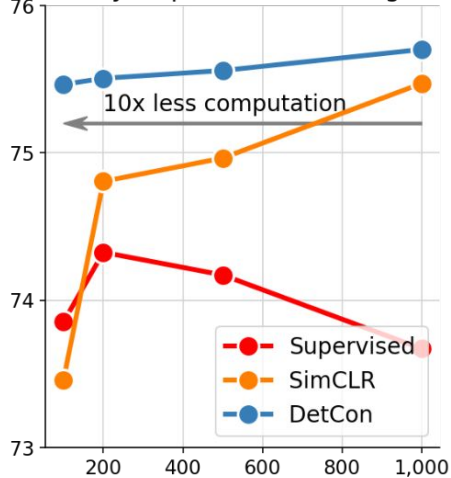
COCO Instance Segmentation



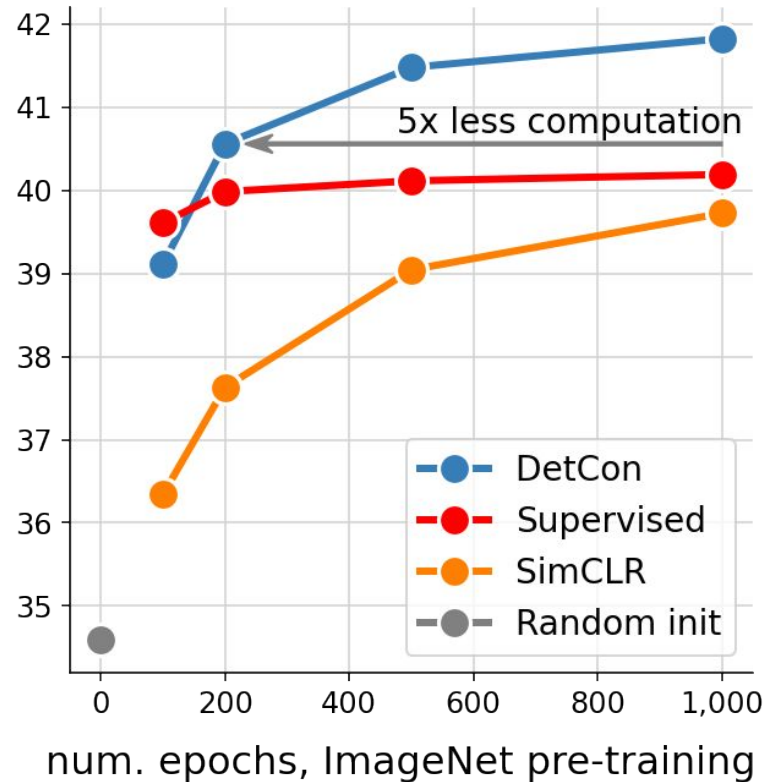
PASCAL Semantic Segmentation



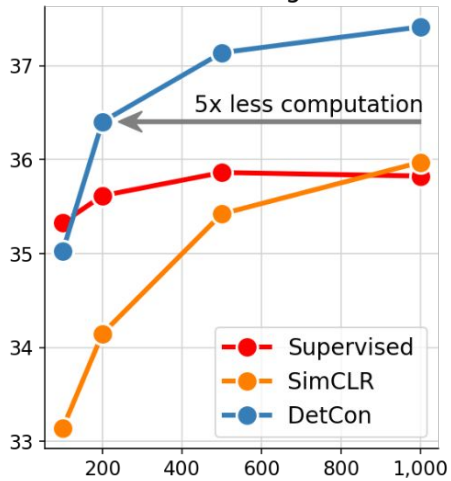
Cityscapes Semantic seg.



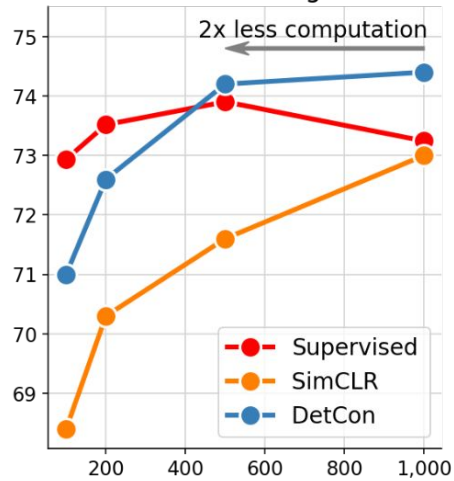
COCO Detection accuracy



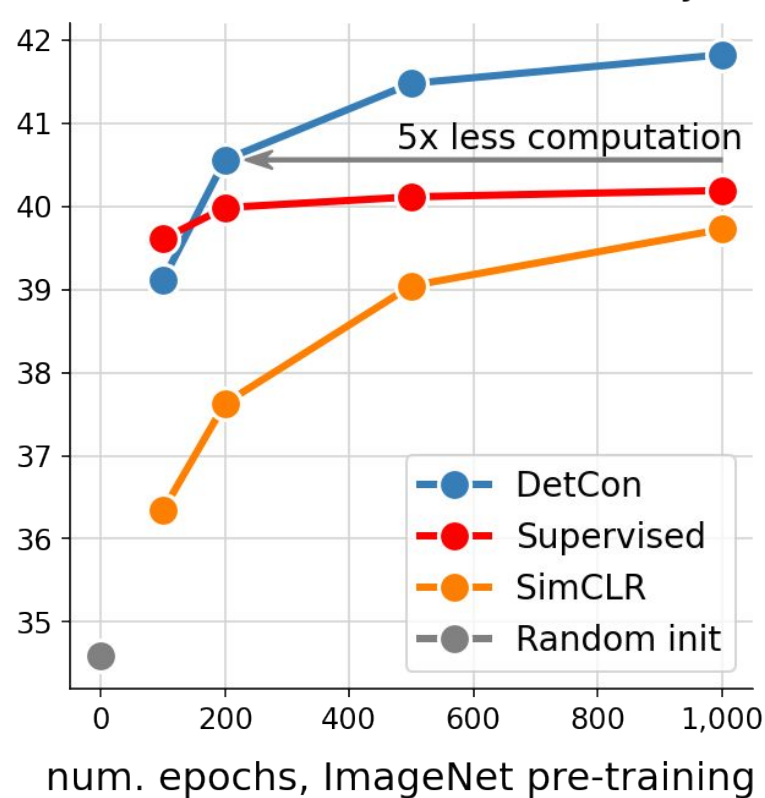
COCO Instance Segmentation



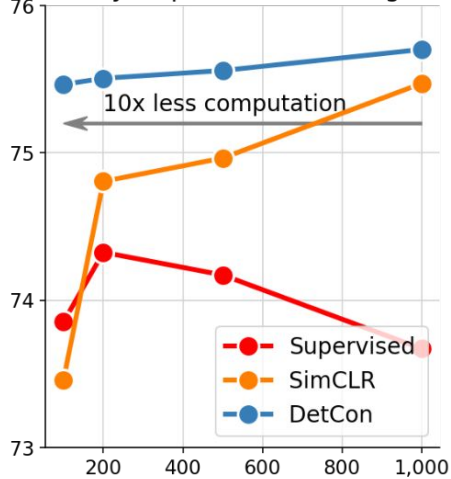
PASCAL Semantic Segmentation



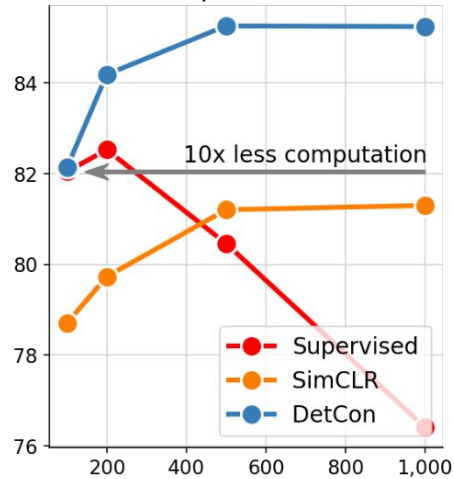
COCO Detection accuracy



Cityscapes Semantic seg.



NYU Depth Estimation



# Experiments

2-10x gain in efficiency

SimCLR  DetCon



# Experiments

2-10x gain in efficiency

SimCLR  $\longrightarrow$  DetCon<sub>S</sub>

+ prediction heads  
+ target networks  
+ better augmentations

BYOL  $\longrightarrow$  DetCon<sub>B</sub>

gain in efficiency?

# Experiments

2-10x gain in efficiency

SimCLR  $\longrightarrow$  DetCon<sub>S</sub>

+ prediction heads  
+ target networks  
+ better augmentations

BYOL  $\longrightarrow$  DetCon<sub>B</sub>

gain in efficiency?

Pretrain epochs	Detection COCO	
	300	1000
BYOL	41.2	41.6
<b>DetCon<sub>B</sub></b>	<b>42.0</b>	42.7
Efficiency Gain	> 3x	

# Experiments

2-10x gain in efficiency

SimCLR  $\longrightarrow$  DetCon<sub>S</sub>

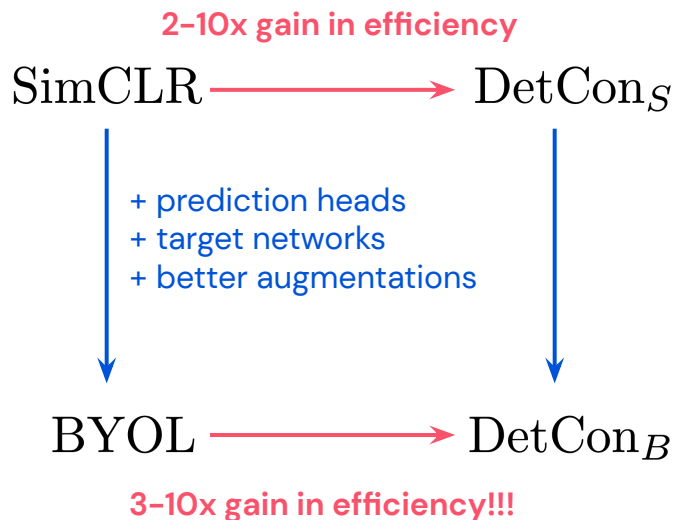
+ prediction heads  
+ target networks  
+ better augmentations

BYOL  $\longrightarrow$  DetCon<sub>B</sub>

3-10x gain in efficiency!!!

	Detection COCO		Instance Segmentation COCO		Semantic Segmentation PASCAL		Semantic Segmentation Cityscapes		Depth Estimation NYU v2	
	300	1000	300	1000	300	1000	300	1000	100	1000
Pretrain epochs	300	1000	300	1000	300	1000	300	1000	100	1000
BYOL	41.2	41.6	37.1	37.2	74.7	75.7	73.4	74.6	83.7	84.2
DetCon <sub>B</sub>	42.0	42.7	37.8	38.2	75.6	77.3	75.1	77.0	85.1	86.3
Efficiency Gain	> 3x		> 3x		≈ 3x		> 3x		> 10x	

# Experiments

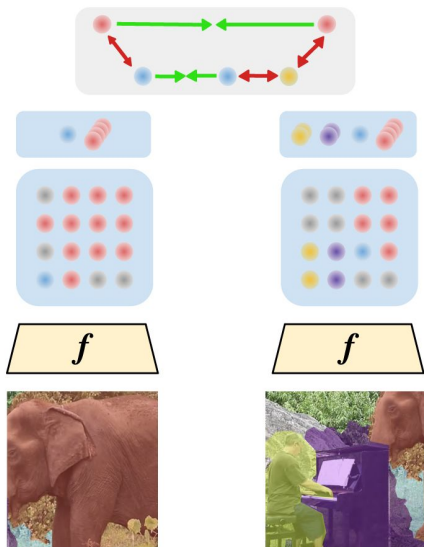


method	Fine-tune 1×		Fine-tune 2×	
	AP <sup>bb</sup>	AP <sup>mk</sup>	AP <sup>bb</sup>	AP <sup>mk</sup>
Supervised	39.6	35.6	41.6	37.6
VADeR	39.2	35.6	-	-
MoCo	39.4	35.6	41.7	37.5
SimCLR	39.7	35.8	41.6	37.4
MoCo v2	40.1	36.3	41.7	37.6
InfoMin	40.6	36.7	42.5	38.4
PixPro	41.4	-	-	-
BYOL	41.6	37.2	42.4	38.0
SwAV	41.6	37.8	-	-
<b>DetCon<sub>S</sub></b>	41.8	37.4	42.9	38.1
<b>DetCon<sub>B</sub></b>	<b>42.7</b>	<b>38.2</b>	<b>43.4</b>	<b>38.7</b>

	Detection COCO		Instance Segmentation COCO		Semantic Segmentation PASCAL		Semantic Segmentation Cityscapes		Depth Estimation NYU v2	
	300	1000	300	1000	300	1000	300	1000	100	1000
Pretrain epochs	300	1000	300	1000	300	1000	300	1000	100	1000
BYOL	41.2	41.6	37.1	37.2	74.7	75.7	73.4	74.6	83.7	84.2
<b>DetCon<sub>B</sub></b>	<b>42.0</b>	42.7	<b>37.8</b>	38.2	<b>75.6</b>	77.3	<b>75.1</b>	77.0	<b>85.1</b>	86.3
Efficiency Gain	> 3×		> 3×		≈ 3×		> 3×		> 10×	

# Knowledge of objects improves representation learning

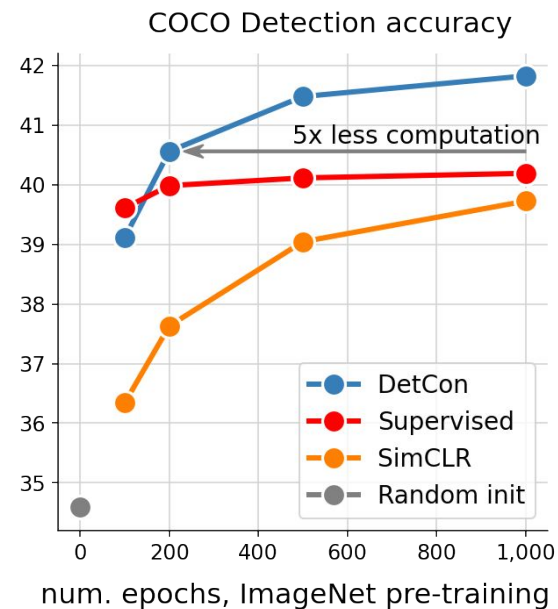
Pretrain a ResNet-50 on ImageNet



Neural representations

- Enable intelligent behavior
- Require minimal supervision
- Are generally applicable

Fine-tune for object detection, segmentation



# Knowledge of objects improves representation learning

Pretrain a ResNet-50 on ImageNet

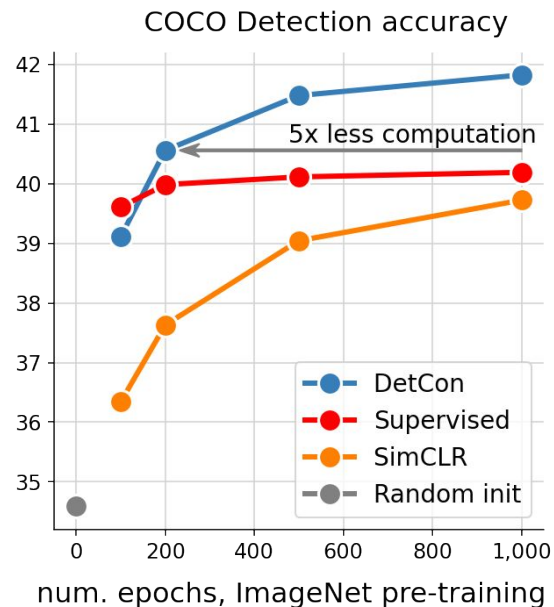


0. Sort  $E$  into  $\pi = (o_1, \dots, o_m)$ , by non-decreasing
1. Start with a segmentation  $S^0$ , where each vert
2. Repeat step 3 for  $q = 1, \dots, m$ .
3. Construct  $S^q$  given  $S^{q-1}$  as follows. Let  $v_i$  and by the  $q$ -th edge in the ordering, i.e.,  $o_q = (v_i$  components of  $S^{q-1}$  and  $w(o_q)$  is small compa both those components, then merge the two co More formally, let  $C_i^{q-1}$  be the component of  $i$  component containing  $v_j$ . If  $C_i^{q-1} \neq C_j^{q-1}$  and  $S^q$  is obtained from  $S^{q-1}$  by merging  $C_i^{q-1}$  and
4. Return  $S = S^m$ .

## Neural representations

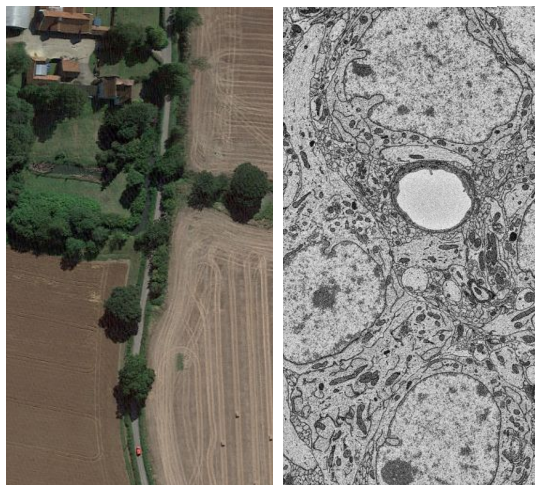
- Enable intelligent behavior
- Require minimal supervision
- Are generally applicable

Fine-tune for object detection, segmentation



# Can the self-supervised paradigm stay general?

## Self-supervised pretraining across domains

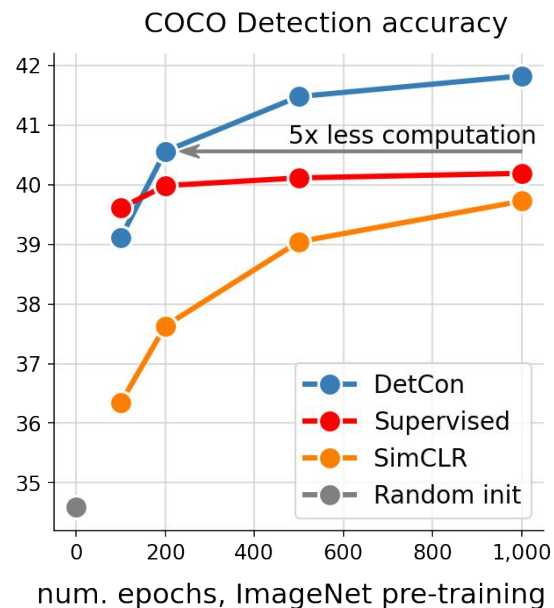


0. Sort  $E$  into  $\pi = (o_1, \dots, o_m)$ , by non-decreasing
1. Start with a segmentation  $S^0$ , where each vertex
2. Repeat step 3 for  $q = 1, \dots, m$ .
3. Construct  $S^q$  given  $S^{q-1}$  as follows. Let  $v_i$  and  $w_j$  be the  $q$ -th edge in the ordering, i.e.,  $o_q = (v_i, w_j)$ . If  $w_j$  is small compared to both those components, then merge the two components. More formally, let  $C_i^{q-1}$  be the component of  $S^{q-1}$  containing  $v_i$  and  $C_j^{q-1}$  be the component containing  $w_j$ . If  $C_i^{q-1} \neq C_j^{q-1}$  and  $w_j$  is small compared to both those components, then  $S^q$  is obtained from  $S^{q-1}$  by merging  $C_i^{q-1}$  and  $C_j^{q-1}$ .
4. Return  $S = S^m$ .

### Neural representations

- Enable intelligent behavior
- Require minimal supervision
- Are generally applicable

## Fine-tune for object detection, segmentation



# Can the self-supervised paradigm stay general?

## Self-supervised pretraining across modalities



*But thy eternal summer shall not fade,  
Nor lose possession of that fair thou ow'st,  
Nor shall death brag thou wander'st in his shade,  
When in eternal lines to time thou grow'st,  
So long as men can breathe, or eyes can see,  
So long lives this, and this gives life to thee.*

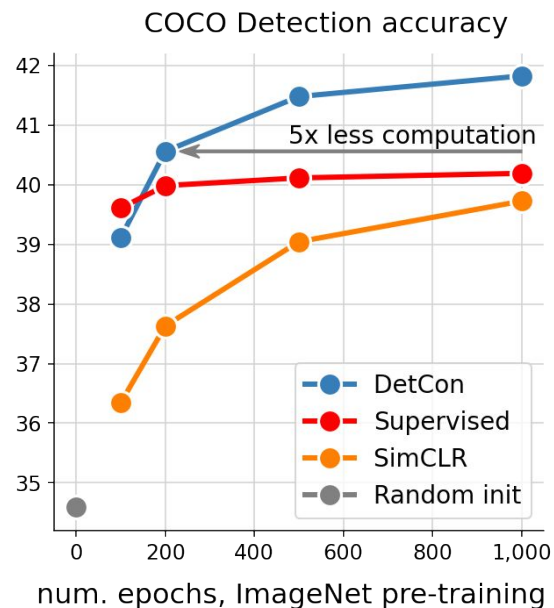


0. Sort  $E$  into  $\pi = (o_1, \dots, o_m)$ , by non-decreasing
1. Start with a segmentation  $S^0$ , where each vertex
2. Repeat step 3 for  $q = 1, \dots, m$ .
3. Construct  $S^q$  given  $S^{q-1}$  as follows. Let  $v_i$  and  $w_j$  be the  $q$ -th edge in the ordering, i.e.,  $o_q = (v_i, w_j)$ . If  $w_j$  is small compared to  $v_i$ , then merge the two components containing both those components, then merge the two components. More formally, let  $C_i^{q-1}$  be the component of  $S^{q-1}$  containing  $v_i$  and  $C_j^{q-1}$  be the component containing  $w_j$ . If  $C_i^{q-1} \neq C_j^{q-1}$  and  $w_j$  is small compared to  $v_i$ , then  $S^q$  is obtained from  $S^{q-1}$  by merging  $C_i^{q-1}$  and  $C_j^{q-1}$ .
4. Return  $S = S^m$ .

### Neural representations

- Enable intelligent behavior
- Require minimal supervision
- Are generally applicable

## Fine-tune for object detection, segmentation





DeepMind

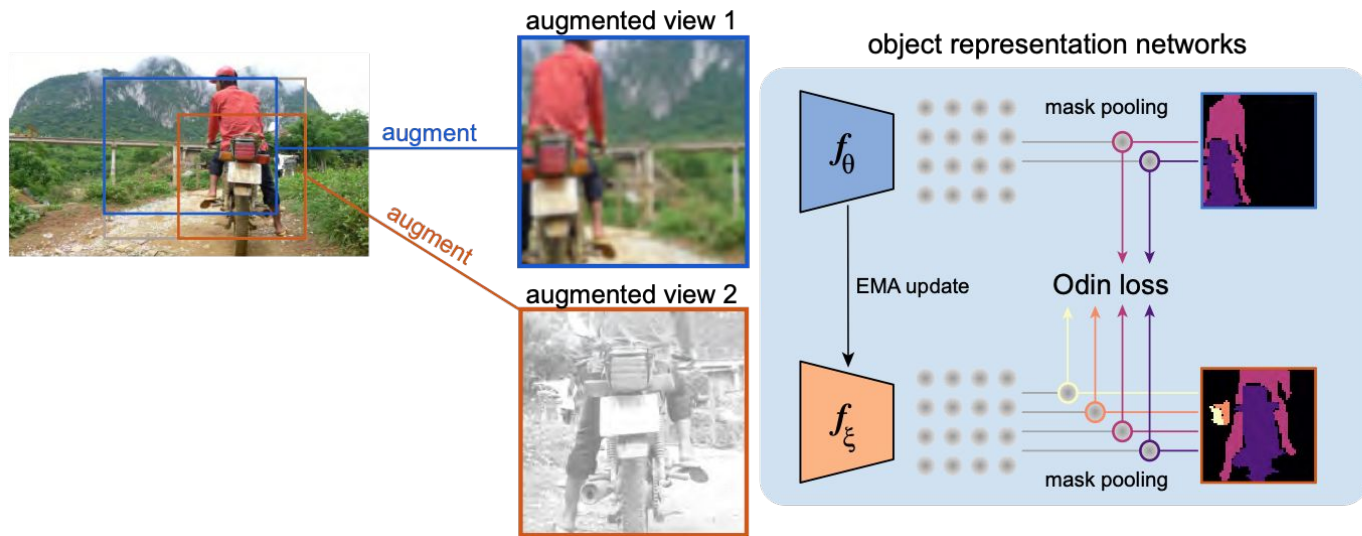
# Object discovery and representation networks

Olivier Hénaff, Skanda Koppula, Evan Shelhamer, Daniel Zoran,  
Drew Jaegle, Andrew Zisserman, João Carreira, Relja Arandjelović

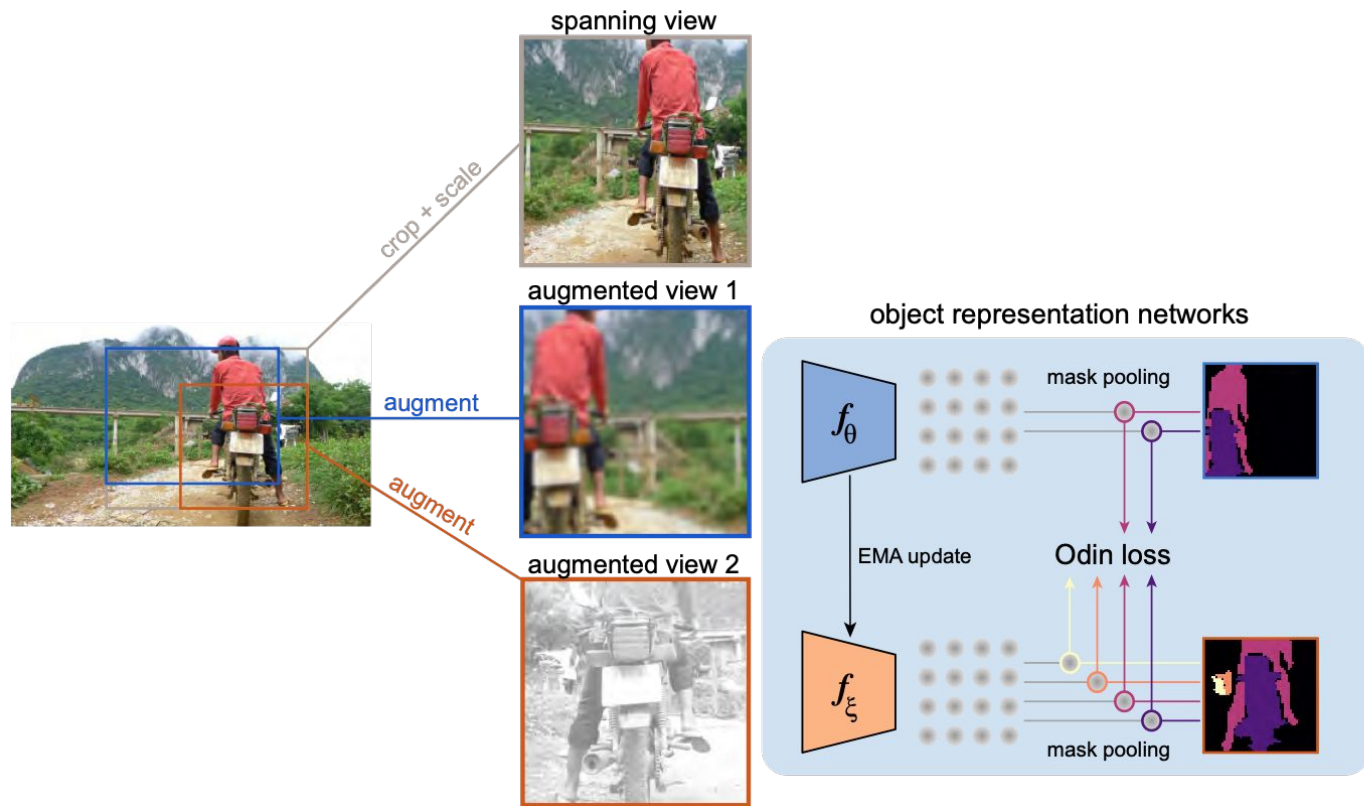
ECCV 2022



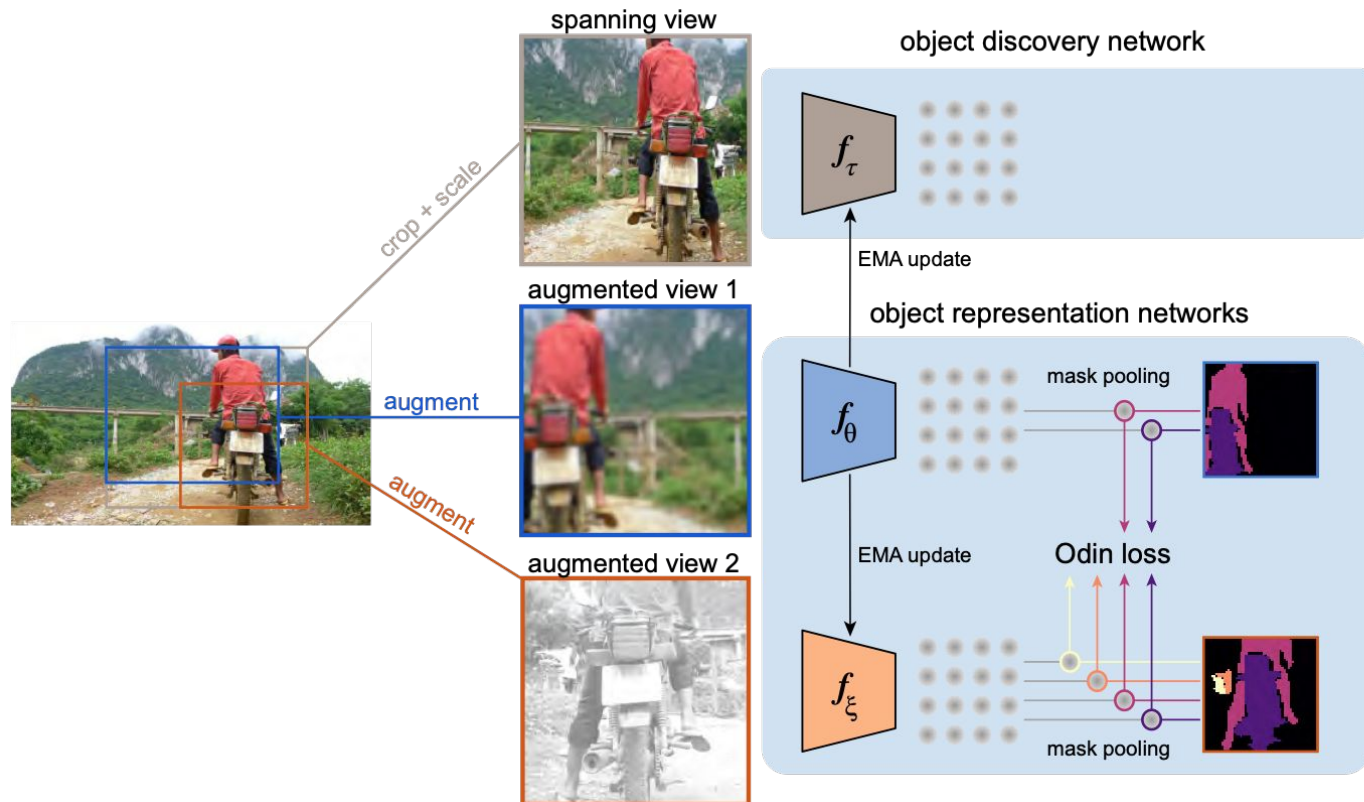
# Odin: Object discovery and representation networks



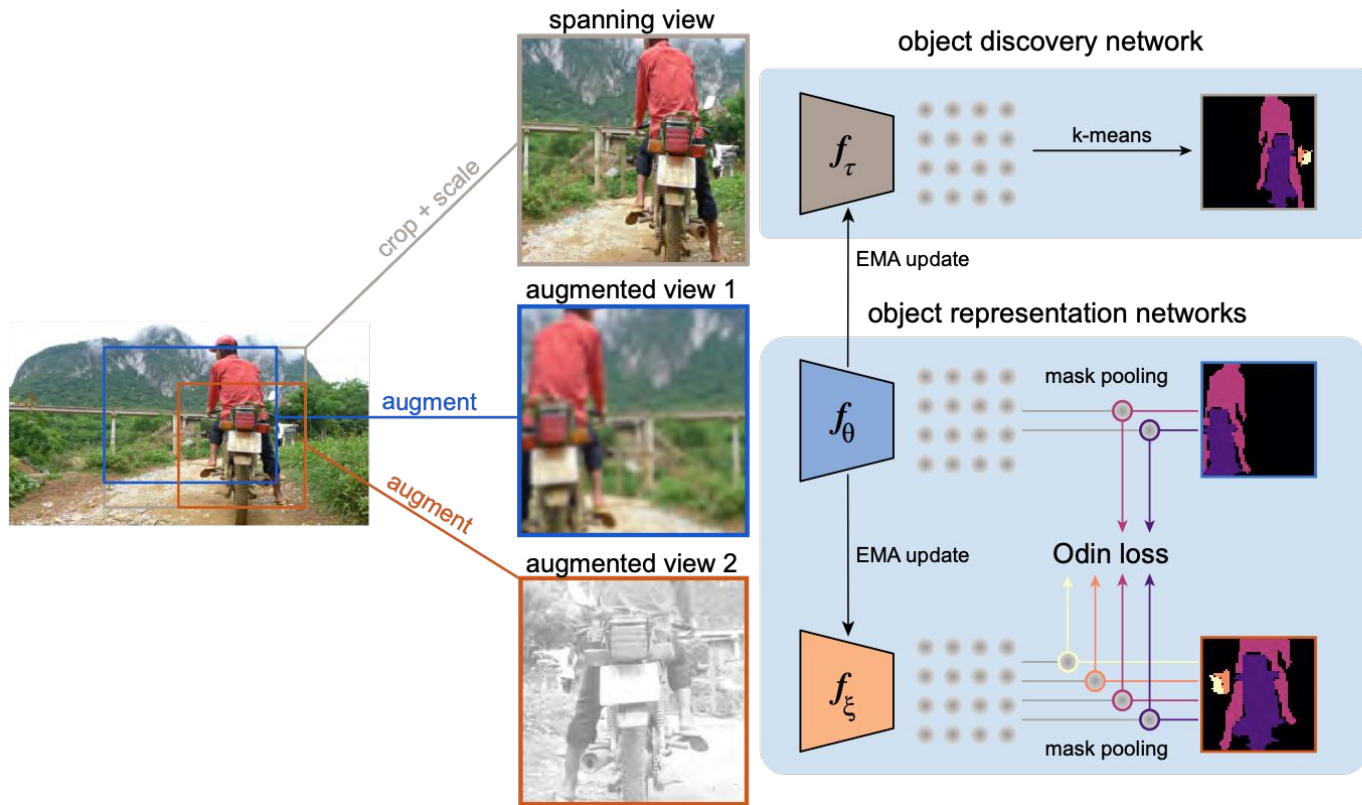
# Odin: Object discovery and representation networks



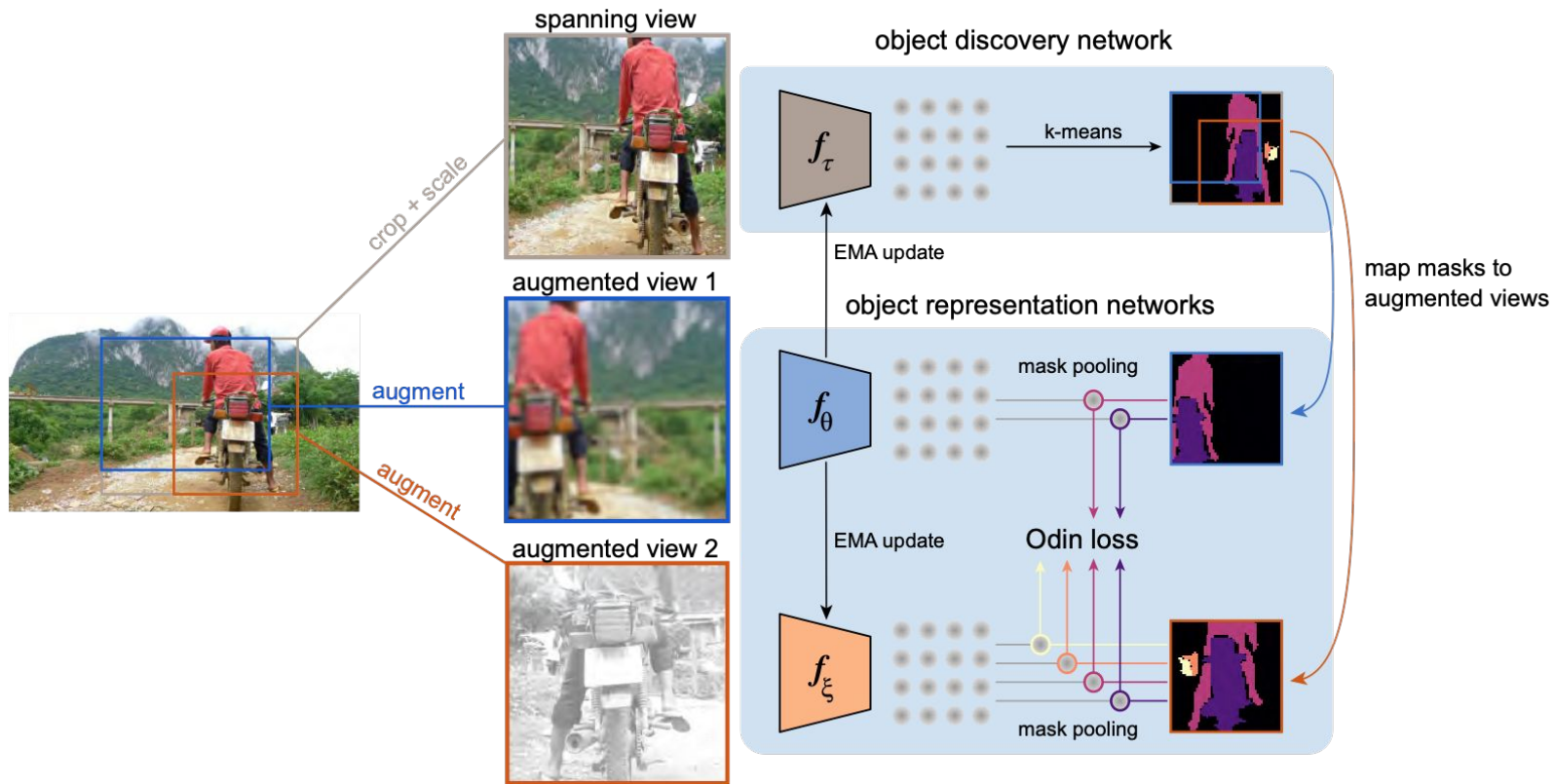
# Odin: Object discovery and representation networks



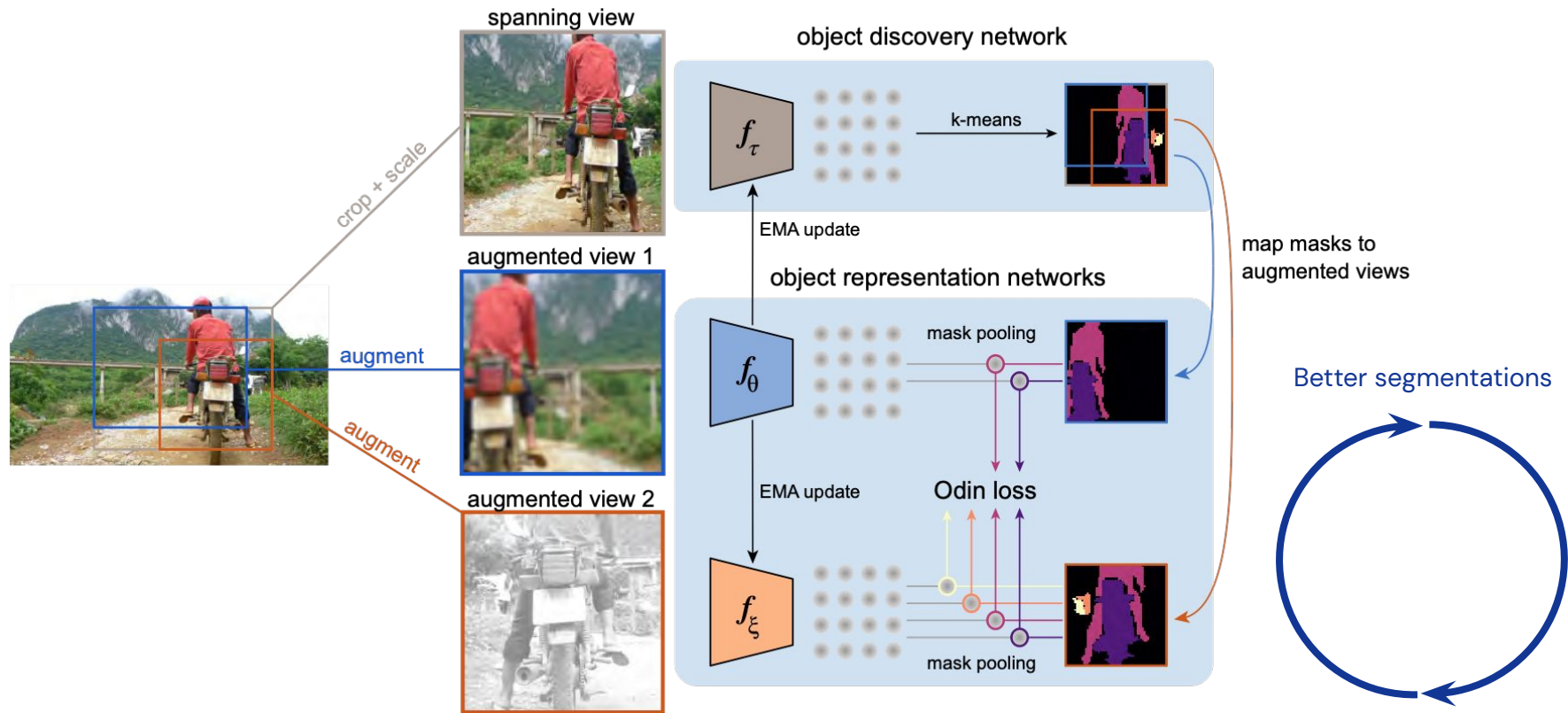
# Odin: Object discovery and representation networks



# Odin: Object discovery and representation networks



# Odin: Object discovery and representation networks



Better representations

# Transfer learning

- Pretrain ResNet-50 on ImageNet using self-supervised objective
- Fine-tune for COCO object detection and instance segmentation using Mask-RCNN

pretraining		fine-tune 1×		fine-tune 2×	
Method	Knows obj?	AP <sup>bb</sup>	AP <sup>mk</sup>	AP <sup>bb</sup>	AP <sup>mk</sup>
Supervised	no	39.6	35.6	41.6	37.6
VADeR [68]	no	39.2	35.6	-	-
MoCo [37]	no	39.4	35.6	41.7	37.5
SimCLR [13]	no	39.7	35.8	41.6	37.4
MoCo v2 [14]	no	40.1	36.3	41.7	37.6
InfoMin [75]	no	40.6	36.7	42.5	38.4
DINO [12]	no	41.2	37.1	42.3	38.1
PixPro [86]	no	41.4	-	-	-
BYOL [33]	no	41.6	37.2	42.4	38.0
SwAV [11]	no	41.6	37.8	-	-
ReLIC v2 [77]	<b>yes</b>	42.5	38.0	43.3	38.6
DetCon <sub>B</sub> [40]	<b>yes</b>	42.7	38.2	43.4	38.7





# Transfer learning

- Pretrain ResNet-50 on ImageNet using self-supervised objective
- Fine-tune for COCO object detection and instance segmentation using Mask-RCNN ✓

pretraining		fine-tune 1×		fine-tune 2×	
Method	Knows obj?	AP <sup>bb</sup>	AP <sup>mk</sup>	AP <sup>bb</sup>	AP <sup>mk</sup>
Supervised	no	39.6	35.6	41.6	37.6
VADeR [68]	no	39.2	35.6	-	-
MoCo [37]	no	39.4	35.6	41.7	37.5
SimCLR [13]	no	39.7	35.8	41.6	37.4
MoCo v2 [14]	no	40.1	36.3	41.7	37.6
InfoMin [75]	no	40.6	36.7	42.5	38.4
DINO [12]	no	41.2	37.1	42.3	38.1
PixPro [86]	no	41.4	-	-	-
BYOL [33]	no	41.6	37.2	42.4	38.0
SwAV [11]	no	41.6	37.8	-	-
ReLIC v2 [77]	yes	42.5	38.0	43.3	38.6
DetCon <sub>B</sub> [40]	yes	42.7	38.2	43.4	38.7
<b>Odin</b>	no	<b>42.9</b>	<b>38.4</b>	<b>43.8</b>	<b>39.1</b>



# Transfer learning

- Pretrain ResNet-50 on ImageNet using self-supervised objective
- Fine-tune for COCO object detection and instance segmentation using Mask-RCNN ✓
- Fine-tune for PASCAL and Cityscapes semantic segmentation ✓

Method	Knows obj?	PASCAL	Cityscapes
Supervised	no	74.4	74.9
BYOL [33]	no	75.7	74.6
DINO [12]	no	76.9	75.6
DetCon <sub>B</sub> [40]	yes	77.3	77.0
ReLIC v2 [77]	yes	77.9	75.2
<b>Odin</b>	no	<b>78.6</b>	<b>77.1</b>



# Transfer learning

- Pretrain ResNet-50 on ImageNet using self-supervised objective
- Fine-tune for COCO object detection and instance segmentation using Mask-RCNN ✓
- Fine-tune for PASCAL and Cityscapes semantic segmentation ✓
- Fine-tune for COCO object detection using FCOS\* ✓

Method	Knows obj?	PASCAL	Cityscapes
Supervised	no	74.4	74.9
BYOL [33]	no	75.7	74.6
DINO [12]	no	76.9	75.6
DetCon <sub>B</sub> [40]	yes	77.3	77.0
ReLIC v2 [77]	yes	77.9	75.2
<b>Odin</b>	no	<b>78.6</b>	<b>77.1</b>

Pretraining	Knows obj?	ResNet-50	Swin-T	Swin-S
Supervised	no	44.2	46.7	48.3
DINO [12]	no	44.3	-	-
MOBY [85]	no	-	47.6	-
DetCon <sub>B</sub> [40]	yes	45.4	48.4	<b>50.4</b>
<b>Odin</b>	no	<b>45.6</b>	<b>48.5</b>	<b>50.4</b>



# Object discovery



# Object discovery

- Pretrain ResNet-50 on ImageNet using self-supervised objective
- Evaluate on COCO, cluster features using k-means, report best overlap

k-means segmentation of CNN features

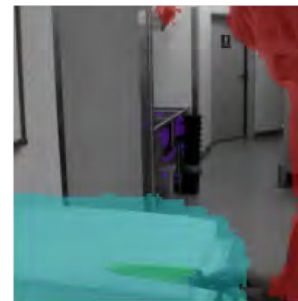
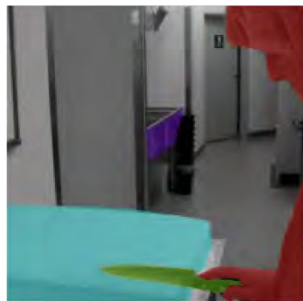
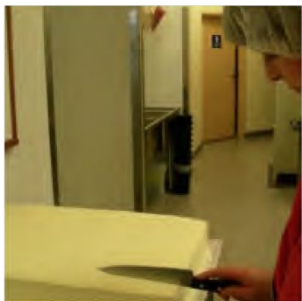
Original image



Human-annotated



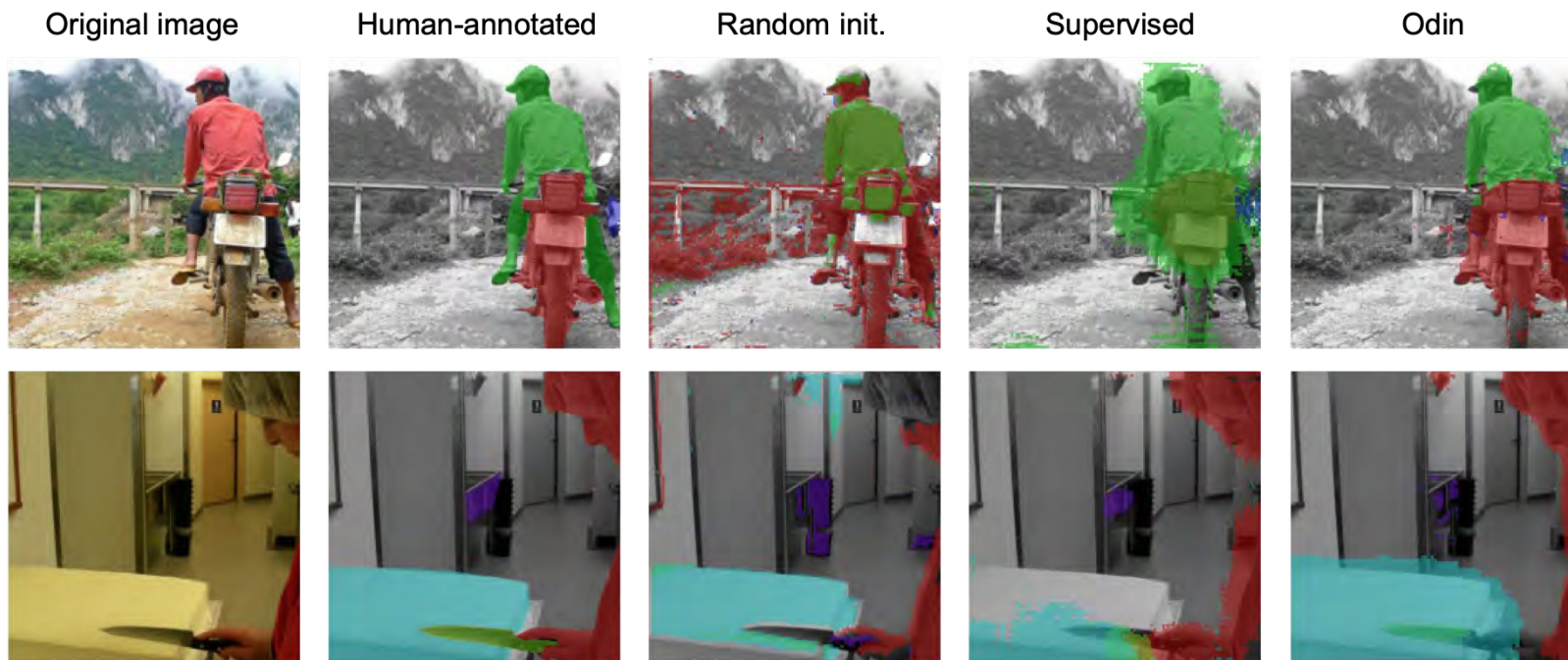
Odin



# Object discovery

- Pretrain ResNet-50 on ImageNet using self-supervised objective
- Evaluate on COCO, cluster features using k-means, report best overlap

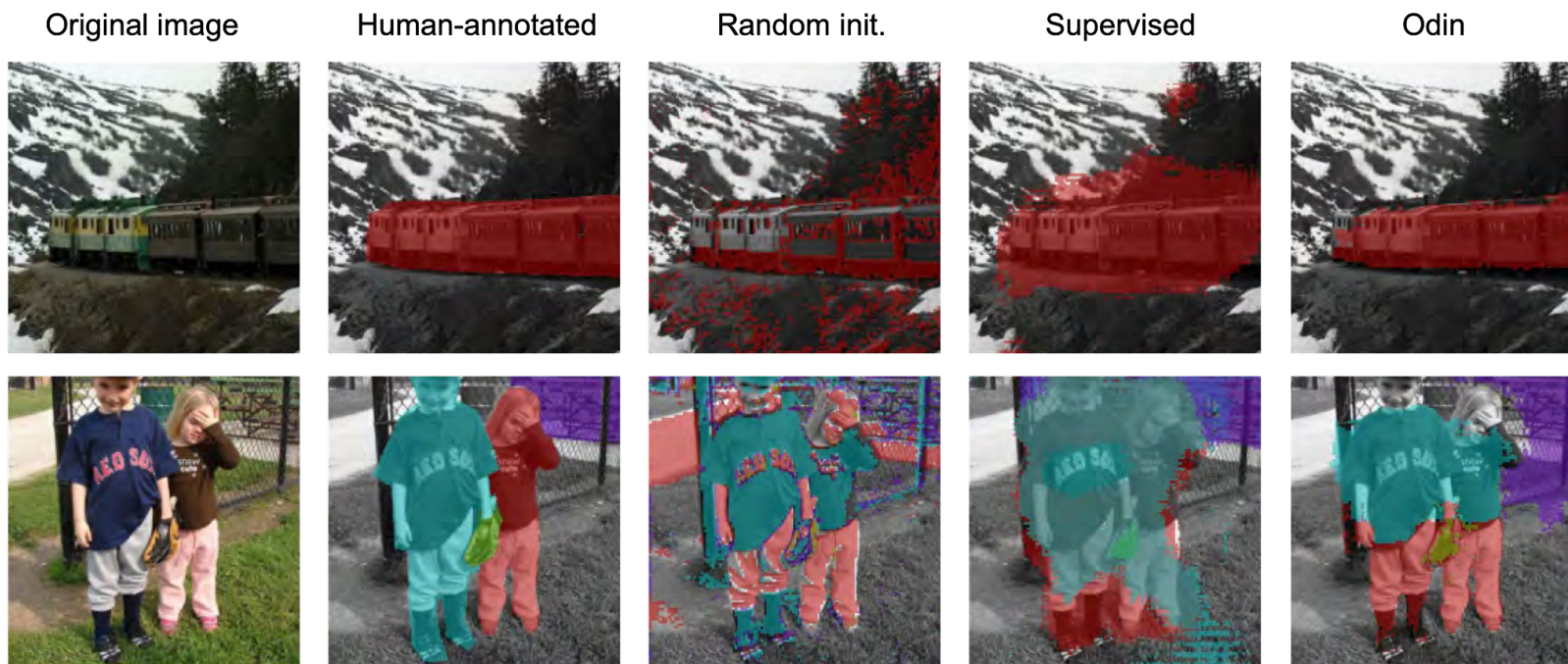
k-means segmentation of CNN features



# Object discovery

- Pretrain ResNet-50 on ImageNet using self-supervised objective
- Evaluate on COCO, cluster features using k-means, report best overlap

k-means segmentation of CNN features



# Object discovery

- Pretrain ResNet-50 on ImageNet using self-supervised objective
- Evaluate on COCO, cluster features using k-means, report best overlap

Pretraining	ResNet-50		
	ABO <sup>i</sup>	ABO <sup>c</sup>	OR
Random init	28.1	33.6	17.0
DetCon <sub>B</sub> [42]	34.1	40.0	20.4
Supervised	35.8	41.1	23.8
DINO [12]	38.3	46.5	30.8
<b>Odin</b>	<b>41.5</b>	<b>48.6</b>	<b>36.5</b>
<b>Odin<sup>†</sup></b>	<b>43.0</b>	<b>53.0</b>	<b>42.3</b>





# Object discovery

- Pretrain ViT-B/8 on ImageNet using self-supervised objective
- Evaluate on COCO, cluster features using k-means, report best overlap

Pretraining	ResNet-50			ViT-B		
	$ABO^i$	$ABO^c$	OR	$ABO^i$	$ABO^c$	OR
Random init	28.1	33.6	17.0	27.8	33.6	17.5
DetCon <sub>B</sub> [42]	34.1	40.0	20.4	-	-	-
Supervised	35.8	41.1	23.8	43.9	53.6	41.9
DINO [12]	38.3	46.5	30.8	42.7	51.7	39.7
<b>Odin</b>	<b>41.5</b>	<b>48.6</b>	<b>36.5</b>	<b>45.9</b>	<b>53.9</b>	<b>44.1</b>
<b>Odin<sup>†</sup></b>	<b>43.0</b>	<b>53.0</b>	<b>42.3</b>			



# Object discovery

- Pretrain ViT-B/8 on ImageNet using self-supervised objective
- Evaluate on COCO, cluster features using k-means, report best overlap

k-means segmentation of ViT features

Original image



Human-annotated



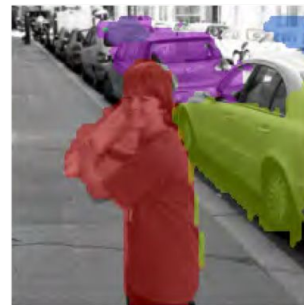
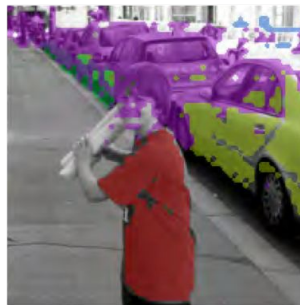
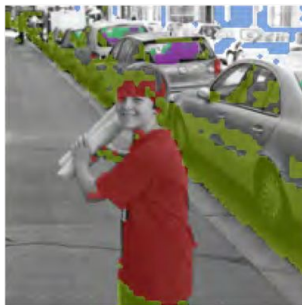
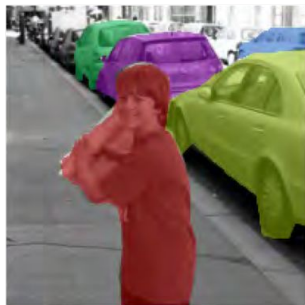
Random init.



DINO

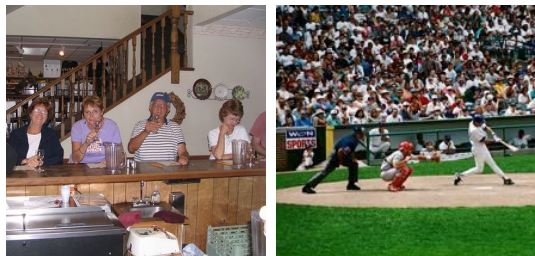


Odin

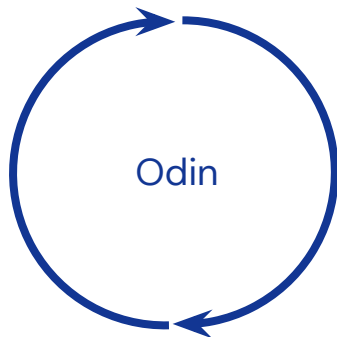


# Object knowledge can be extracted from learned representations

Self-supervised pretraining



Better segmentations



Better representations

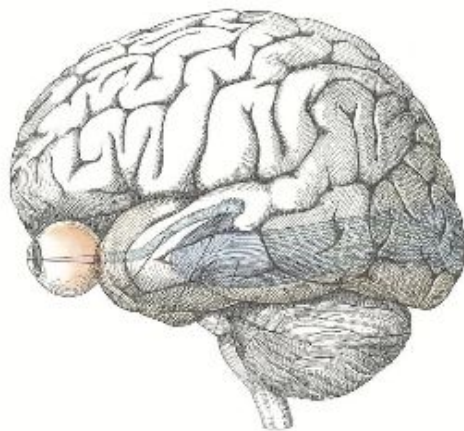
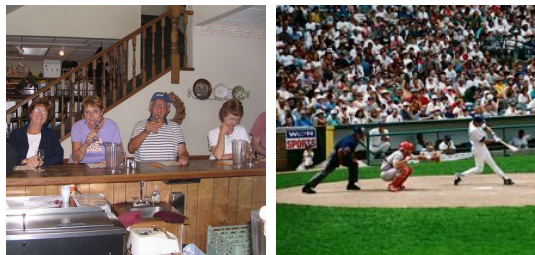
Fine-tune for object detection, segmentation



Method	pretraining	fine-tune 1×	
	Knows obj?	AP <sup>bb</sup>	AP <sup>mk</sup>
Supervised	no	39.6	35.6
VADeR [68]	no	39.2	35.6
MoCo [37]	no	39.4	35.6
SimCLR [13]	no	39.7	35.8
MoCo v2 [14]	no	40.1	36.3
InfoMin [75]	no	40.6	36.7
DINO [12]	no	41.2	37.1
ReLIC v2 [77]	<b>yes</b>	42.5	38.0
DetCon <sub>B</sub> [40]	<b>yes</b>	42.7	38.2
<b>Odin</b>	no	<b>42.9</b>	<b>38.4</b>

# Can SSL leverage videos to learn good image representations?

## Self-supervised pretraining

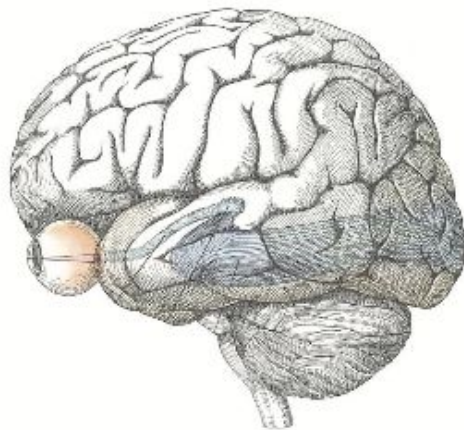
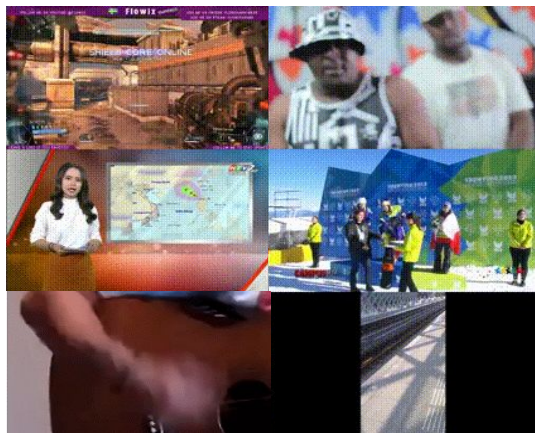
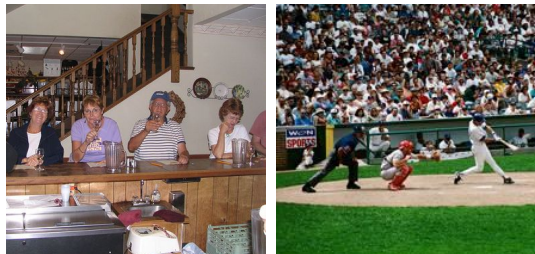


Natural videos provide

- Rich image augmentations
- Strong segmentation cues from motion

# Can SSL leverage videos to learn good image representations?

## Self-supervised pretraining



Natural videos provide

- Rich image augmentations
- Strong segmentation cues from motion

## Fine-tune for object detection, segmentation



Yet they have yet to yield strong image representations!  
(as evaluated by scene understanding tasks)

DeepMind

# Self-supervised video pretraining yields strong image representations

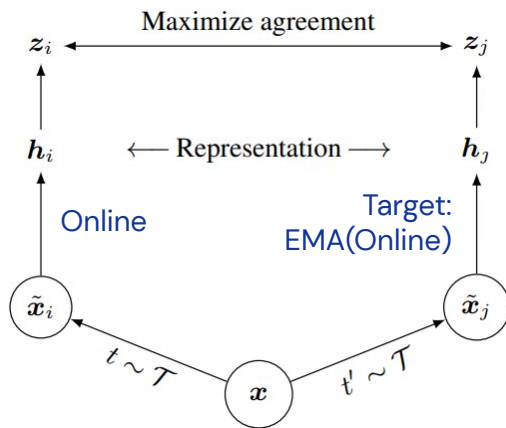
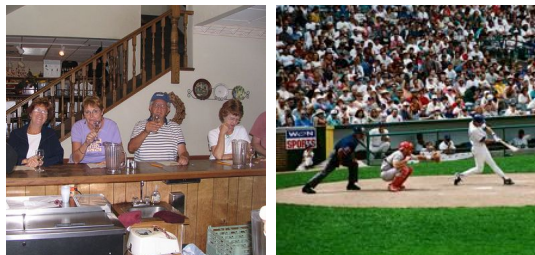
Nikhil Parthasarathy, Ali Eslami, João Carreira, Olivier Hénaff

arXiv 2022



# Can SSL leverage videos to learn good image representations?

Self-supervised pretraining on images or video



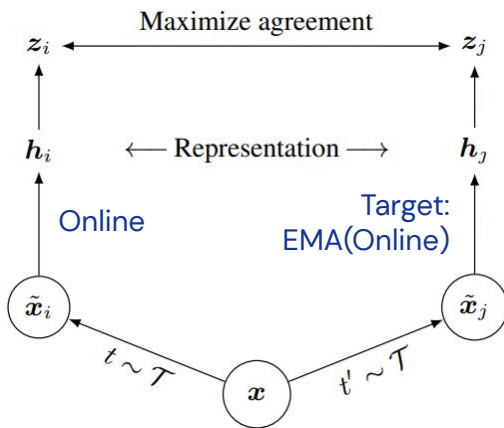
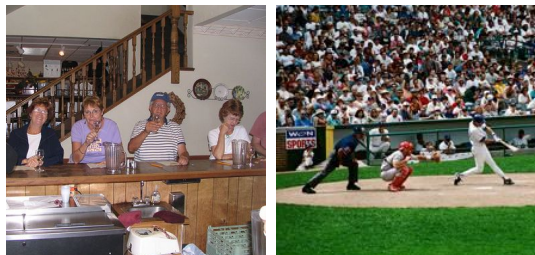
A strong contrastive baseline

- MoCLR: best of SimCLR, MoCo, and BYOL (Tian, 2021)
- Frame  $x$  sampled from image datasets or videos



# Can SSL leverage videos to learn good image representations?

Self-supervised pretraining on images or video



A strong contrastive baseline

- MoCLR: best of SimCLR, MoCo, and BYOL (Tian, 2021)
- Frame  $x$  sampled from image datasets or videos

Fine-tune for object detection, segmentation



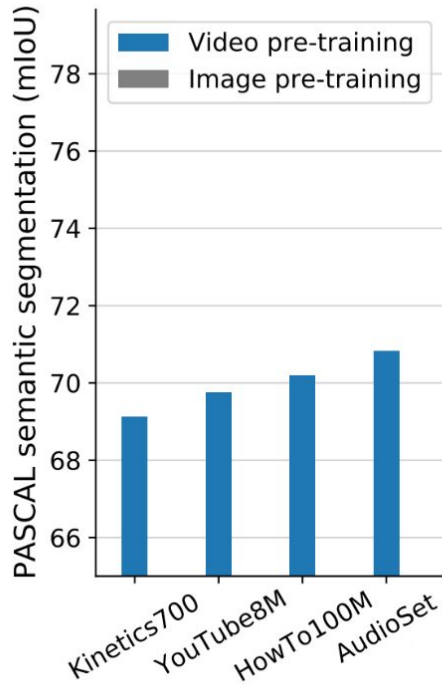
Fine-tune on

- Semantic segmentation (PASCAL or ADE20K)
- Object detection (COCO or LVIS)

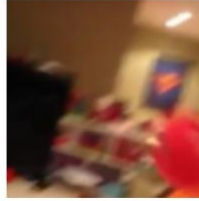


# Dataset curation matters

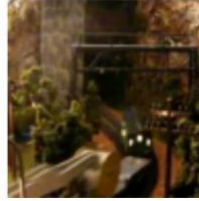
Treating video datasets as collections of independent frames



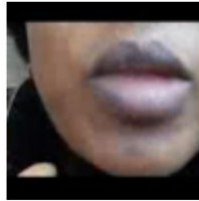
Kinetics700



Audioset

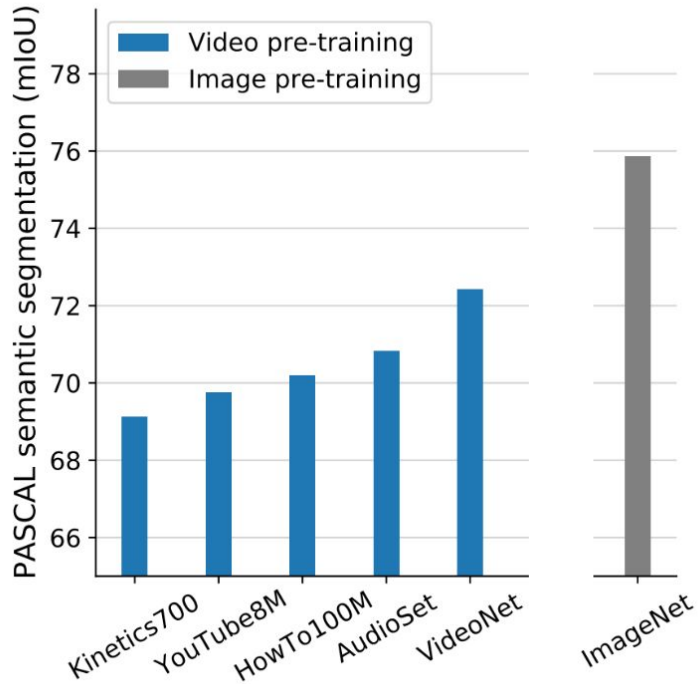


ImageNet

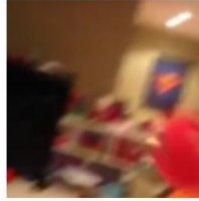


# Dataset curation matters

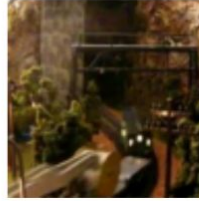
Treating video datasets as collections of independent frames



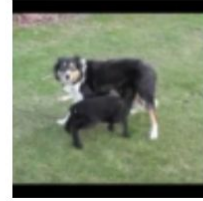
Kinetics700



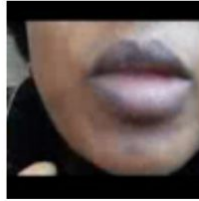
Audioset



VideoNet



ImageNet



# VideoNet data pipeline

Procedure for filtering uncurated video datasets

- 1) Query internet for ImageNet categories (**~5 million videos**)
- 2) Filter out videos less than 10s long
- 3) Run an ImageNet classifier on the first 100 frames of each video to verify they contain stated category (**~1.2 million videos**)

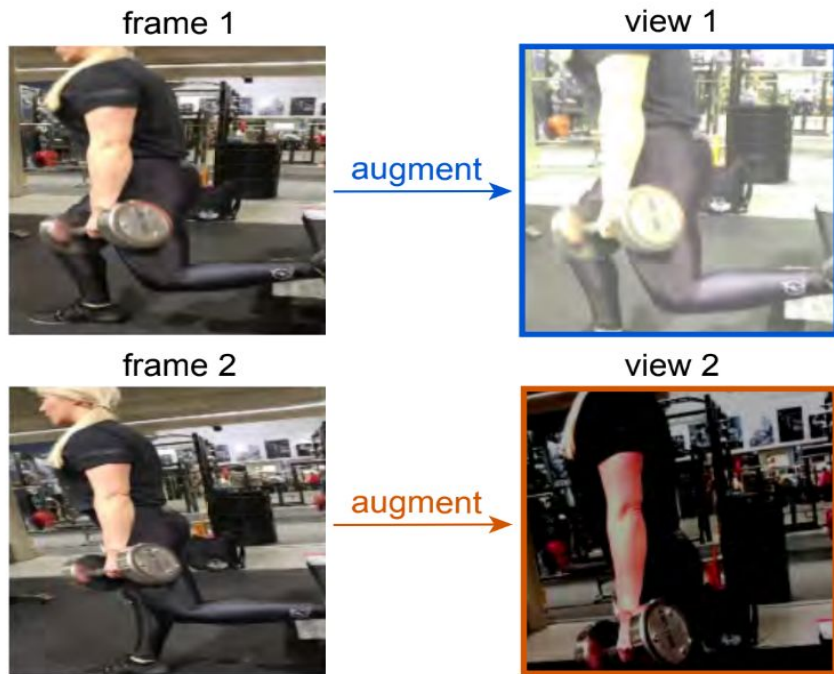


# VITO: a better video-to-image model

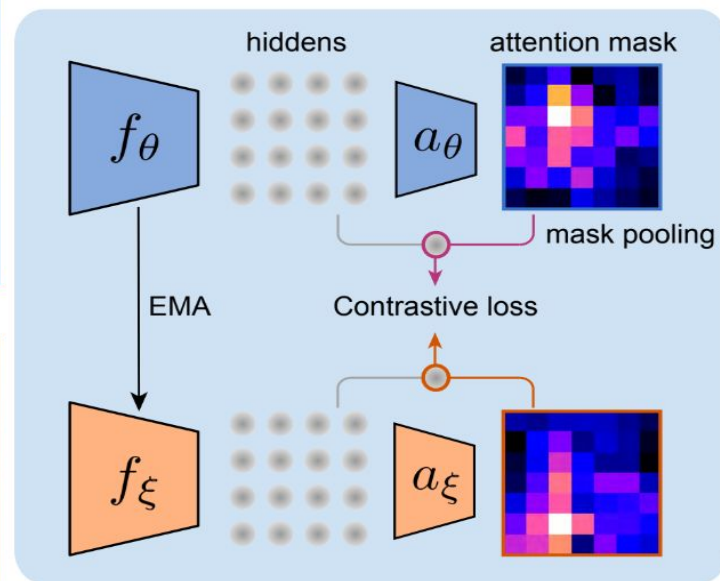
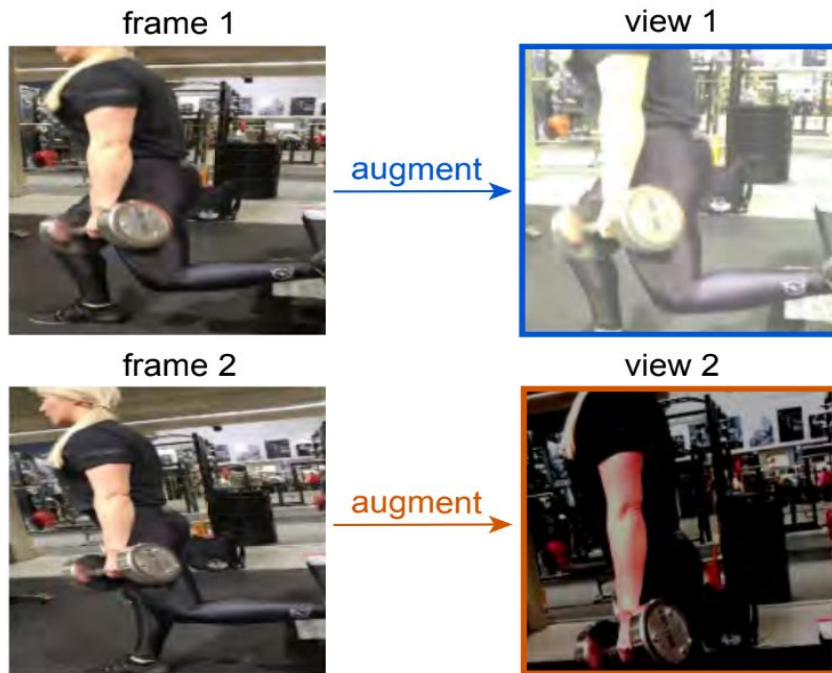
Private & Confidential



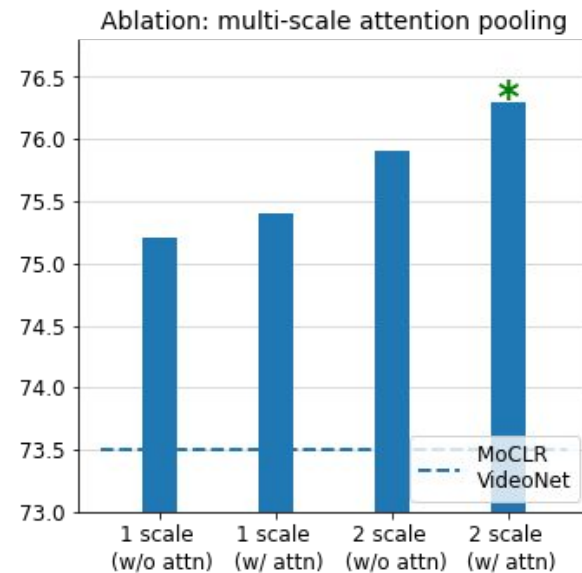
# VITO: a better video-to-image model



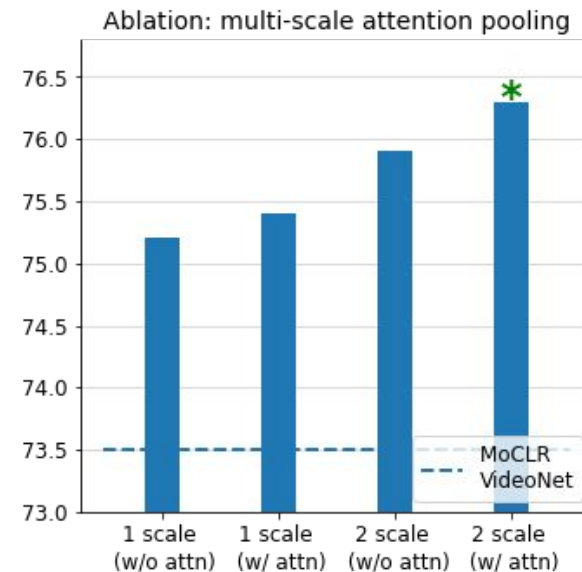
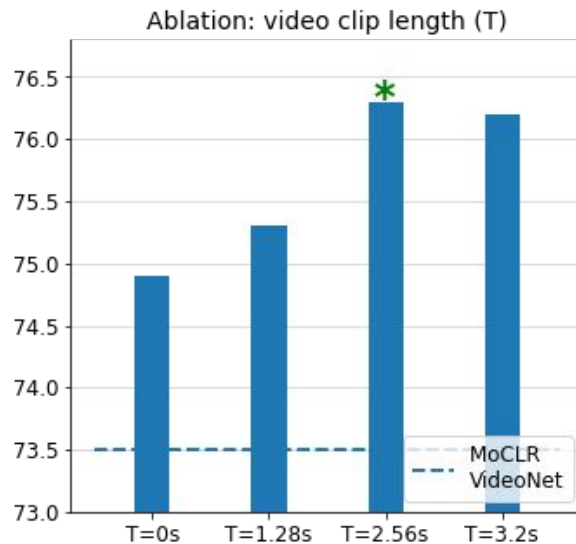
# VITO: a better video-to-image model



# Results: ablating the components of VITO

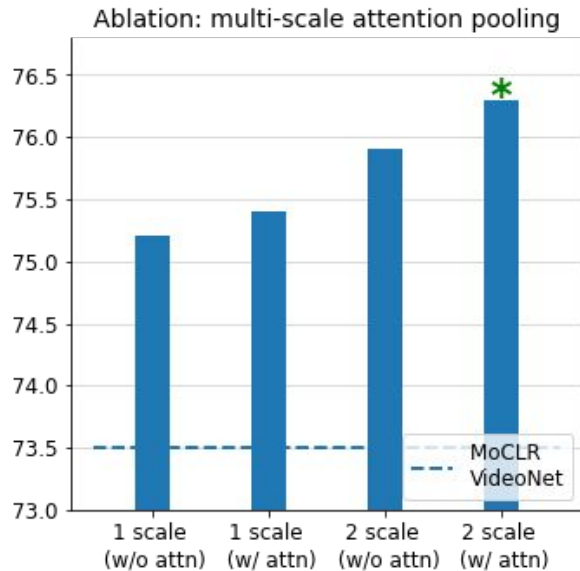
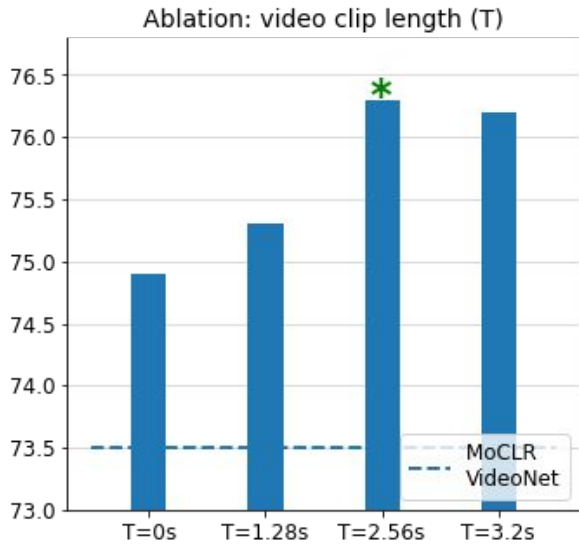
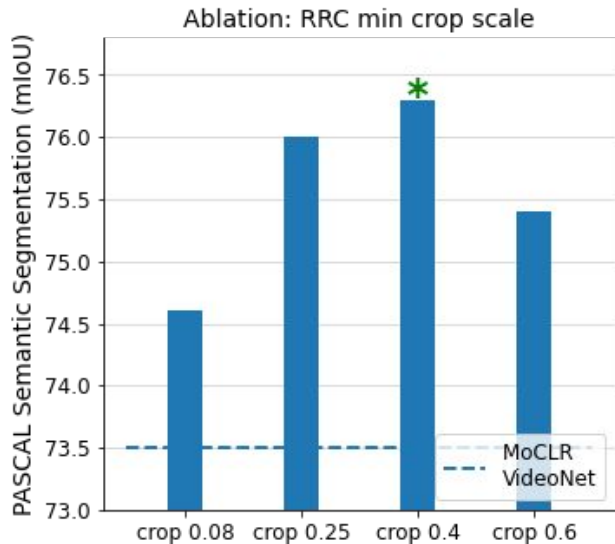


# Results: ablating the components of VITO

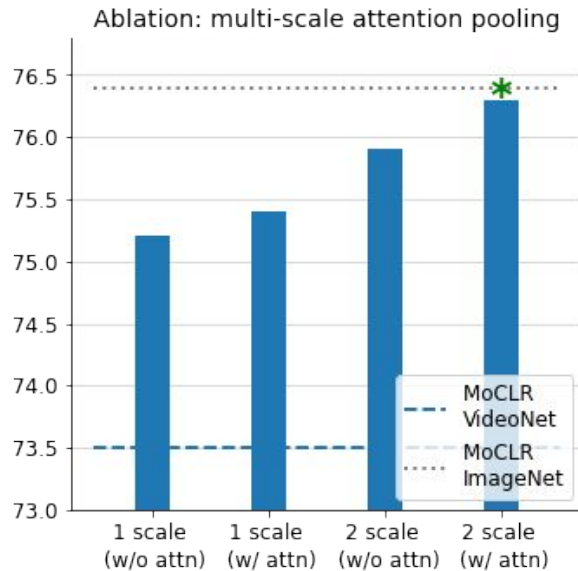
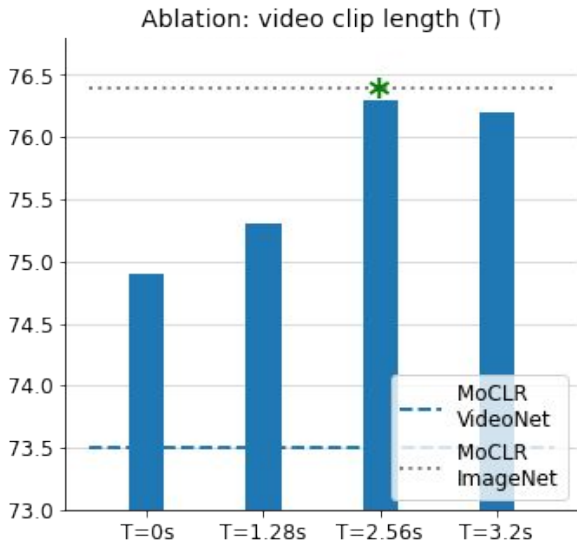
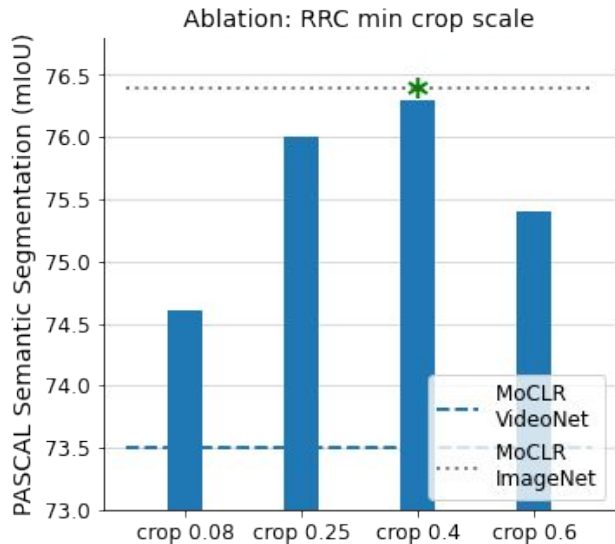




# Results: ablating the components of VITO



# Results: ablating the components of VITO



*VITO closes the gap with ImageNet MoCLR*



# Results: what does the attention pooling learn?

Due to the contrastive loss, VITO must learn to attend to content that is:

- **Stable (or predictable)** over time
- **Unique or discriminative** relative to content from other videos



# Results: what does the attention pooling learn?

Due to the contrastive loss, VITO must learn to attend to content that is:

- **Stable (or predictable)** over time
- **Unique or discriminative** relative to content from other videos

This also enables **semantic binding of content**

View 1



View 2



# Results: what does the attention pooling learn?



# Results: comparison to prior video-to-image methods

Pretraining	Dataset	Epochs	Semantic segmentation		Object detection	
			PASCAL	ADE20K	COCO	LVIS
Random Init			53.0	27.9	39.0	21.1
<i>Methods pretraining on video datasets</i>						
VFS (Xu & Wang, 2021)	K400	100	63.9	31.4	41.6	23.2
VIVI (Tschannen et al., 2020)	YouTube8M	192	65.8	34.2	41.3	23.2
VINCE (Gordon et al., 2020)	R2V2	200	69.0	35.7	42.4	24.4
CycleContrast (Wu & Wang, 2021)	R2V2	200	69.2	35.6	42.8	24.5
<b>VITO</b>	<b>VideoNet</b>	<b>200</b>	<b>75.5</b>	<b>39.2</b>	<b>43.6</b>	<b>25.6</b>

VITO vs. other frame-level SSL objectives highlights importance of VITO components

→ data curation, attention pooling, random-crop scale



# Results: comparison to prior video-to-image methods

Pretraining	Dataset	Epochs	Semantic segmentation		Object detection	
			PASCAL	ADE20K	COCO	LVIS
Random Init			53.0	27.9	39.0	21.1
<i>Methods pretraining on video datasets</i>						
VFS (Xu & Wang, 2021)	K400	100	63.9	31.4	41.6	23.2
VIVI (Tschannen et al., 2020)	YouTube8M	192	65.8	34.2	41.3	23.2
VINCE (Gordon et al., 2020)	R2V2	200	69.0	35.7	42.4	24.4
CycleContrast (Wu & Wang, 2021)	R2V2	200	69.2	35.6	42.8	24.5
MMV (Alayrac et al., 2020)	AS + HT	1600	70.6	32.5	41.3	24.2
<b>VITO</b>	<b>VideoNet</b>	<b>200</b>	<b>75.5</b>	<b>39.2</b>	<b>43.6</b>	<b>25.6</b>

VITO vs. other frame-level SSL objectives highlights importance of VITO components

→ data curation, attention pooling, random-crop scale

VITO also outperforms recent multimodal methods



# Results: comparison to prior video-to-image methods

Pretraining	Dataset	Epochs	Semantic segmentation		Object detection	
			PASCAL	ADE20K	COCO	LVIS
Random Init			53.0	27.9	39.0	21.1
<i>Methods pretraining on video datasets</i>						
VFS (Xu & Wang, 2021)	K400	100	63.9	31.4	41.6	23.2
VIVI (Tschannen et al., 2020)	YouTube8M	192	65.8	34.2	41.3	23.2
VINCE (Gordon et al., 2020)	R2V2	200	69.0	35.7	42.4	24.4
CycleContrast (Wu & Wang, 2021)	R2V2	200	69.2	35.6	42.8	24.5
MMV (Alayrac et al., 2020)	AS + HT	1600	70.6	32.5	41.3	24.2
VITO	VideoNet	200	<b>75.5</b>	<b>39.2</b>	<b>43.6</b>	<b>25.6</b>
VITO	AudioSet	300	73.6	38.5	43.2	25.0
VITO	VideoNet	300	<b>76.3</b>	<b>39.4</b>	<b>44.0</b>	<b>25.7</b>

VITO also outperforms prior art with existing video datasets (Audioset),  
but VideoNet provides further gains





# Results: VITO closes the gap with ImageNet pretraining

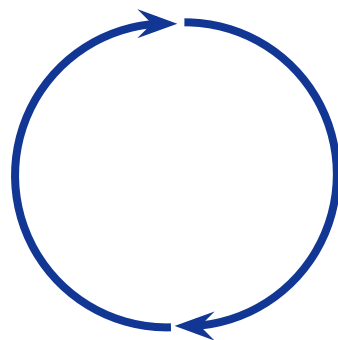
Pretraining	Dataset	Epochs	Semantic segmentation		Object detection	
			PASCAL	ADE20K	COCO	LVIS
Random Init			53.0	27.9	39.0	21.1
<i>Methods pretraining on video datasets</i>						
VFS (Xu & Wang, 2021)	K400	100	63.9	31.4	41.6	23.2
VIVI (Tschannen et al., 2020)	YouTube8M	192	65.8	34.2	41.3	23.2
VINCE (Gordon et al., 2020)	R2V2	200	69.0	35.7	42.4	24.4
CycleContrast (Wu & Wang, 2021)	R2V2	200	69.2	35.6	42.8	24.5
MMV (Alayrac et al., 2020)	AS + HT	1600	70.6	32.5	41.3	24.2
<b>VITO</b>	<b>VideoNet</b>	<b>300</b>	<b>76.3</b>	<b>39.4</b>	<b>44.0</b>	<b>25.7</b>
<i>Methods pretraining on ImageNet</i>						
Supervised	ImageNet	200	71.3	33.5	44.2	25.2
BYOL (Grill et al., 2020)	ImageNet	300	76.1	38.8	43.7	25.5
MoCLR (Tian et al., 2021)	ImageNet	300	76.4	39.2	43.9	25.8
DINO (Caron et al., 2021)	ImageNet	300	76.1	39.0	44.3	26.4



# Conclusion

1. **Knowledge of objects accelerates and improves representation learning**  
→ DetCon objective (ICCV, 2021)
2. **Knowledge of objects can be extracted from learned representations**  
→ Odin framework (ECCV, 2022)
3. **Videos can be used to learn strong image representations**  
→ VITO framework (arXiv, 2022)

Better object knowledge



Better representations

DeepMind

**Thanks!**

