

# Multi-view invariance and grouping for self-supervised learning

Ishan Misra

Facebook AI Research

# Multi-view

Same data sample  
Different ways of looking at it



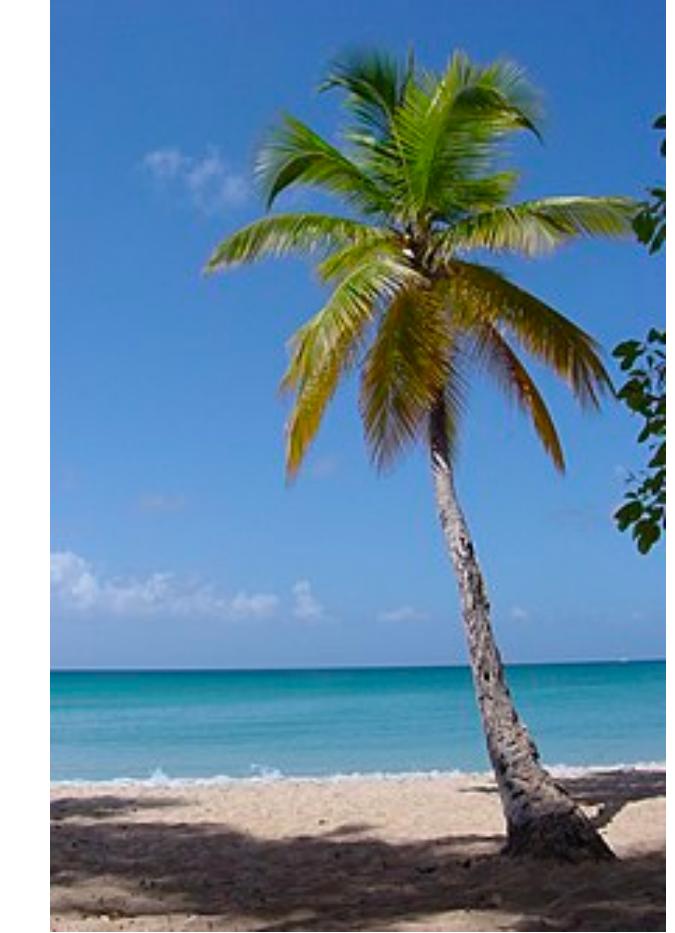
# Grouping

Associating related  
data samples



# Multi-view invariance

# Grouping



Let us fill this table throughout the talk

Method	Multi-view Invariance	Grouping
Method Name	?	?

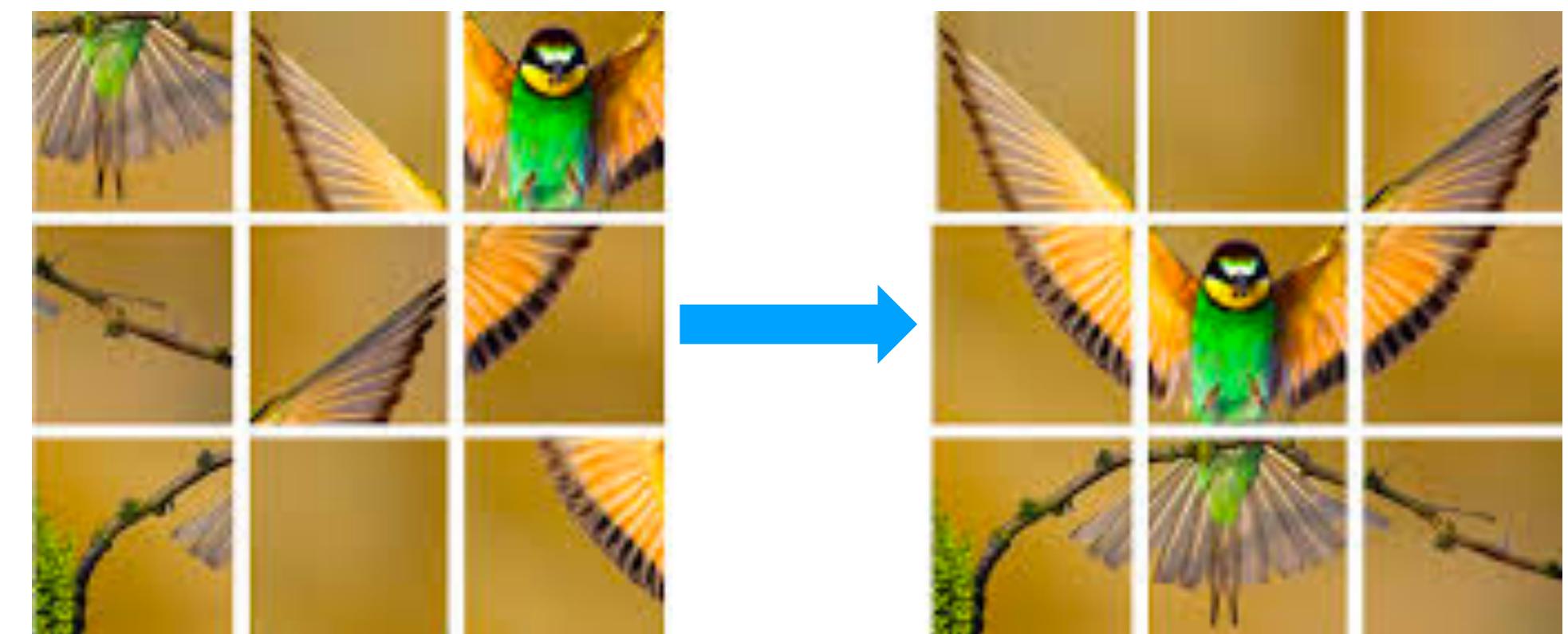
# Pre-2019: "Pretext" tasks

- Create proxy tasks



Rotation

(Gidaris et al., 2018)

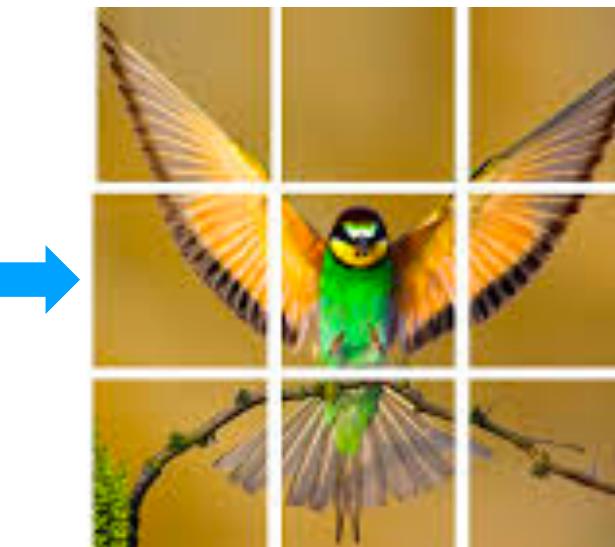
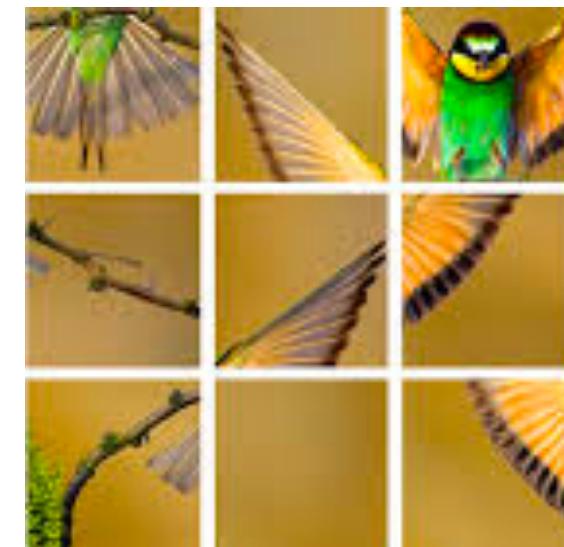


Jigsaw puzzles

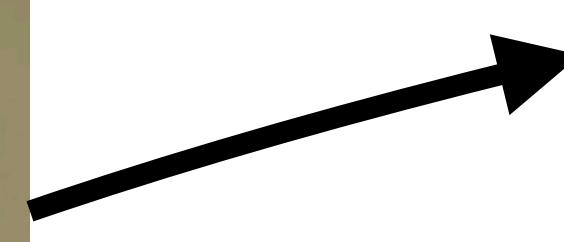
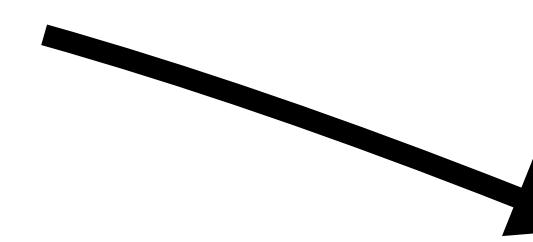
(Noroozi et al., 2016)

# The hope of generalization

- We hope that the pre-training task and the transfer task are "aligned"



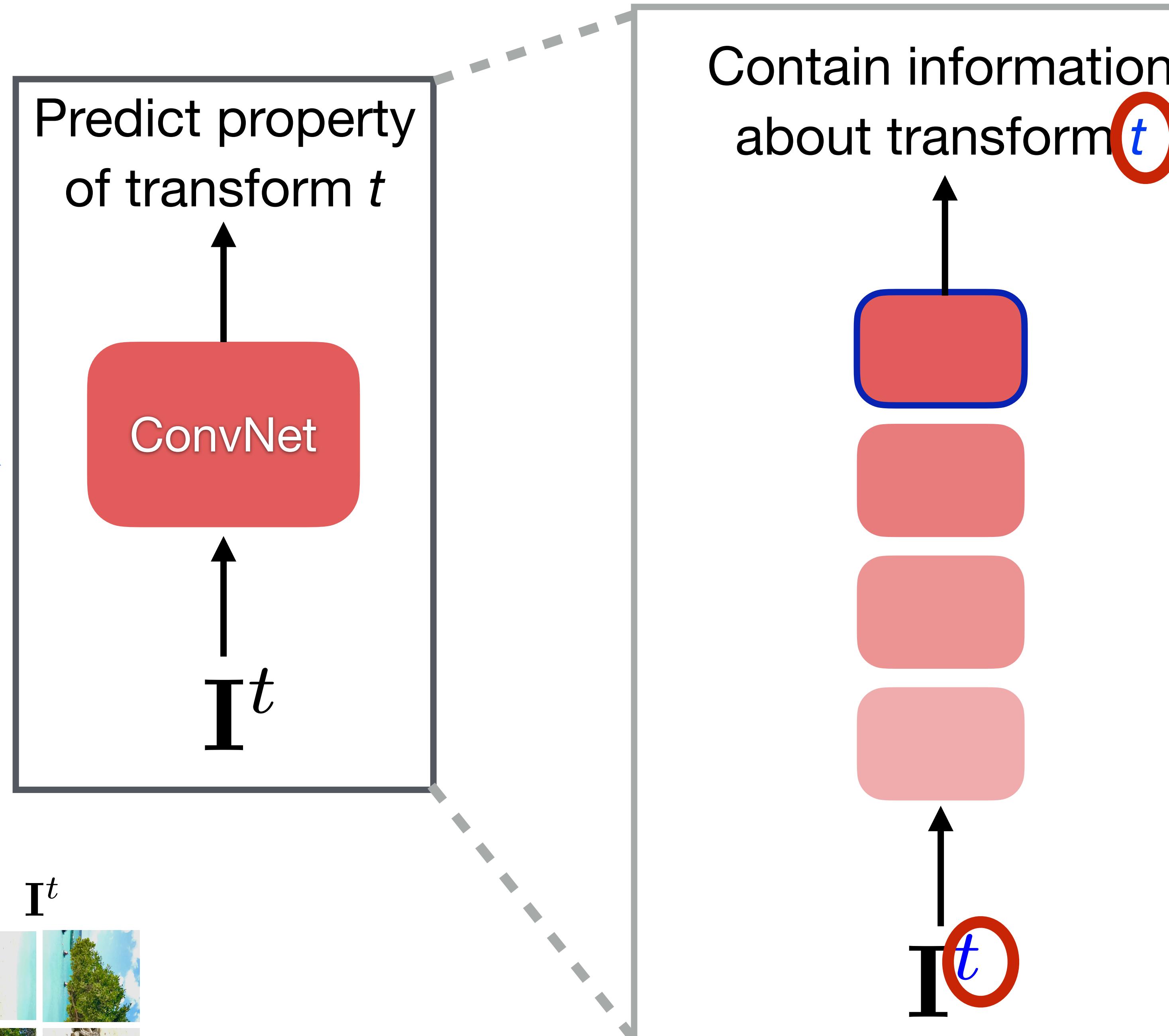
Pre-training



Transfer  
Tasks

# Less Semantic Features

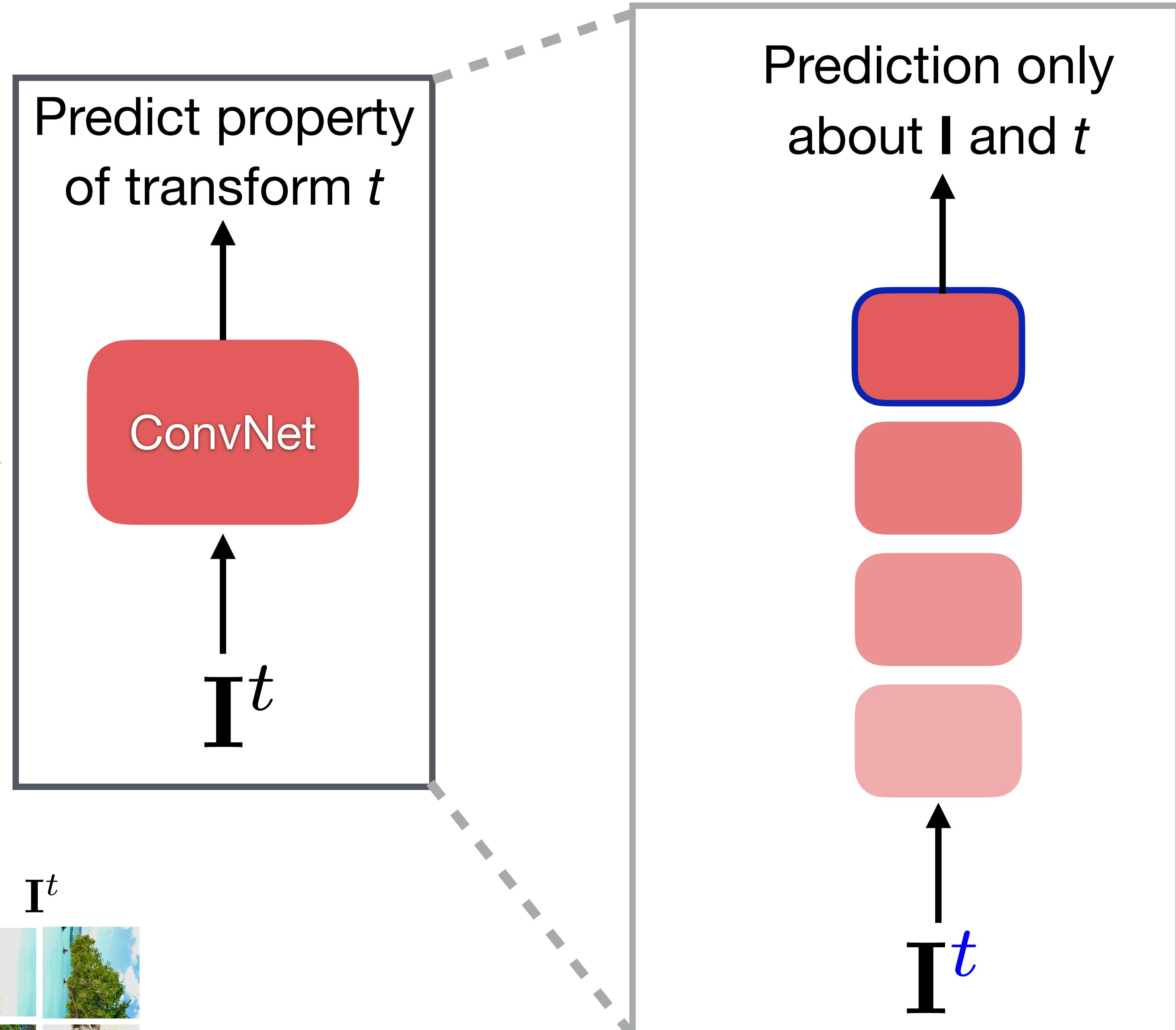
Pretext task



Method	Multi-view Invariance	Grouping
Pretext Task	No	?

No relation  
to other  
images

## Pretext task



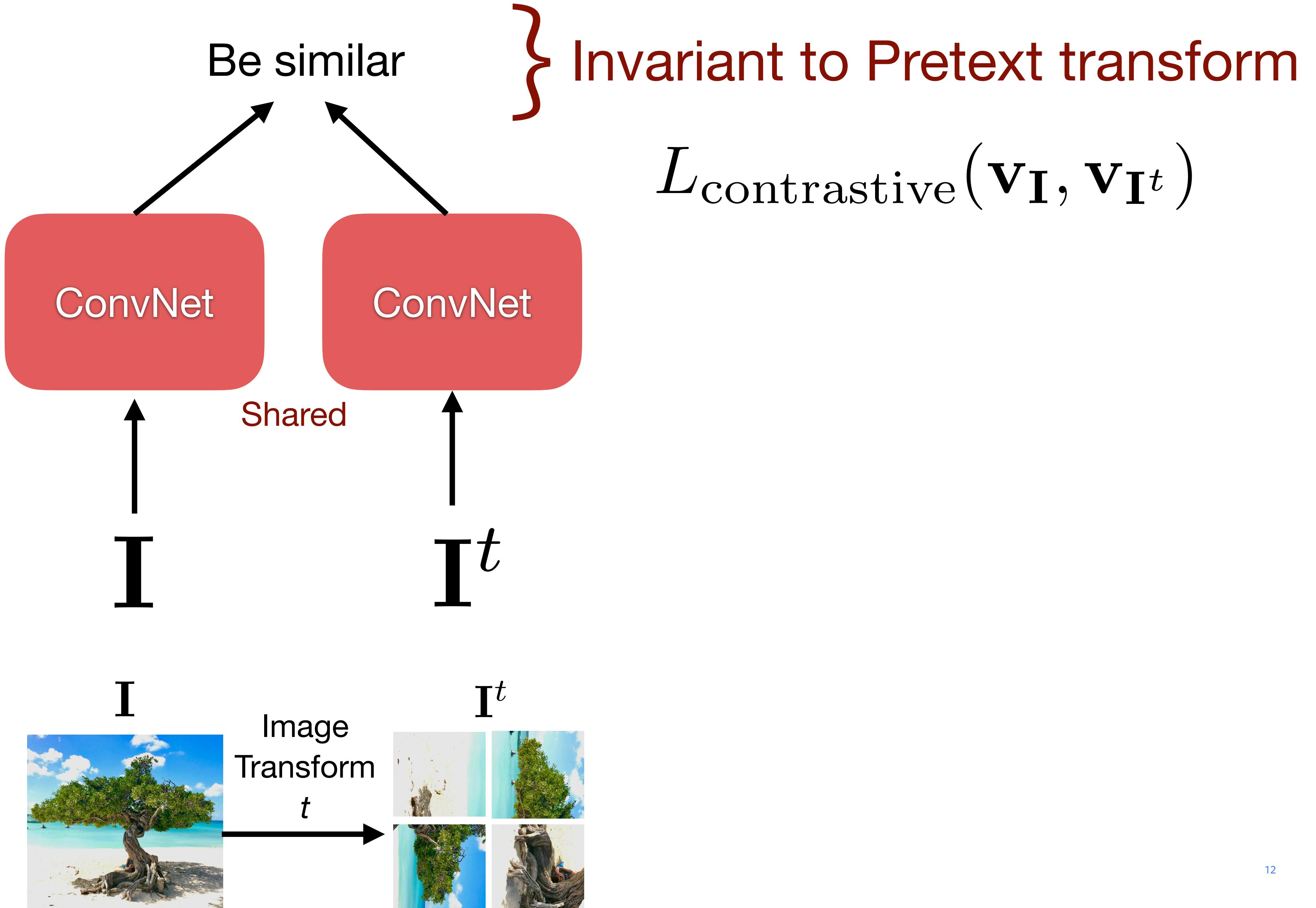
Method	Multi-view Invariance	Grouping	Performance
Pretext Task	No	No	Weak

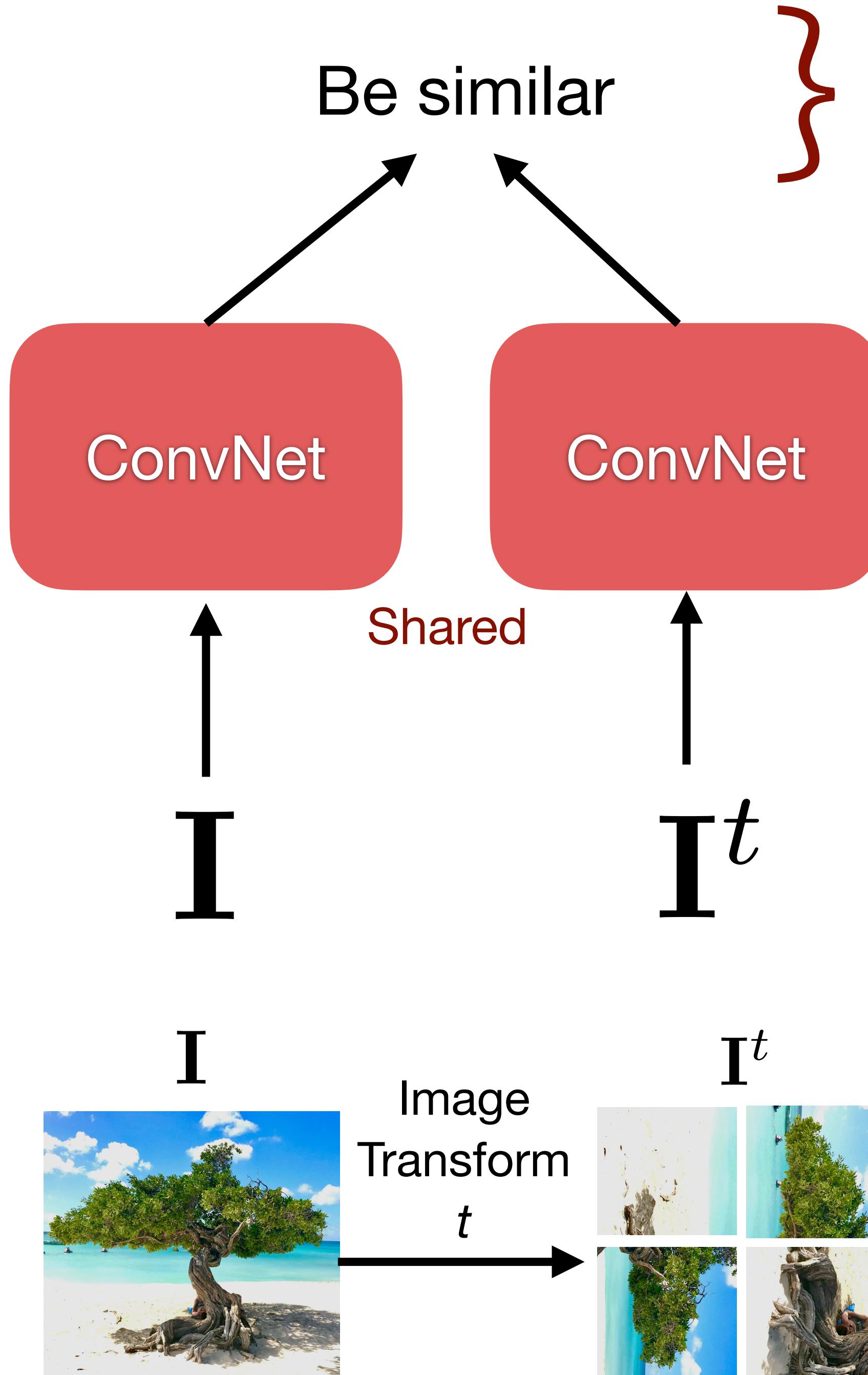
*"Self-supervised learning is a naive/optimistic hope for feature learning"*

# Pretext-Invariant Representation Learning (PIRL)

Ishan Misra, Laurens van der Maaten





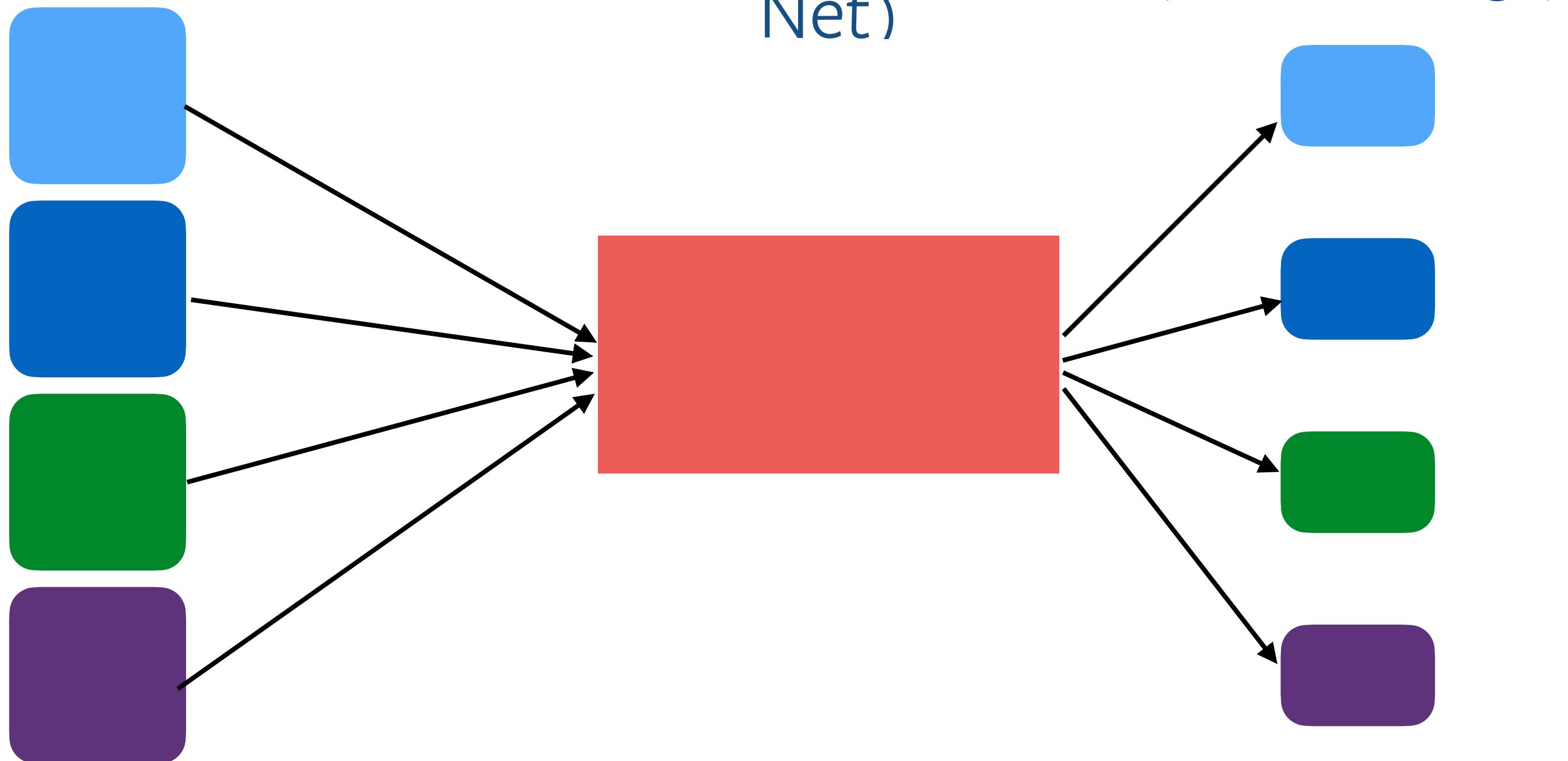


Invariance to

- Data Augmentations
- Multiple views created by pretext task (Jigsaw/Rotation)

# Contrastive Learning

Related and  
Unrelated  
Images



Shared  
network  
(Siamese  
Net)

Image  
Features  
(Embeddings)

## Loss Function

Embeddings from related images should be closer than embeddings from unrelated images

$$d(\text{blue}, \text{blue}) < d(\text{blue}, \text{green})$$
$$d(\text{blue}, \text{blue}) < d(\text{blue}, \text{purple})$$

# Contrastive Learning in PIRL

## Dataset



## Loss Function

$$d(\text{blue}, \text{blue}) < d(\text{blue}, \text{green})$$

$$d(\text{blue}, \text{blue}) < d(\text{blue}, \text{purple})$$

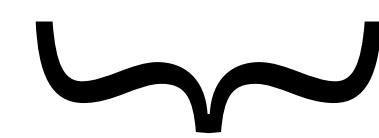


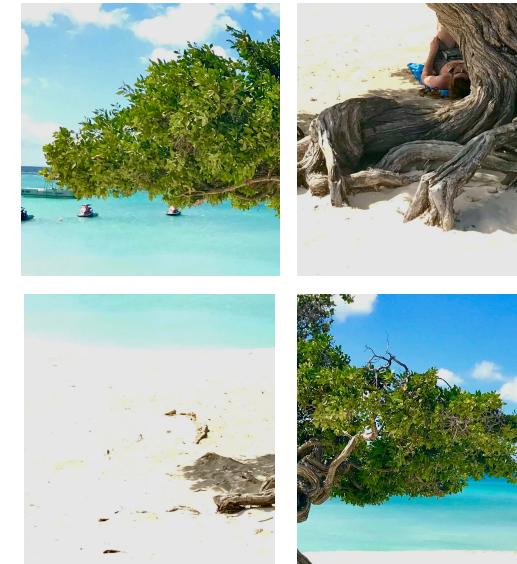
Image Feature &  
Patch Features

Random Images

I

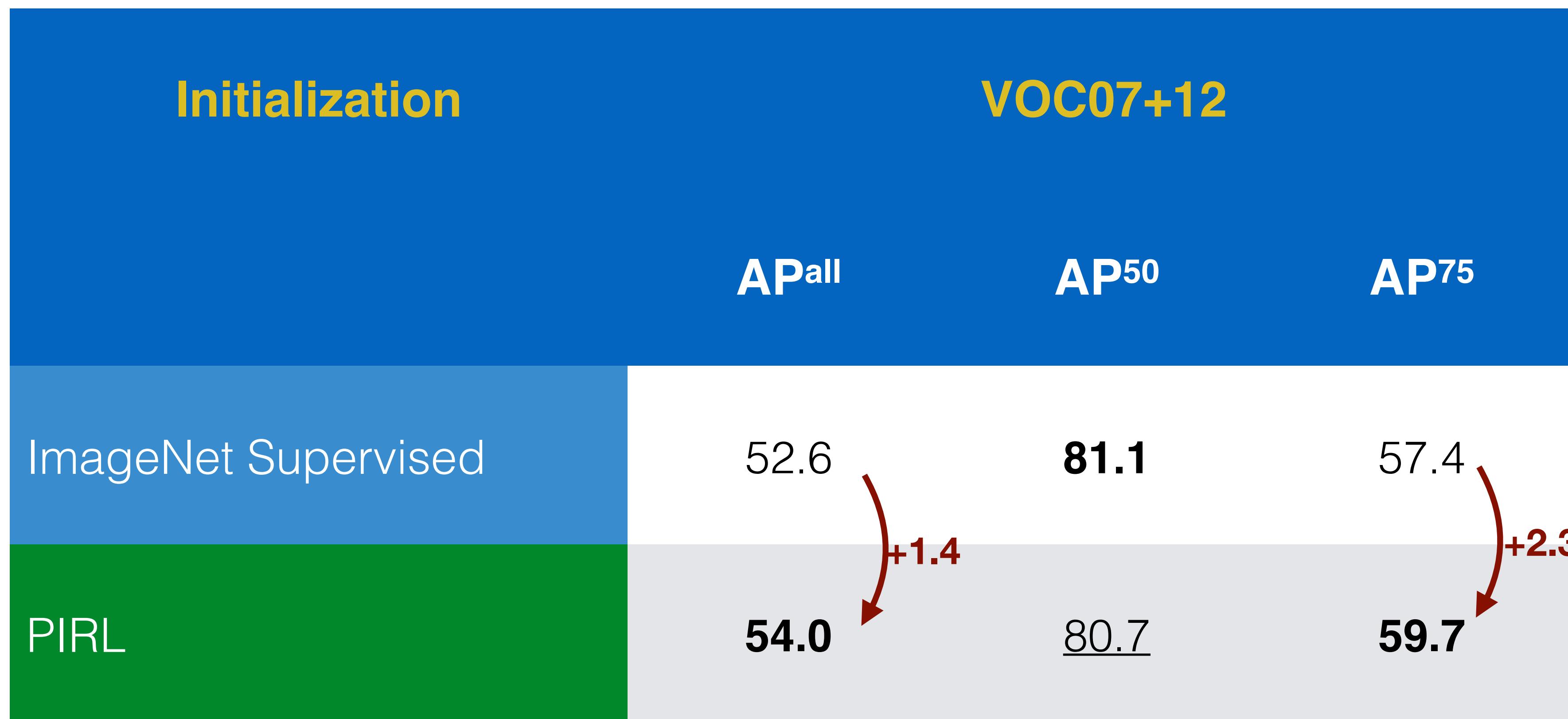


$I^t$



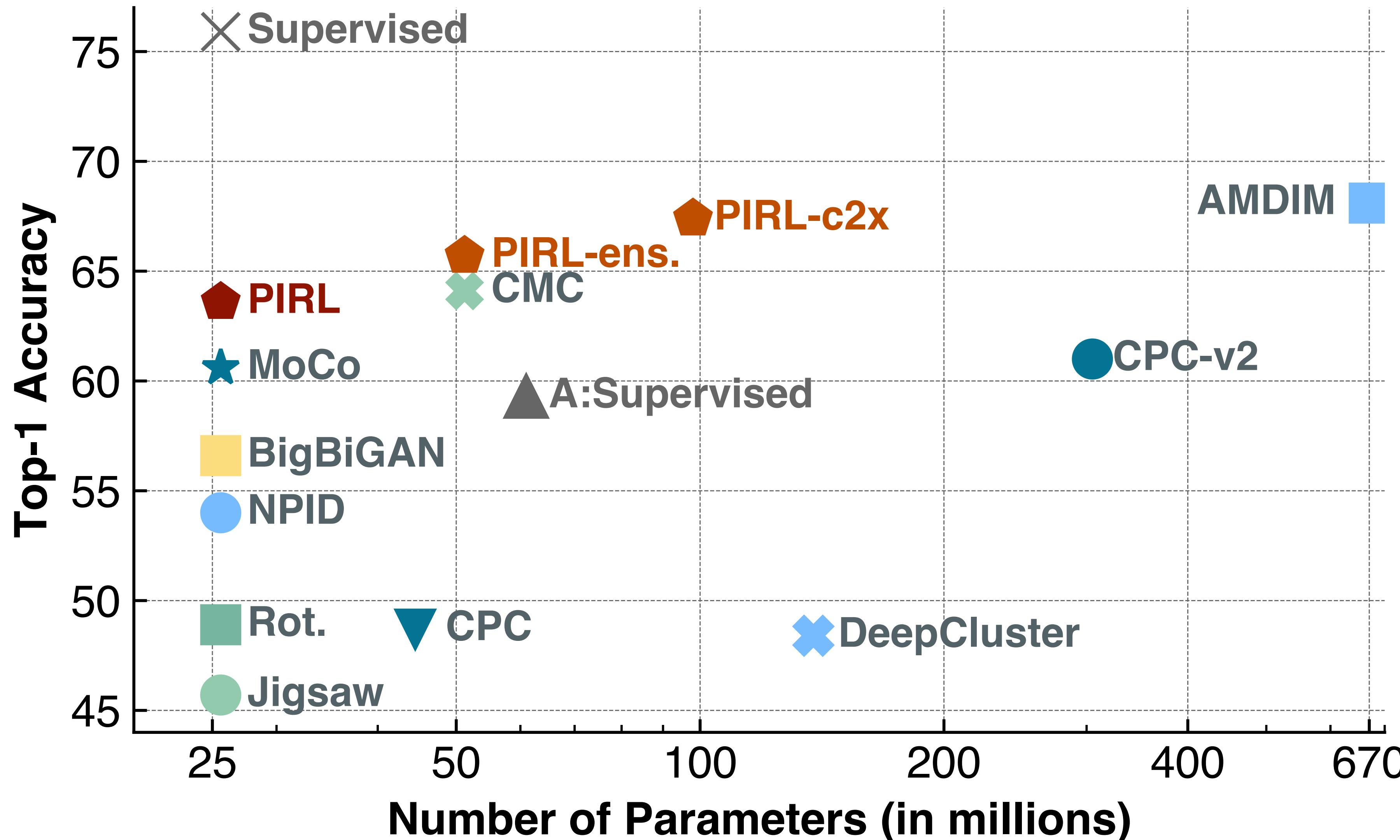
# Object Detection

- **Outperforms** ImageNet supervised pre-trained networks
- Full fine-tuning, no bells & whistles
- No extra data, changes in model architecture, fine-tuning schedule



# ImageNet Classification

- Linear classifiers on fixed features.



Method	Multi-view Invariance	Grouping	Performance
Pretext Task	No	No	Weak
PIRL	Yes	Weak	Moderate

Image & Patch as views      Relate other images using negatives      Outperform supervised features on some tasks

See also - CPCv2, MoCo, BoWNet, SimCLR

# Lack of Grouping in "Instance" based Contrastive Learning

**Positives**

$$d(\text{[blue box]} \text{ [blue box]}) < d(\text{[blue box]} \text{ [green box]})$$

$$d(\text{[blue box]} \text{ [blue box]}) < d(\text{[blue box]} \text{ [purple box]})$$



**Negatives**

Same image

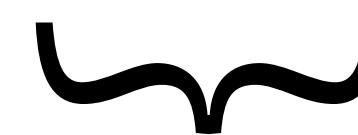
Views created by data augmentation

# Lack of Grouping in "Instance" based Contrastive Learning

## Positives

$$d(\text{blue box}, \text{blue box}) < d(\text{blue box}, \text{green box})$$

$$d(\text{blue box}, \text{blue box}) < d(\text{blue box}, \text{purple box})$$



No "groups"

## Negatives



Relate to other images  
using negatives

# How to add Grouping?

Audio-Video  
(AVID CMA)

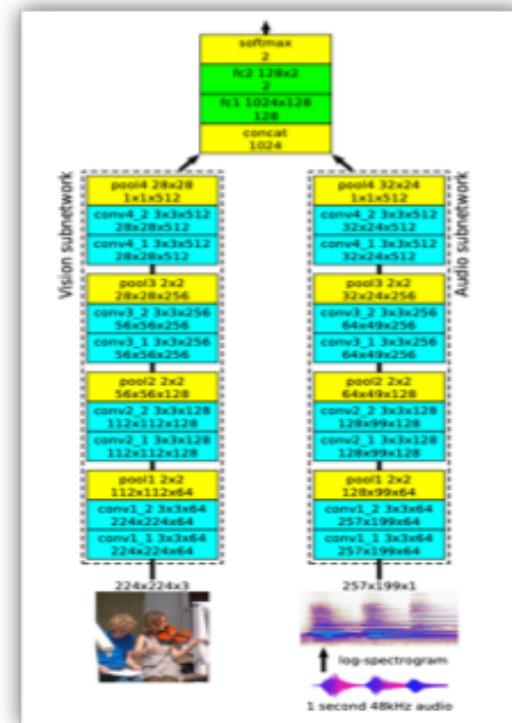
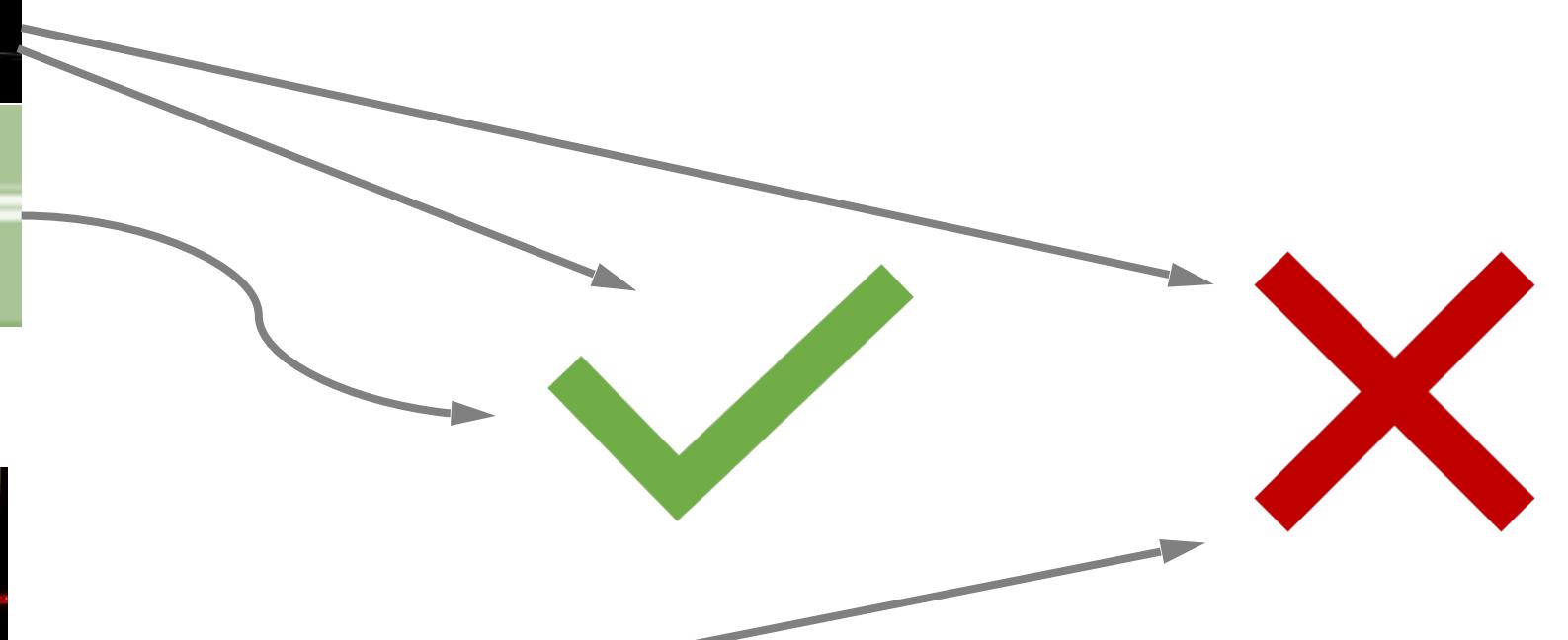
Images  
(SwAV)

# Audio Visual Instance Discrimination with Cross Modal Agreement (AVID + CMA)

Pedro Morgado, Nuno Vasconcelos, Ishan Misra



# Audio-visual correspondence

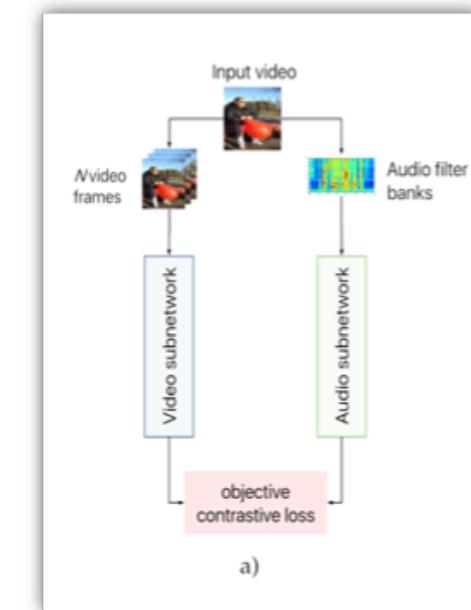


Look, listen & learn  
Arandjelovic et al.,  
2017.

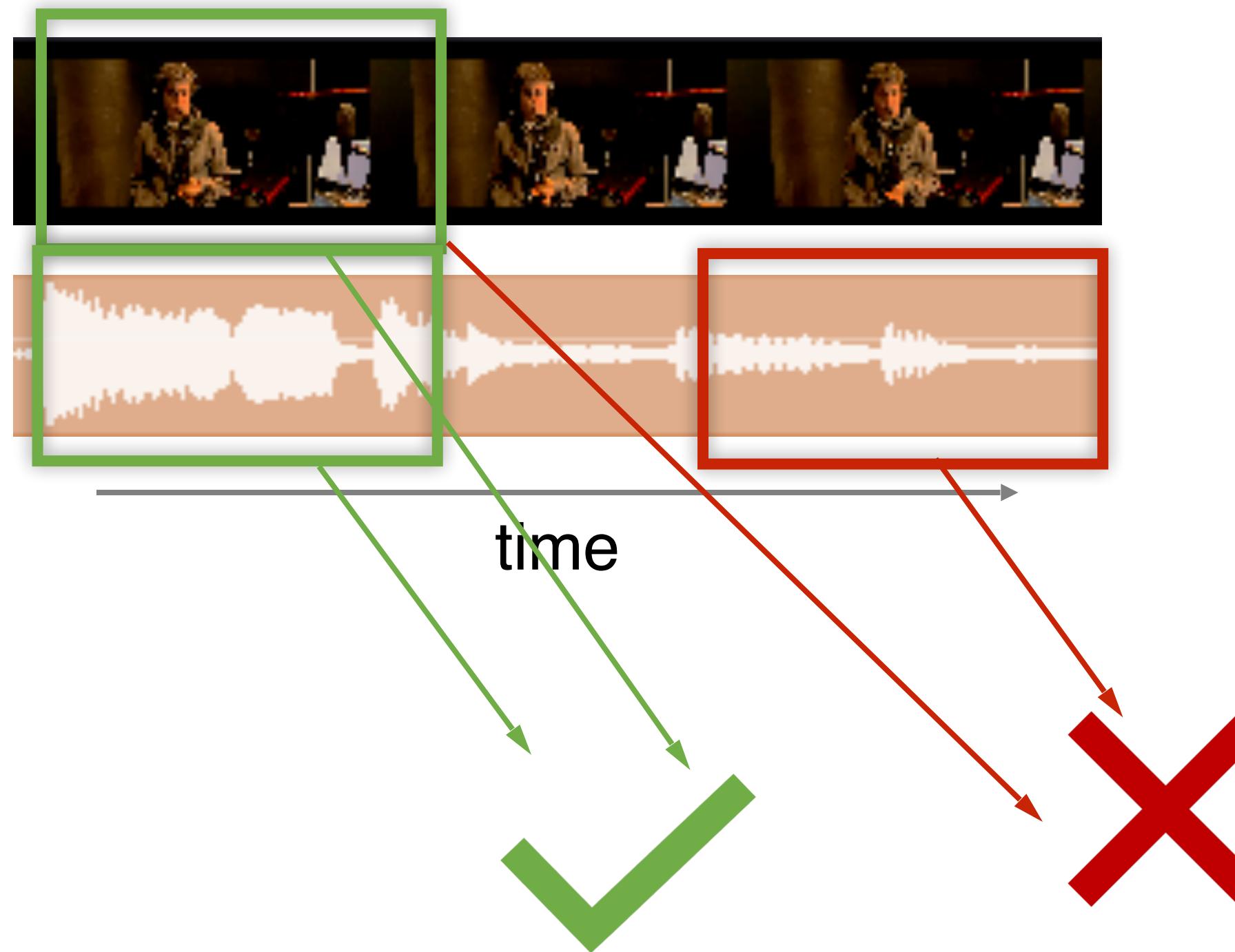
# Audio-visual synchronization



**Multisensory synchronization**  
Owens et al 2018.

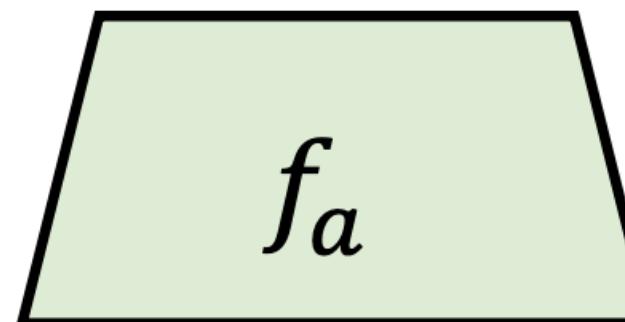
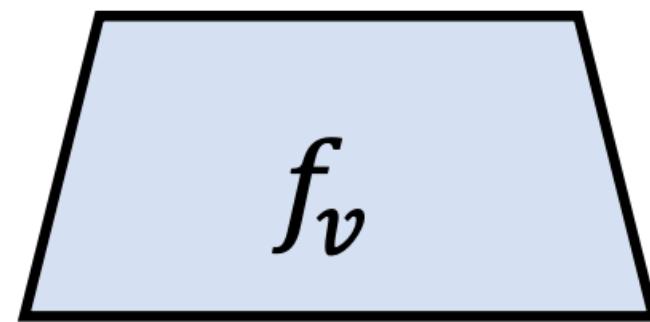


**Synchronization w/ curriculum**  
Korbar et. al 2018.



Method	Multi-view Invariance	Grouping	Performance
AV Sync	Yes	No	Weak

# Contrastive (Audio Video Instance Discrimination)



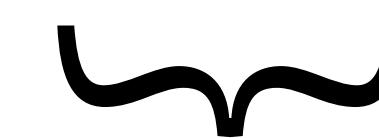
**Positives**

$$d(\text{blue box}, \text{blue box}) < d(\text{blue box}, \text{green box})$$

$$d(\text{blue box}, \text{blue box}) < d(\text{blue box}, \text{purple box})$$



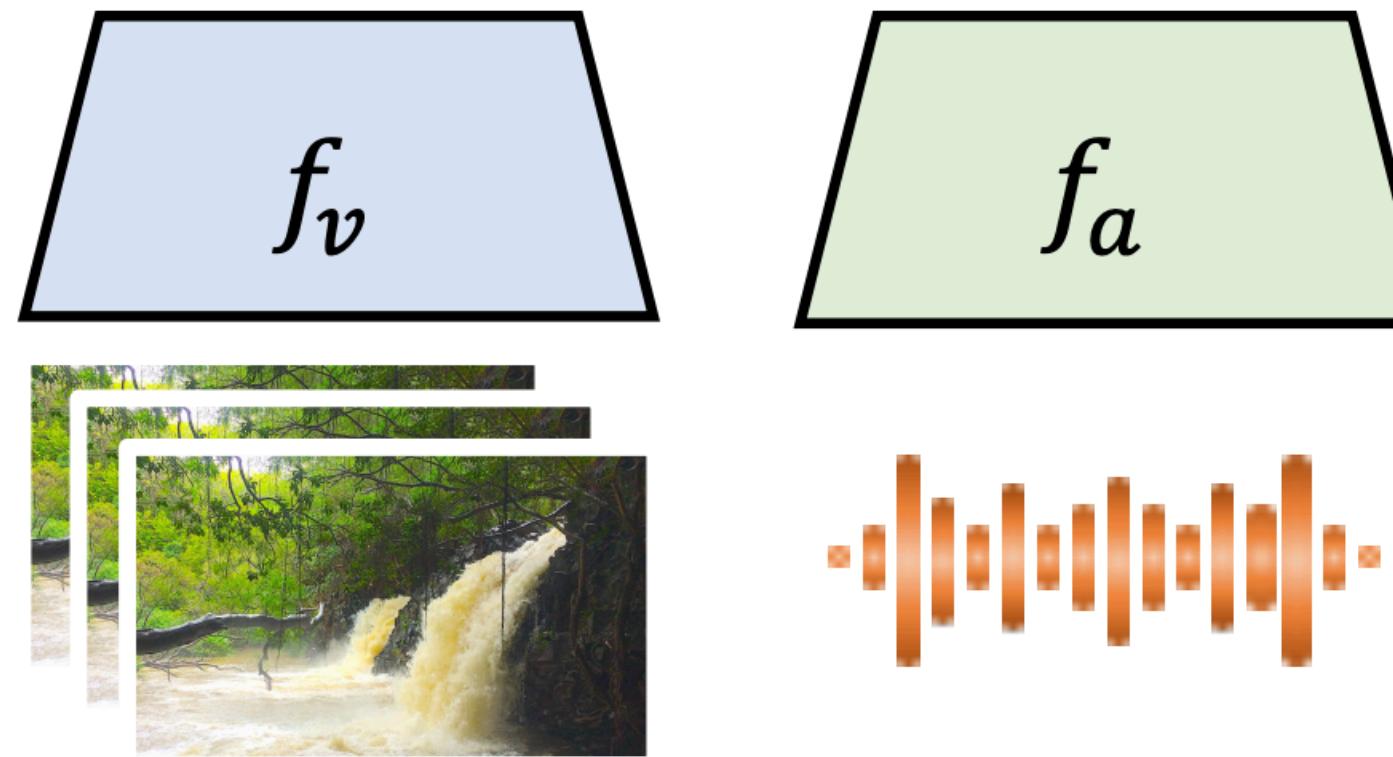
**Negatives**



Audio & Video  
(same sample)

Relate to other video/audio  
using negatives

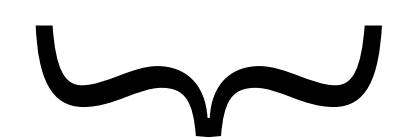
# Contrastive (Audio Video Instance Discrimination)



**Positives**

$$d(\text{blue box}, \text{blue box}) < d(\text{blue box}, \text{green box})$$

$$d(\text{blue box}, \text{blue box}) < d(\text{blue box}, \text{purple box})$$



Audio & Video  
(same sample)

**Negatives**

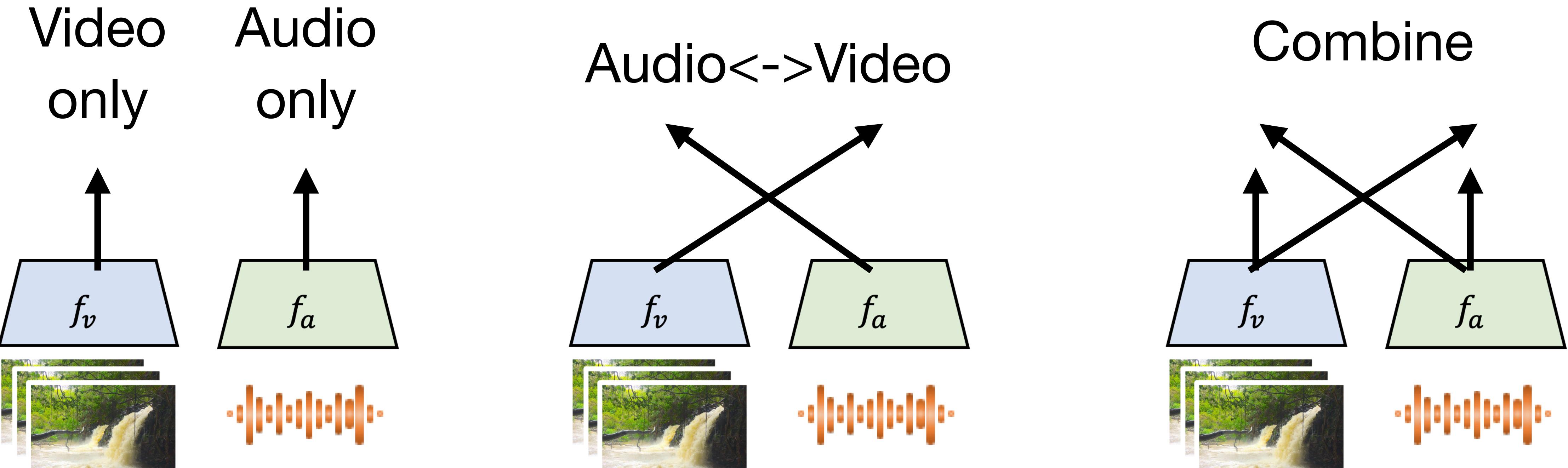
$$d(\text{blue box}, \text{purple box})$$



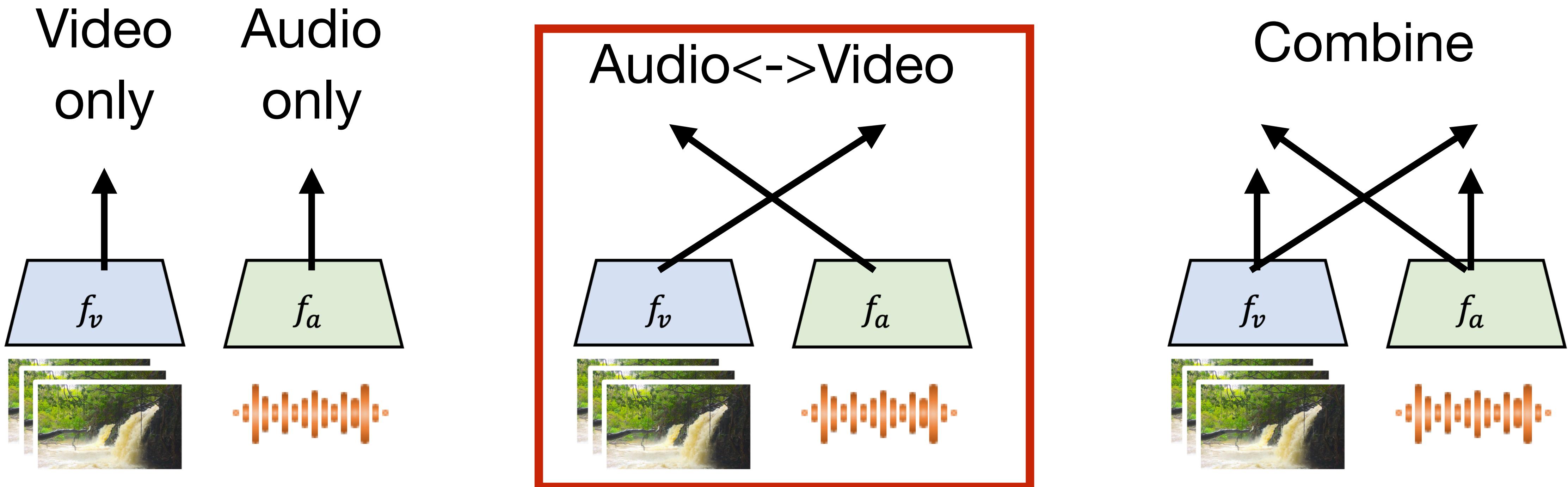
Relate to other video/audio  
using negatives

**Limited Grouping**

# Importance of Multi-view

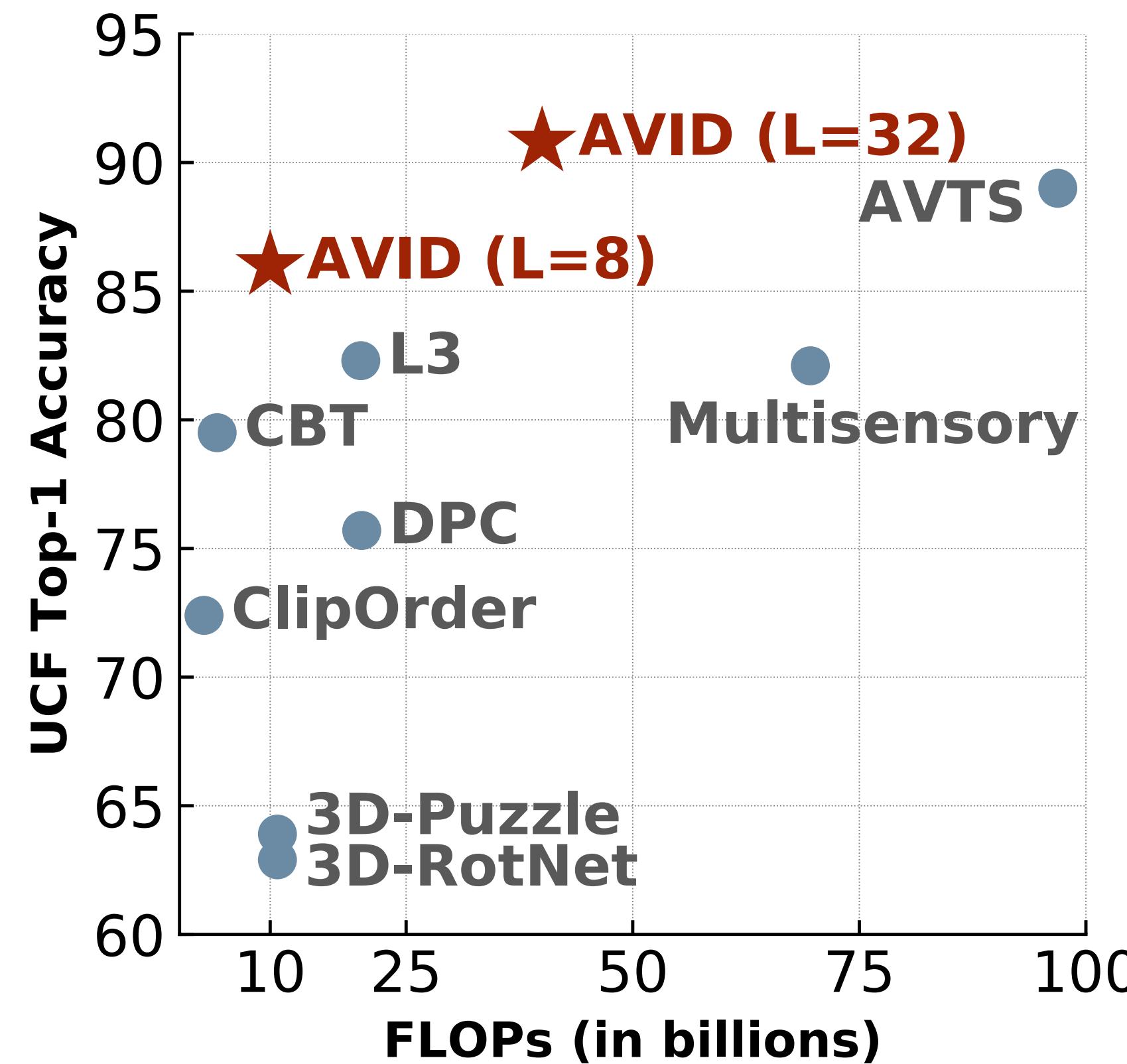


# Importance of Multi-view

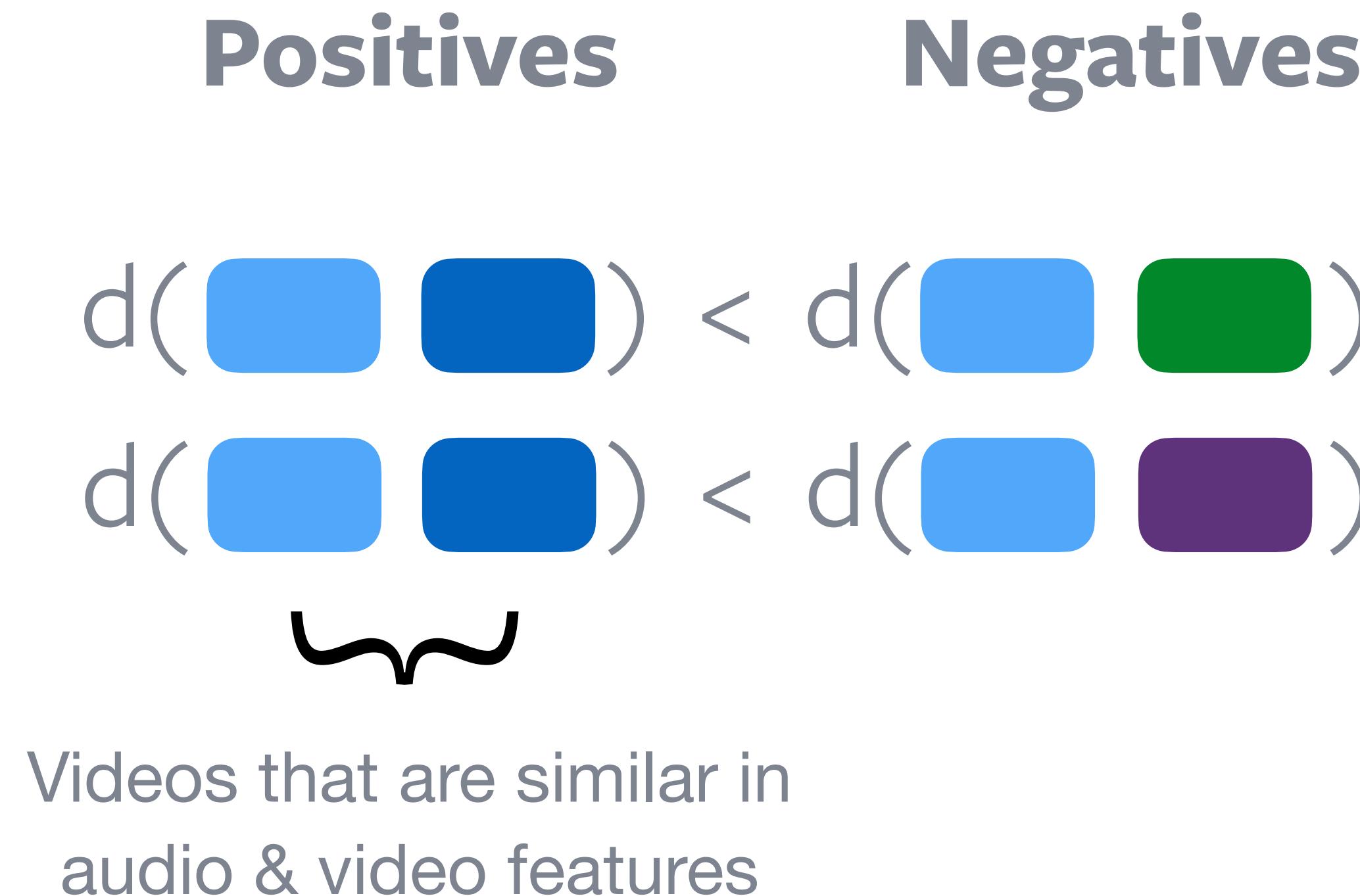
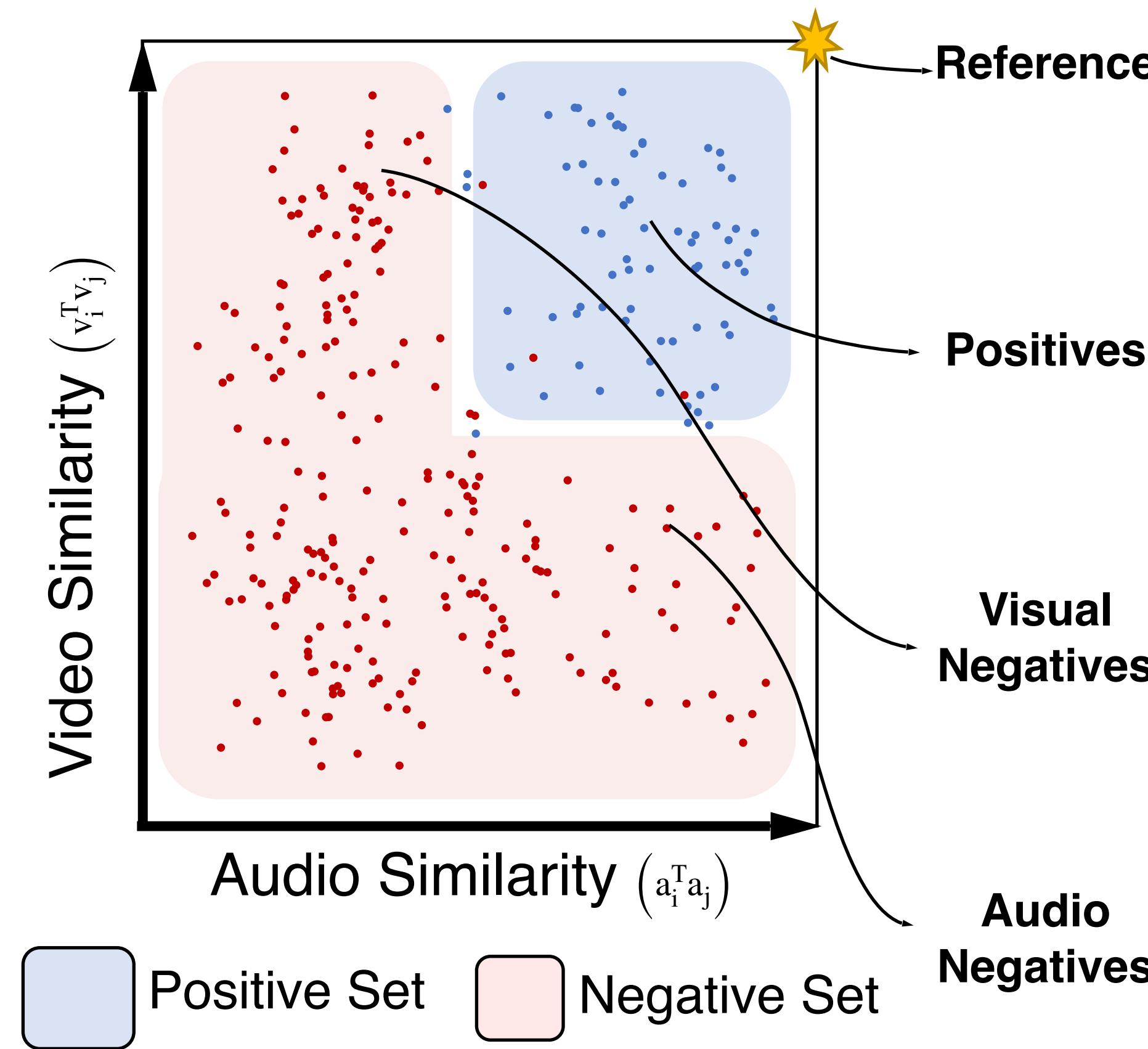


Outperforms all variants

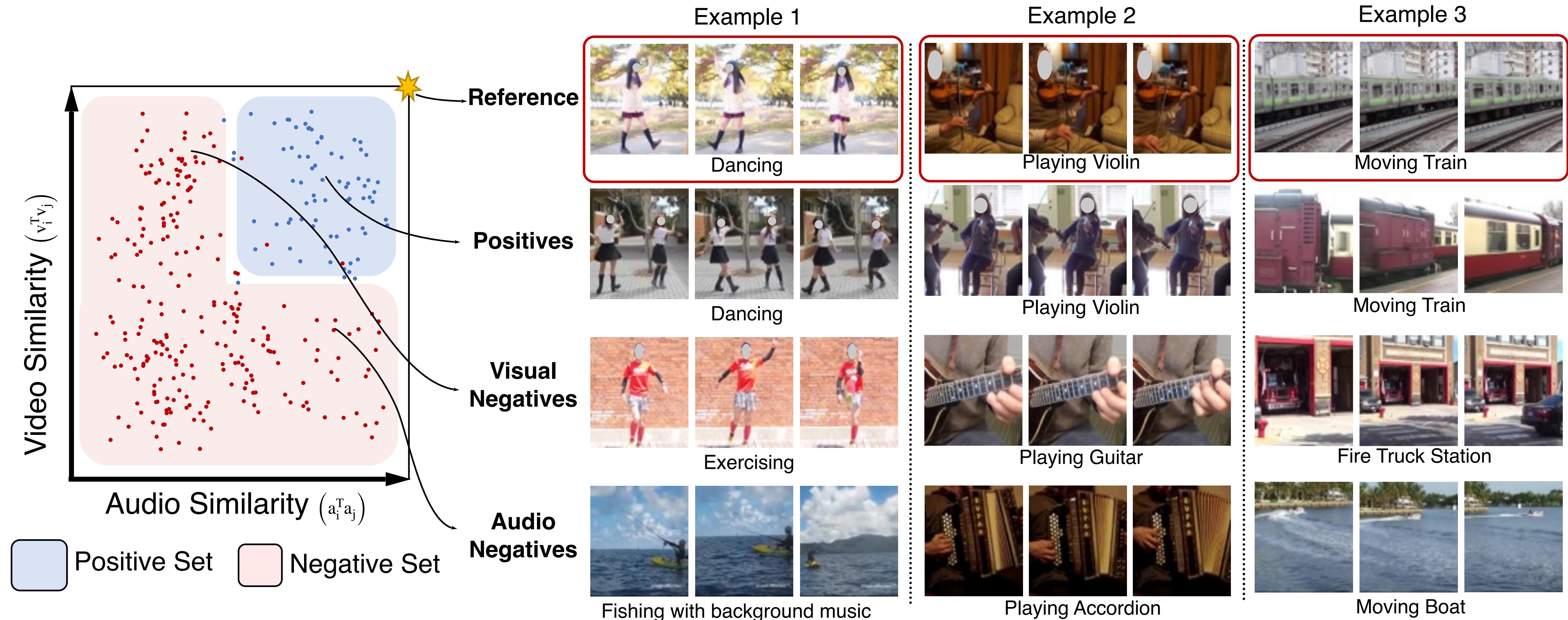
Method	Multi-view Invariance	Grouping	Performance
AV Sync	Yes	No	Weak
AVID	Yes	Weak	Good

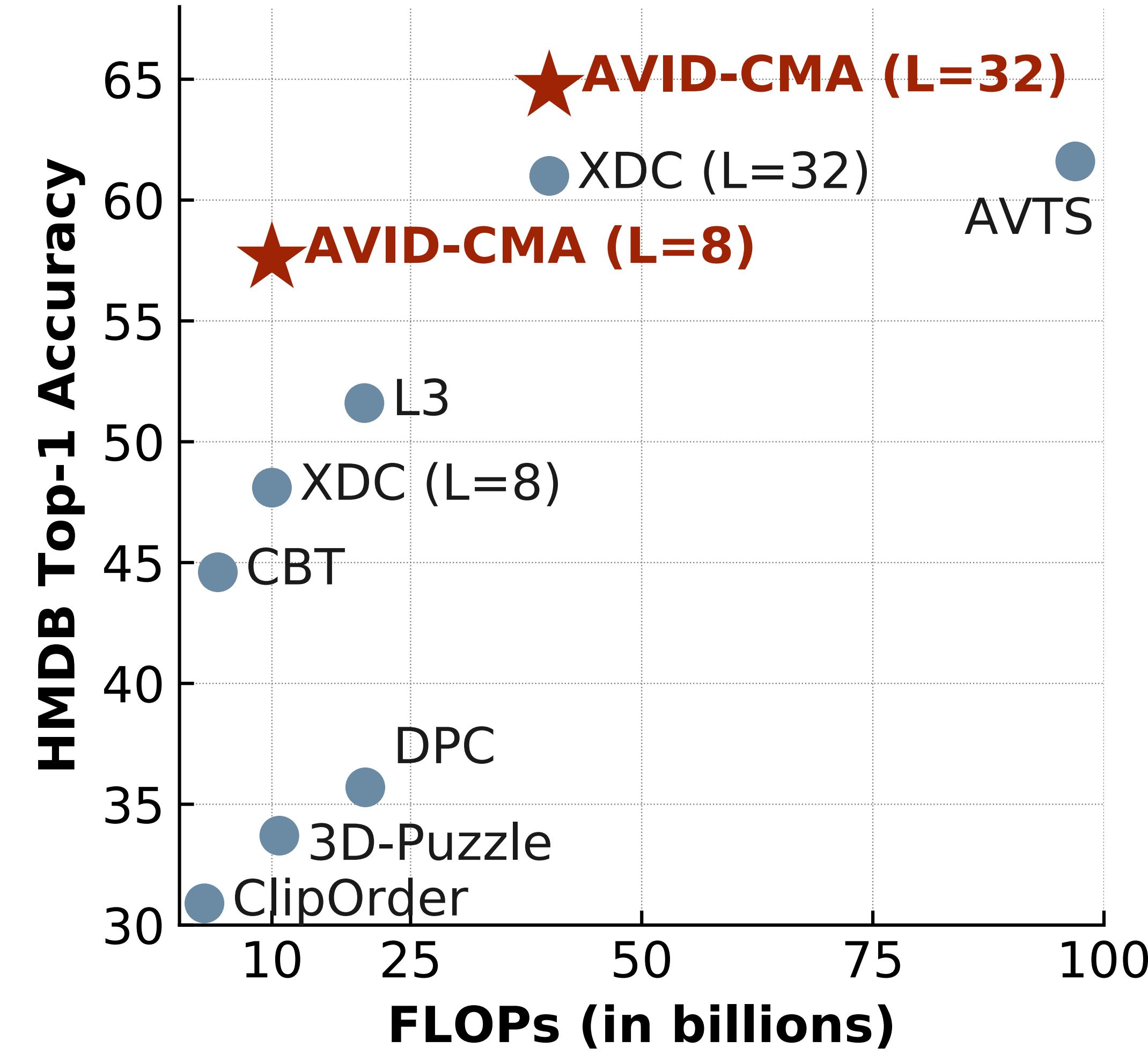
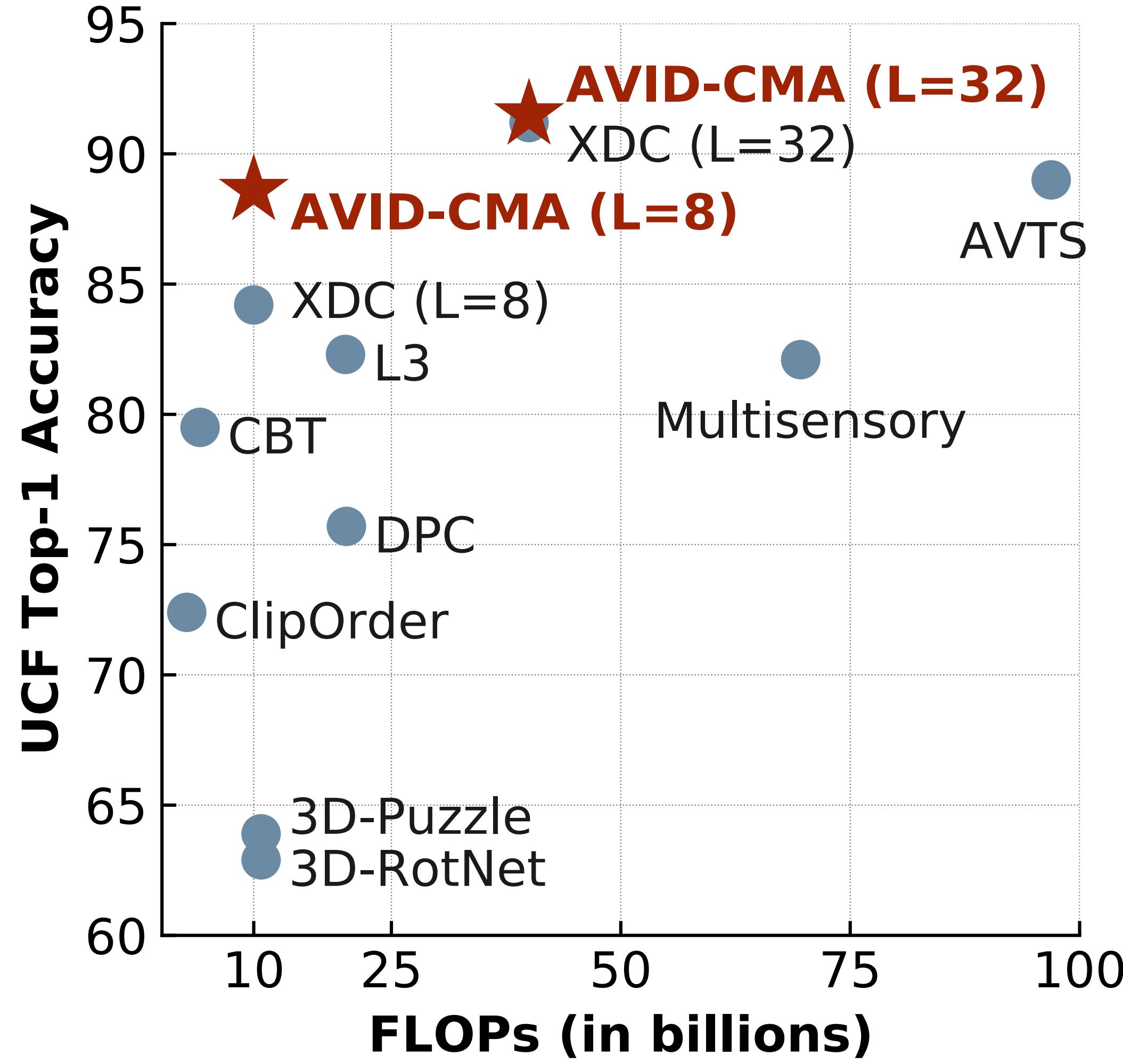


# Grouping using Audio-visual Agreements (CMA)



# Grouping using Audio-visual Agreements (CMA)





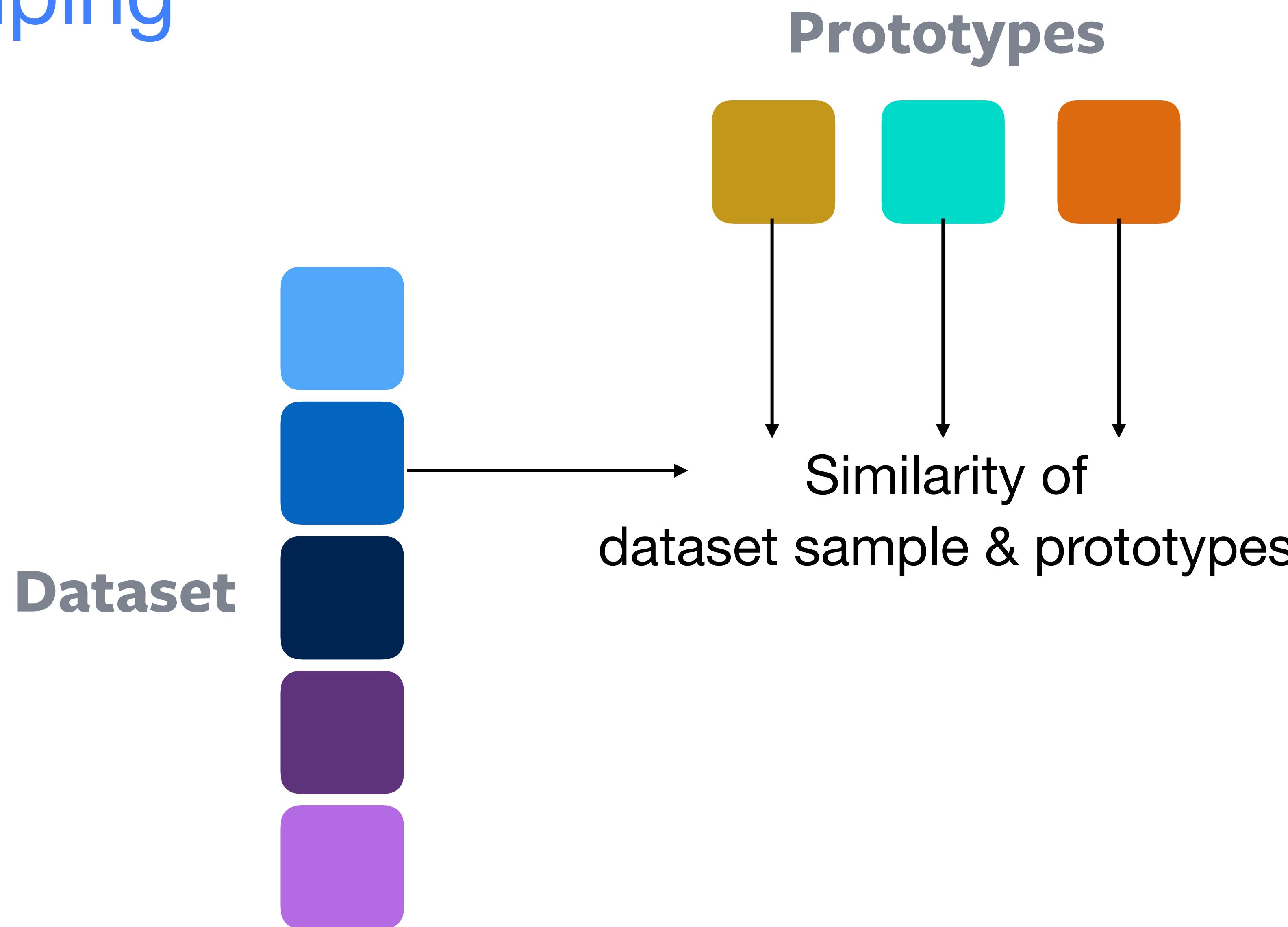
Method	Multi-view Invariance	Grouping	Performance
AV Sync	Yes	No	Weak
Contrastive	Yes	Weak	Good
CMA	Yes	Yes	Better

# Swapping Assignments between Views (SwAV)

Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, Armand Joulin



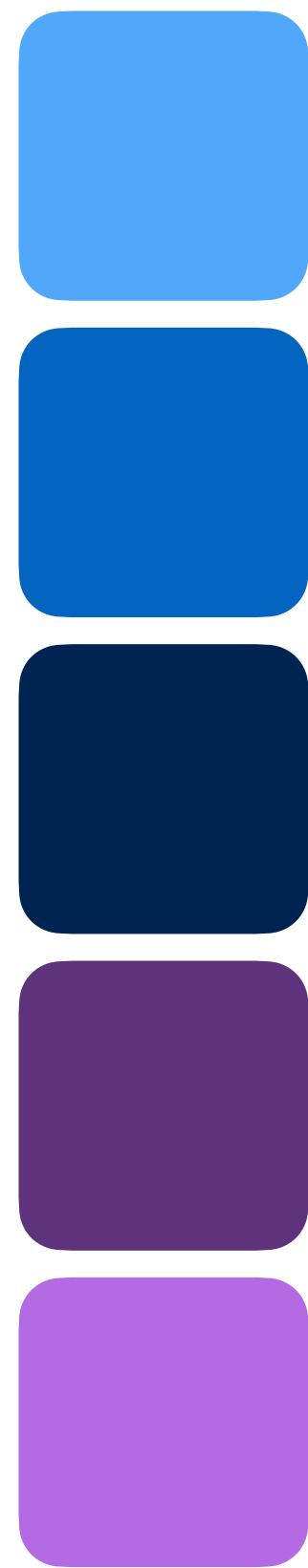
# Grouping



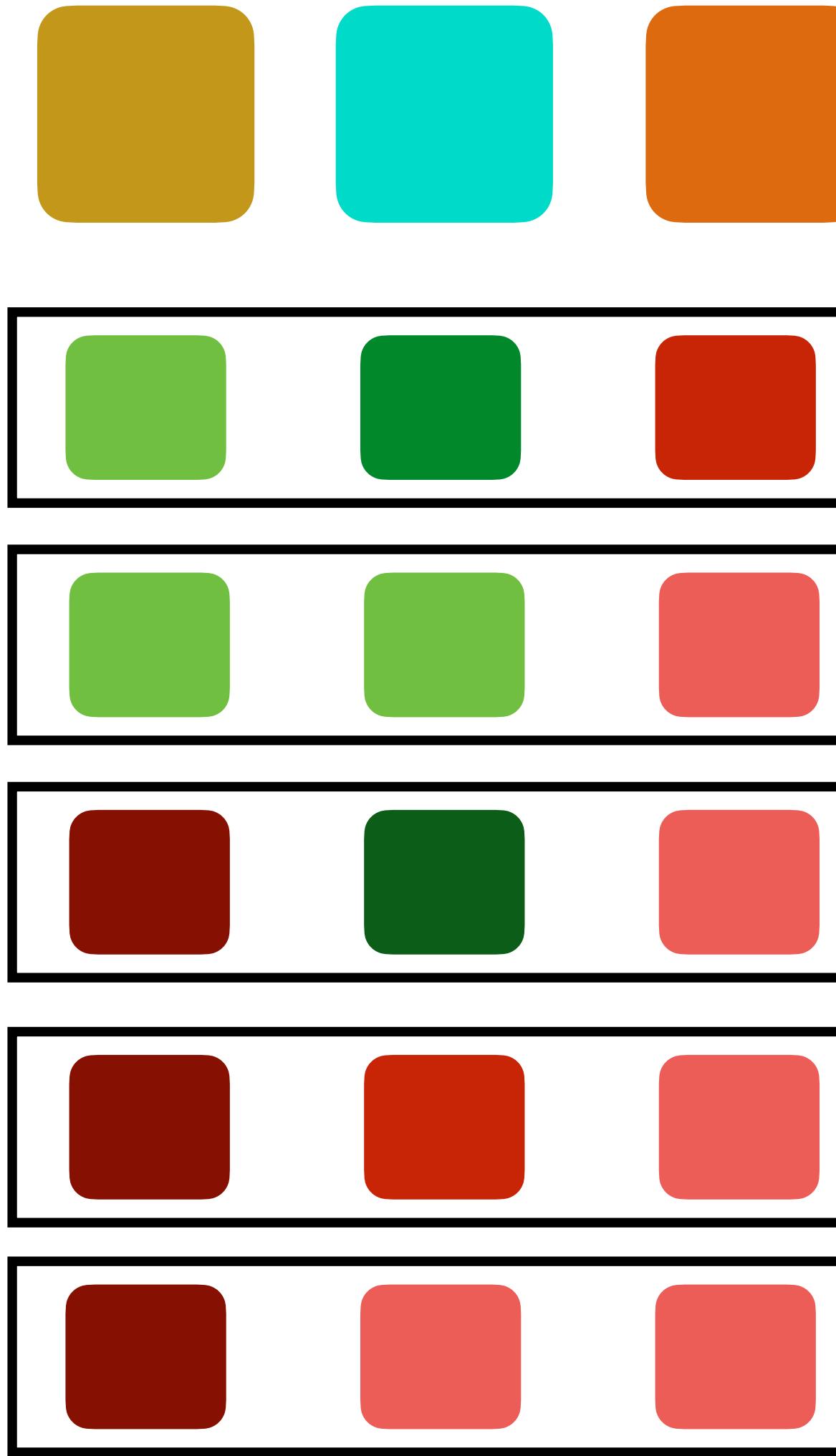
See also - SeLa by Asano et al., 2019 <sup>36</sup>

# Grouping

Dataset

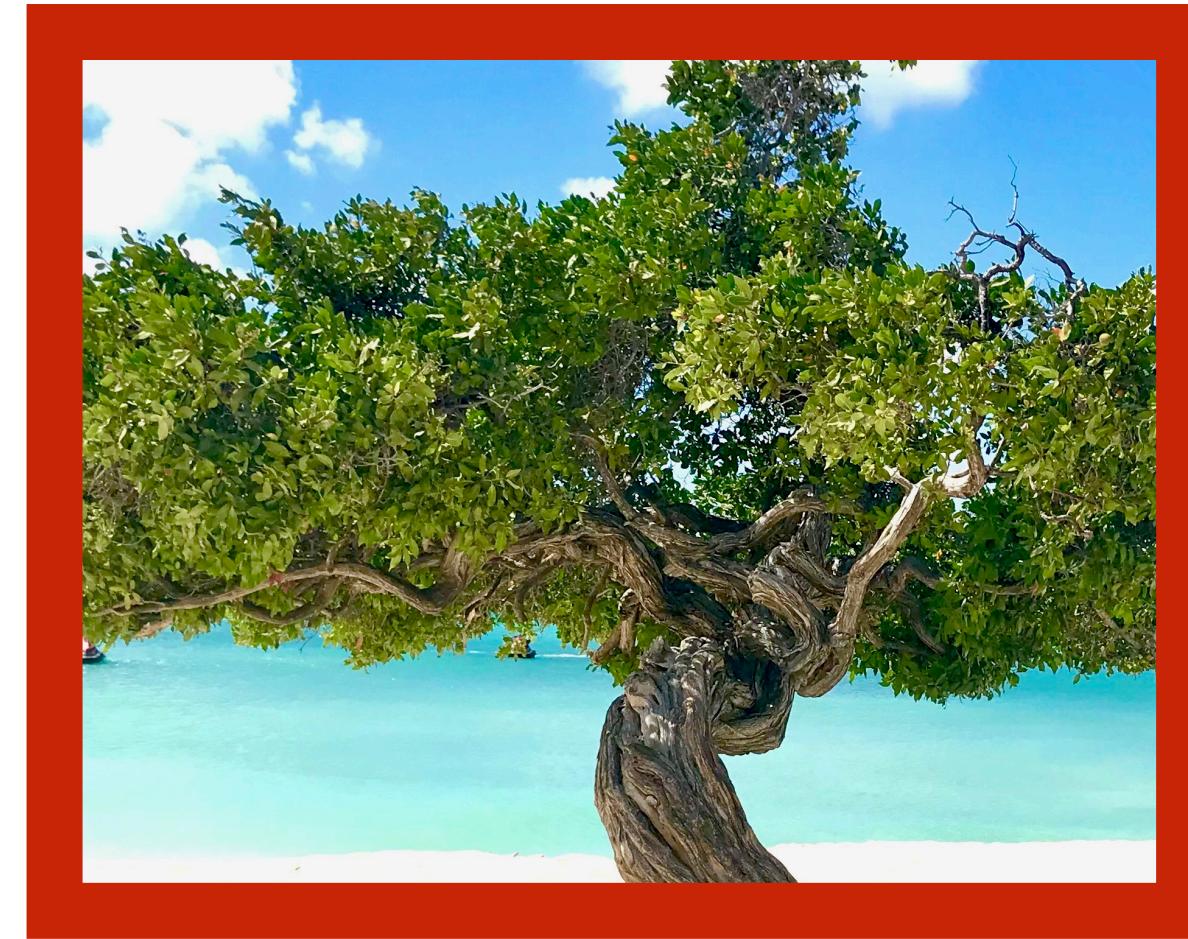
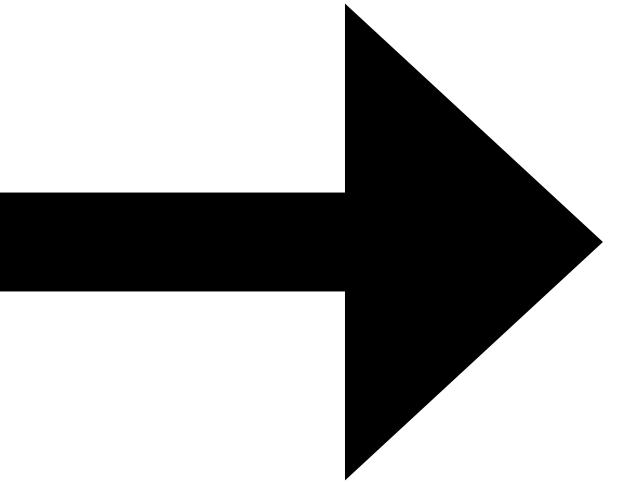
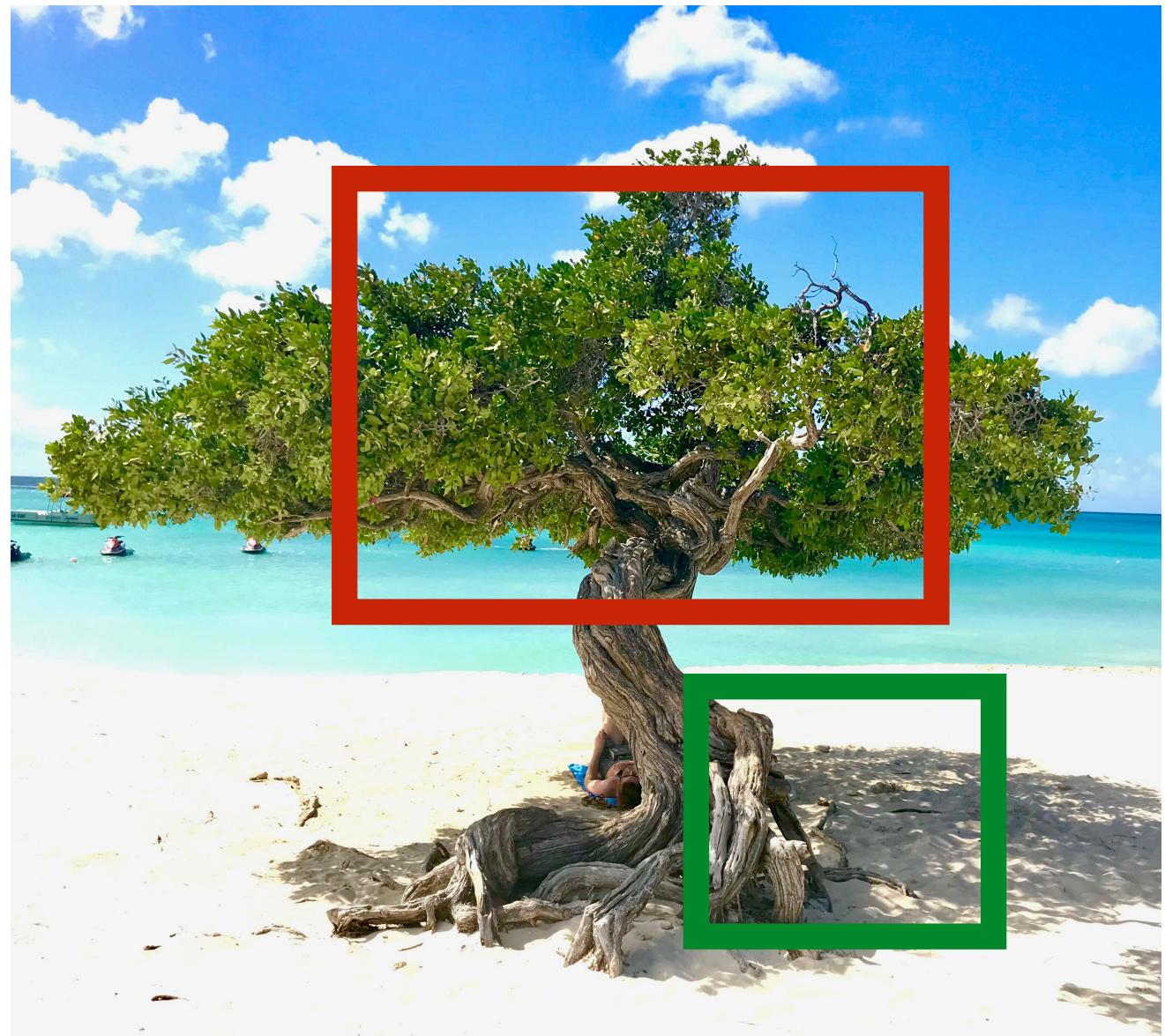


Prototypes



Codes

# Multi-crop (Multi-view)

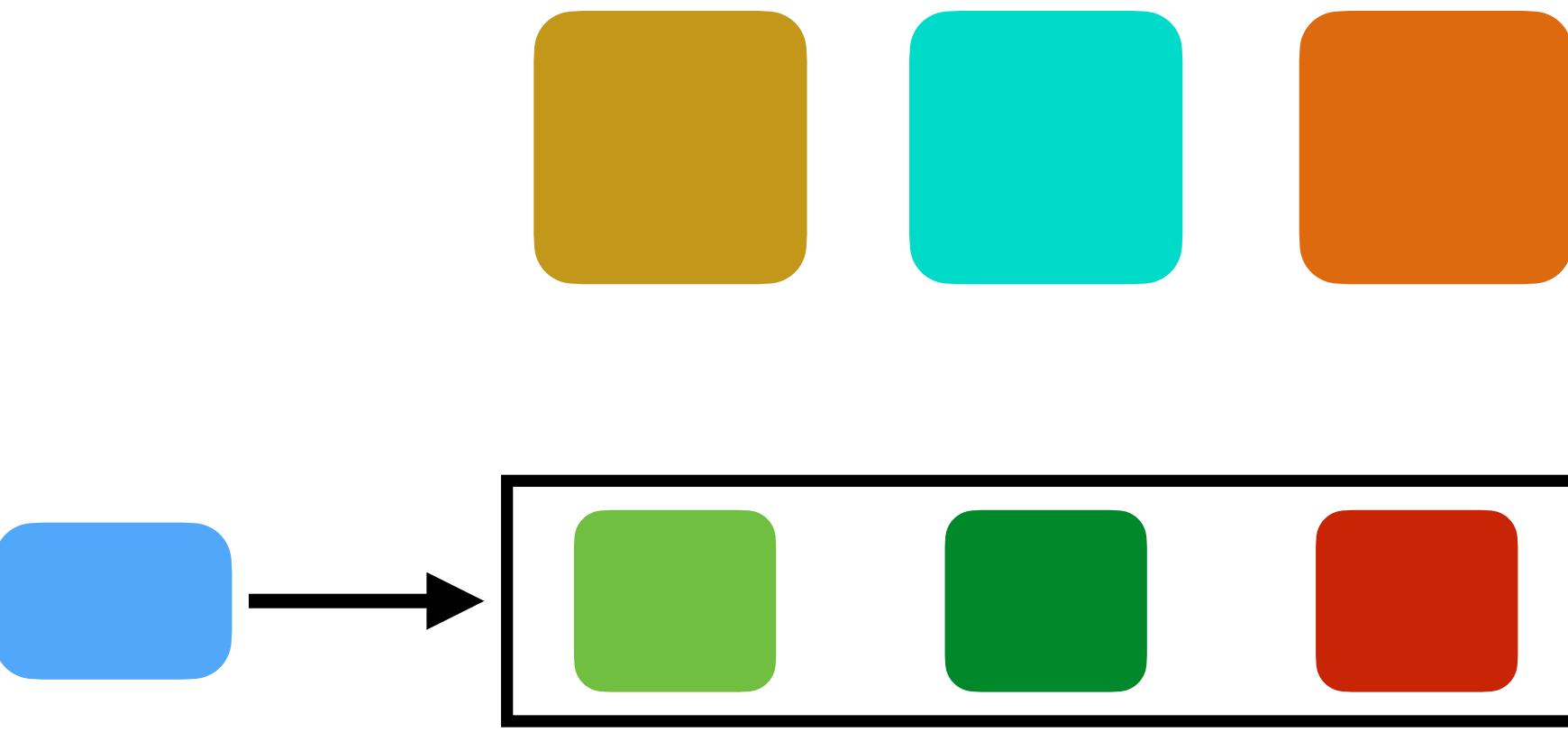


Compare crops of  
different sizes & resolutions

# Prototypes



$f_\theta$



Code 1

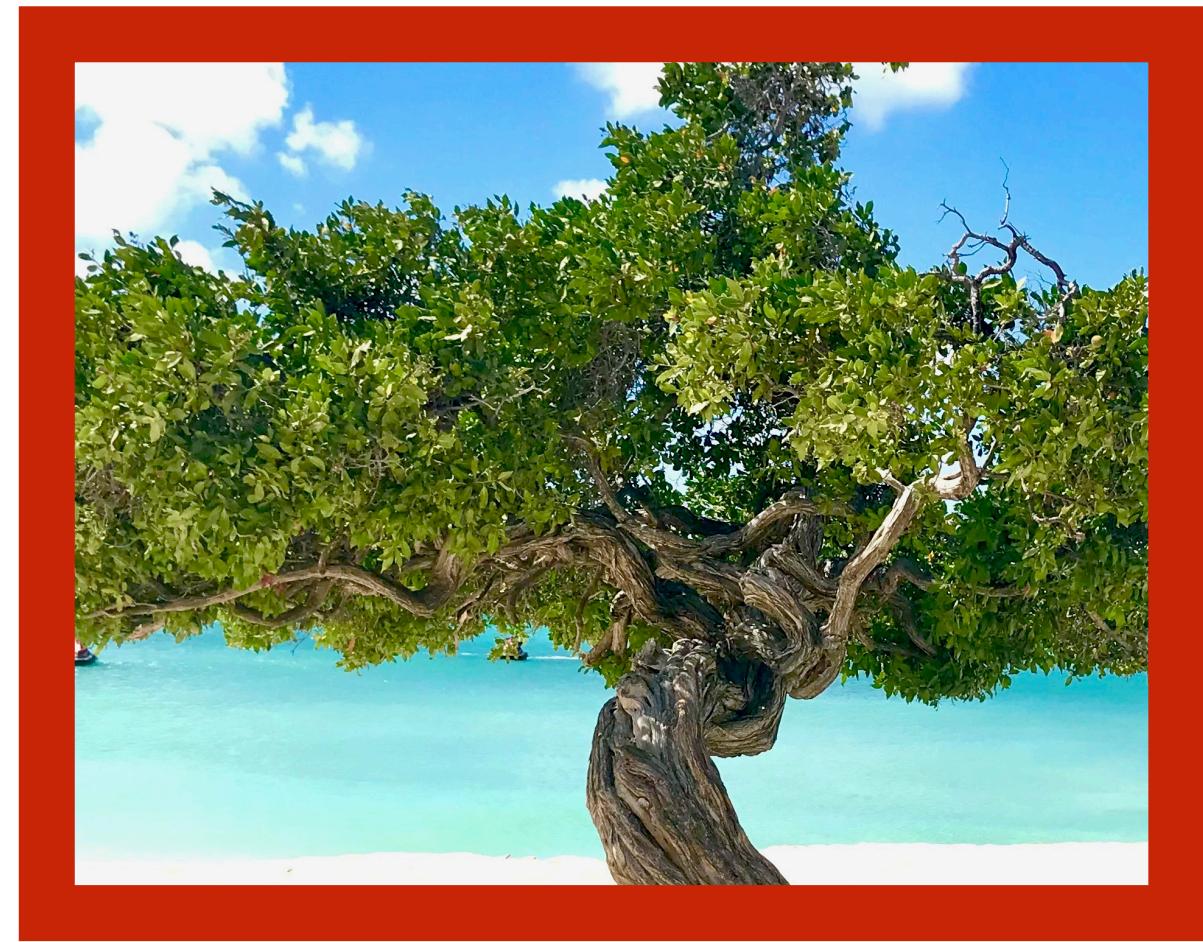


$f_\theta$



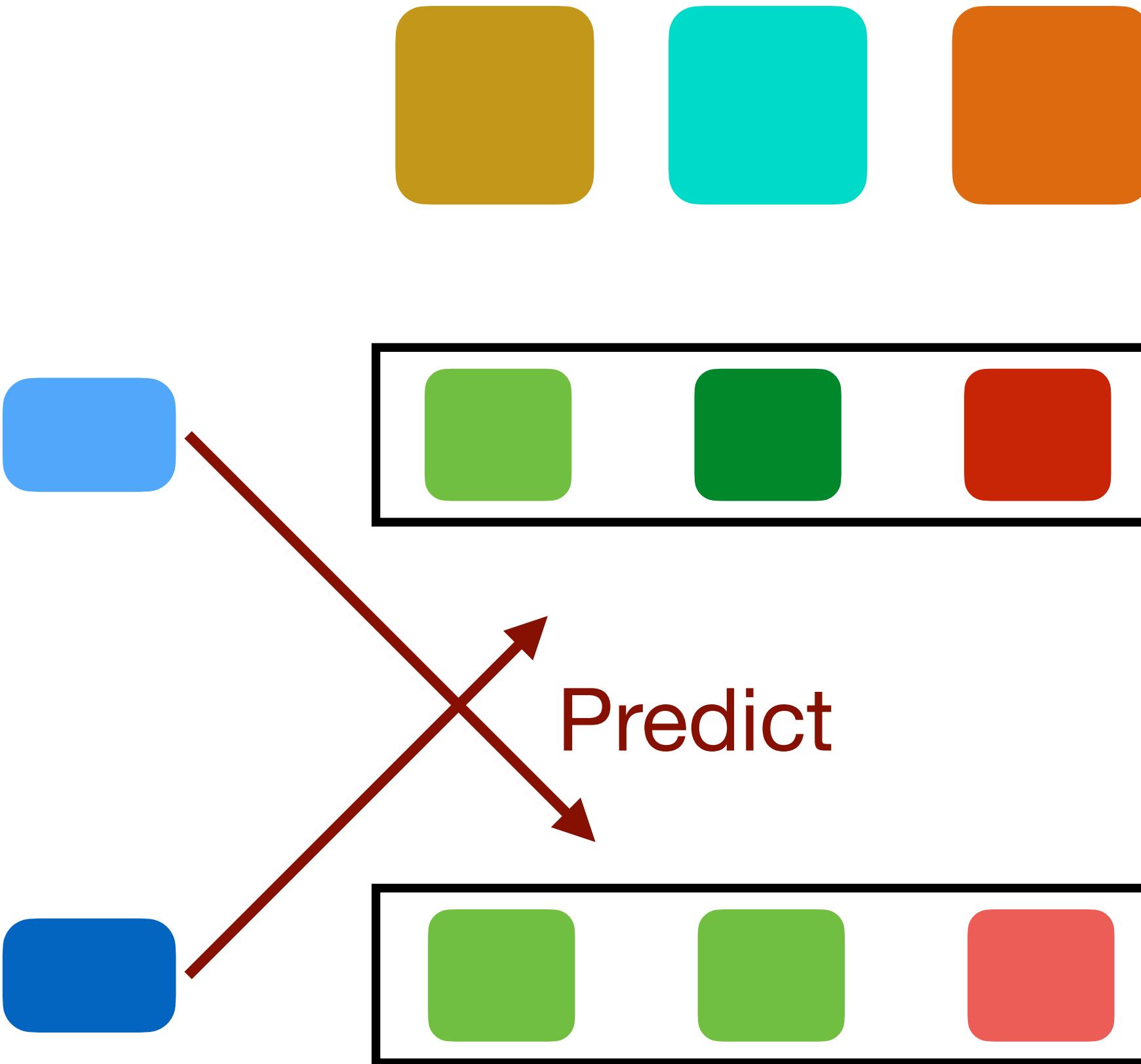
Code 2

# Prototypes



$f_{\theta}$

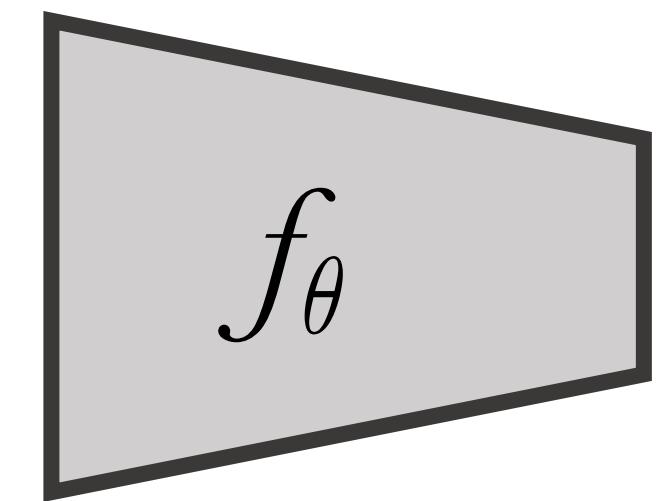
$f_{\theta}$



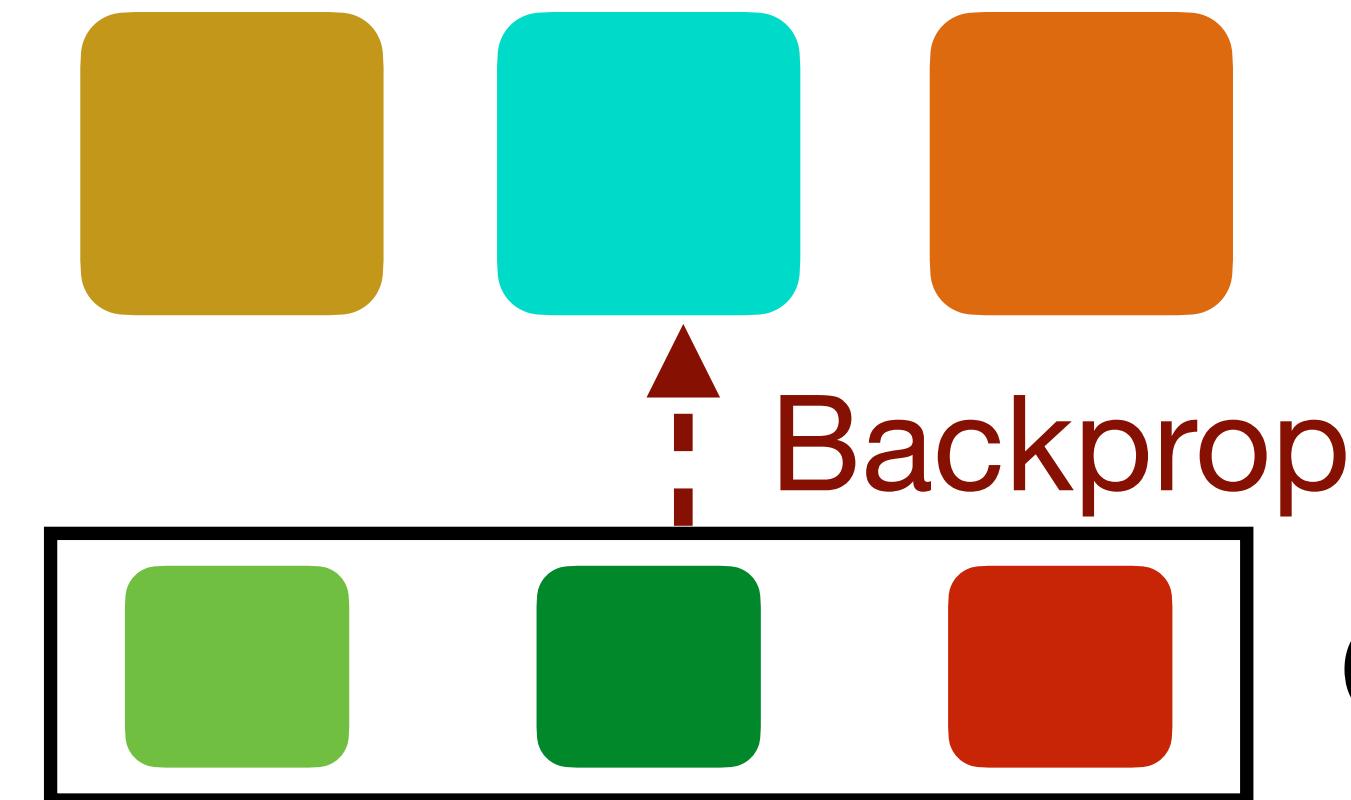
Code 1

Code 2

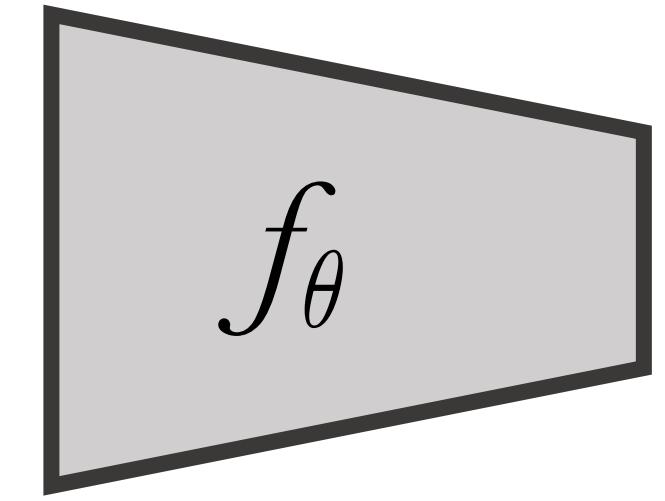
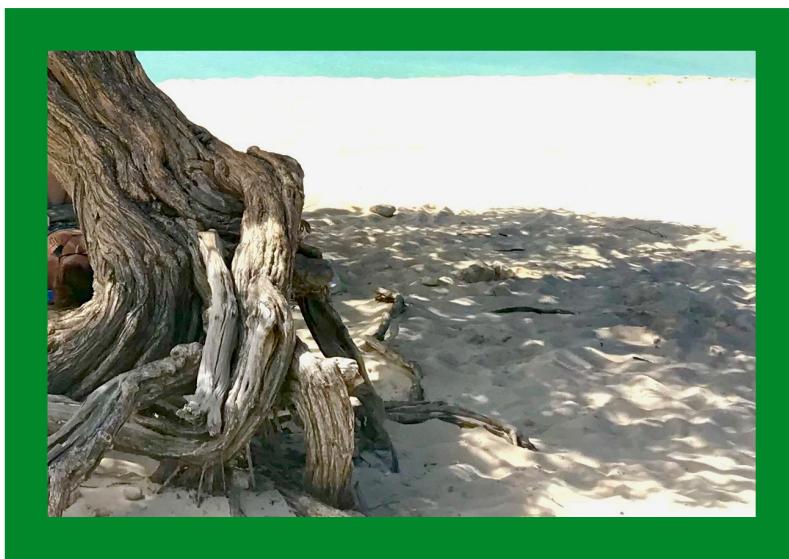
# Prototypes



← - - - Backprop



Code 1



Code 2

Not contrastive!

# Multi-crop (Multi-view)

Method	ImageNet Top-1		$\Delta$
	w/ Multi-crop	Multi-crop	
SimCLR	68.2	70.6	+2.4
SeLa-v2	67.2	71.8	+4.6
DeepCluster-v2	70.2	74.3	+4.1
SwAV	70.1	74.1	+4.0

Multi-view improves ALL methods

# Key Results

## Linear Classifier (Fixed Features)

## Detection

	ImageNet	Places	iNaturalist	VOC07+12	COCO
Supervised	76.5	53.2	46.7	81.3	40.8
Prior self-supervised	71.1 (-5.4)	52.1	38.9	82.5	42.0
SwAV	75.3 (-1.2)	<b>56.7</b>	<b>48.6</b>	<b>82.6</b>	<b>42.1</b>

# Practical advantages of SwAV

- Trains on 4-8 GPUs
- **Faster convergence** than prior work (SimCLR, MoCov2)
  - Smaller compute requirements.
  - **2x faster** than MoCo-v2 on 8 GPUs
    - 72% after 100h vs. 71% after 200h
- Better results



Code & Models - <https://github.com/facebookresearch/swav>  
PyTorch Lightning implementation on the way

Method	Multi-view Invariance	Grouping	Performance
Pretext Task	No	No	Weak
PIRL	Yes	Weak	Moderate
SwAV	Yes	Yes	Good

Outperforms ImageNet supervised pretraining  
on **all transfer tasks**

# What invariances matter?

- Our set of invariances are overfit to ImageNet
- Need to evaluate on **different image distributions (uncurated data)**



Horizontal flipping may not be a good idea ...

# Scalable objectives

- Contrastive learning converges very slowly & scales poorly to large data
- Notion of grouping is important

# The Future ...

Method	Multi-view Invariance	Grouping	???	Performance
Pretext Task	No	No	No	Weak
PIRL	Yes	Weak	No	Moderate
SwAV	Yes	Yes	No	Good
<b>Your method</b>	Yes	Yes	Yes	<b>BEST</b>

# Thanks!