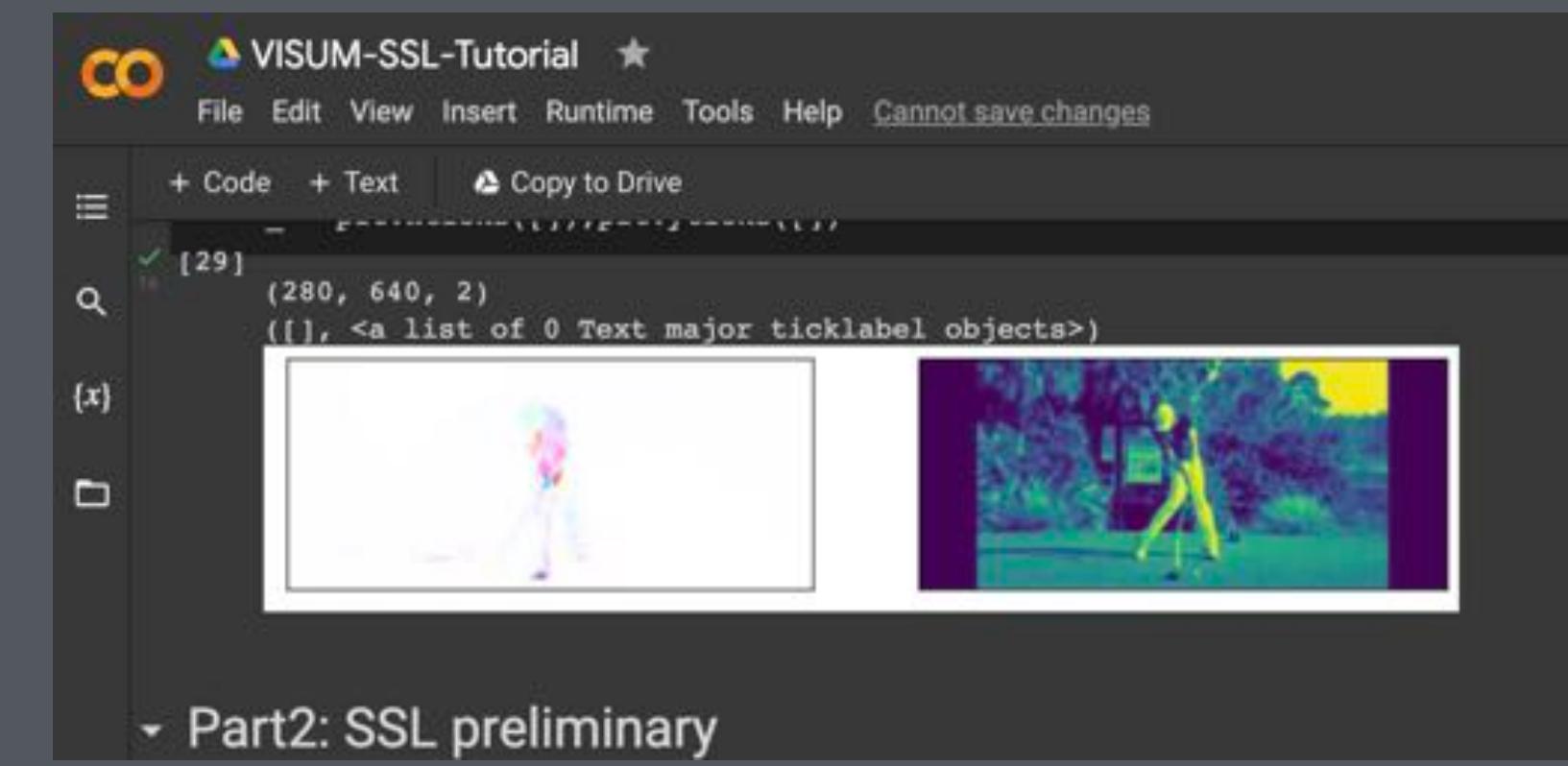
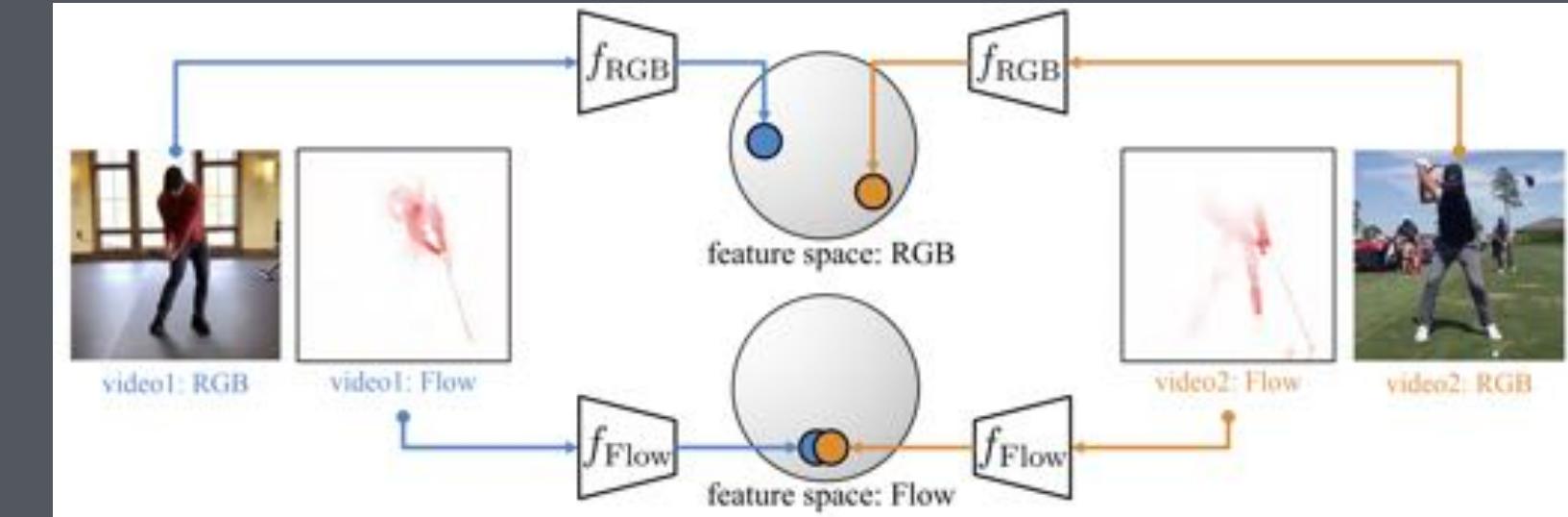
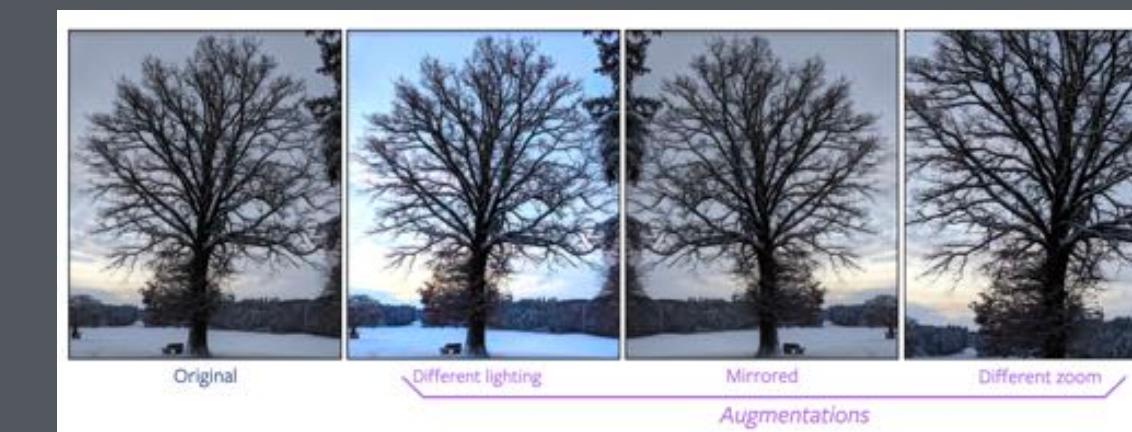
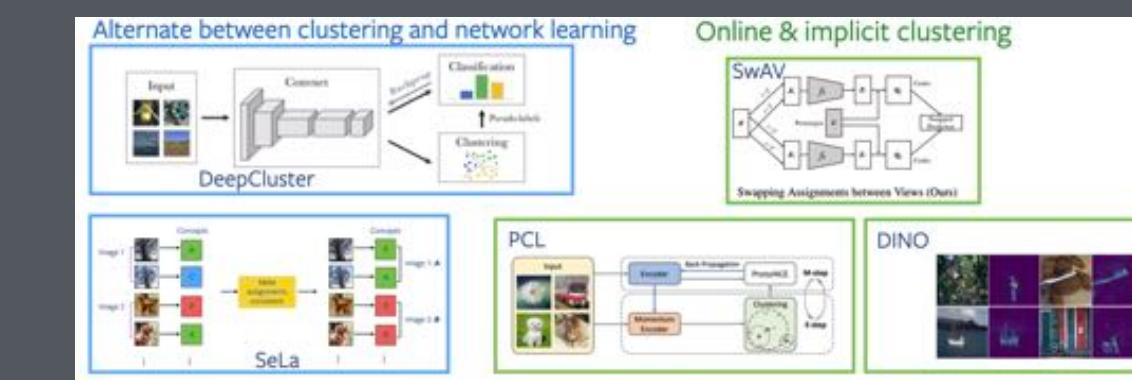


Self-supervised learning for computer vision from images, video and audio. The *why* & *how*.



• Part2: SSL preliminary

@ VISUM SUMMERSCHOOL

YUKI M. ASANO, QUVVA LAB MANAGER, ASSISTANT PROFESSOR, UNIVERSITY OF AMSTERDAM

TENGDA HAN, VISUAL GEOMETRY GROUP, UNIVERSITY OF OXFORD

What is *Self-supervised learning (SSL)*?

Why do we want to do *SSL*?

How to do *SSL*?

Deep-dive into one image-*SSL* method

Why multi-modal data for *SSL*?

Brief intro into two video-*SSL* methods

Hi, I'm Yuki

- Assistant Professor with VISLab
 - Self-supervised Learning
 - Multi-modal Learning
 - Privacy and Bias in Computer Vision
 - Other Interdisciplinary Research
- Prior to this:
 - PhD at VGG in Oxford;
 - Applied Mathematics/Physics/Economics at Oxford/Munich/Hagen
- Scientific Manager (with Prof. Snoek, Prof. Welling, Prof. Gavves)
- Qualcomm-UvA (QUvA) Lab



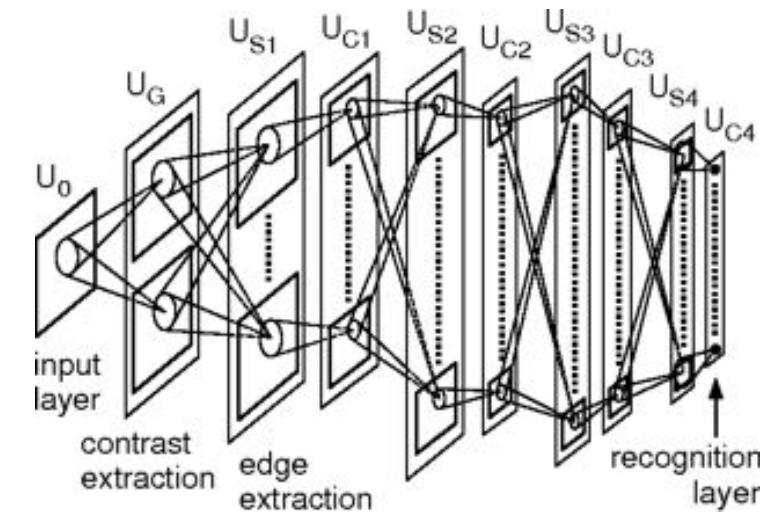
Introduction to self-supervised learning in computer vision

Part: “What”?

The field of AI has made rapid progress, the crucial fuel is data

Algorithms

Deep neural networks



Hardware

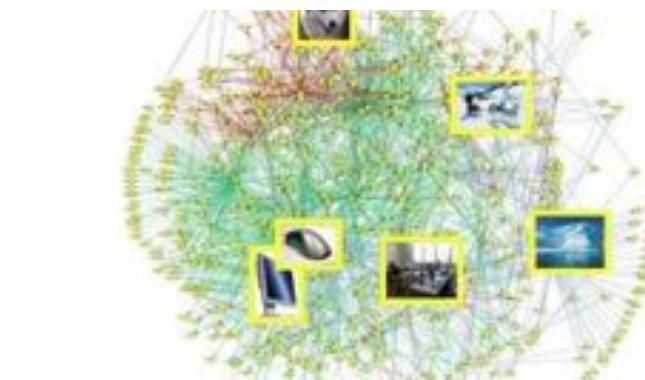
GPUs



Data

Large scale datasets

IM₂GENET



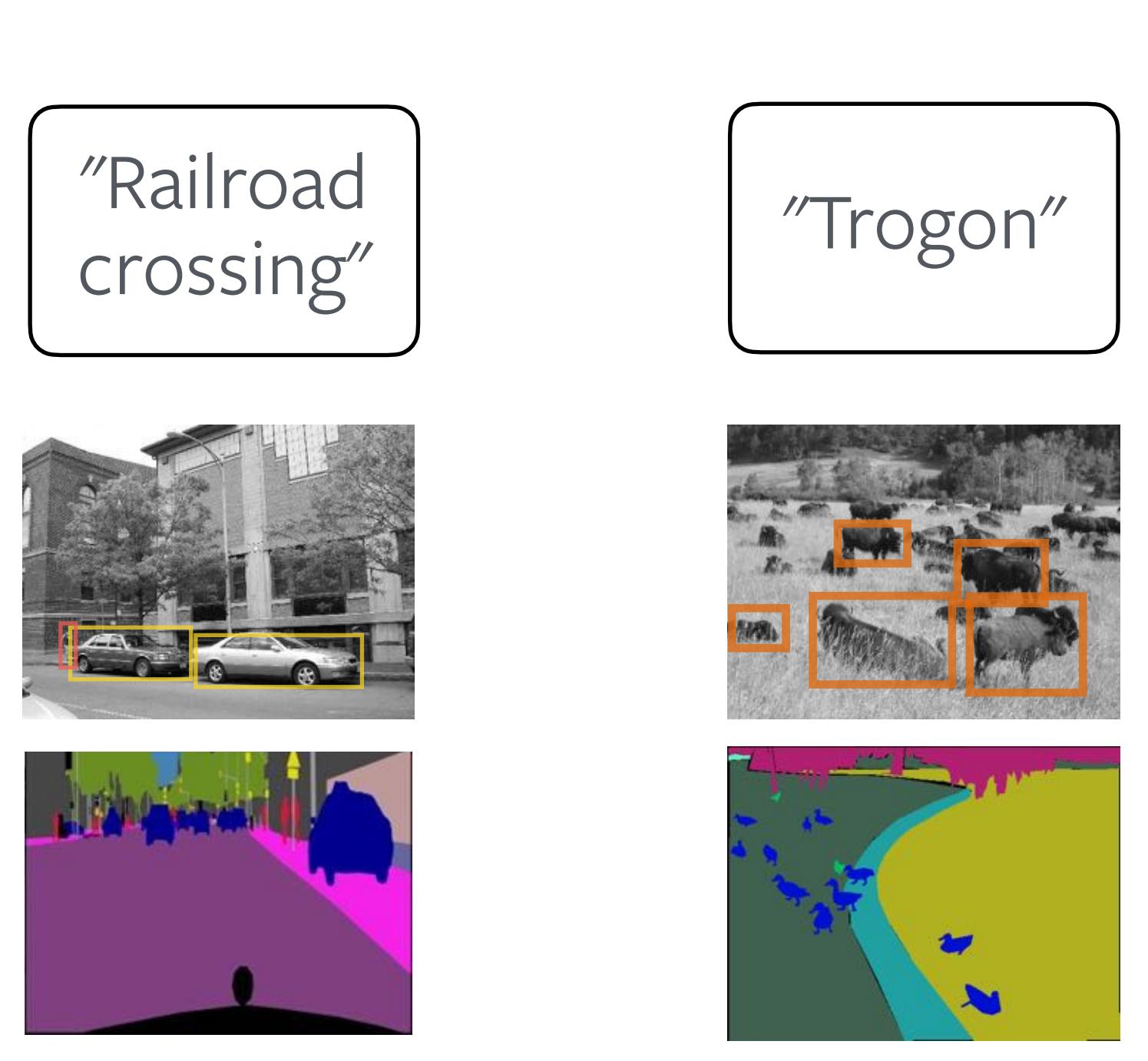
Manual annotations for the data are limiting.

Images are often cheap

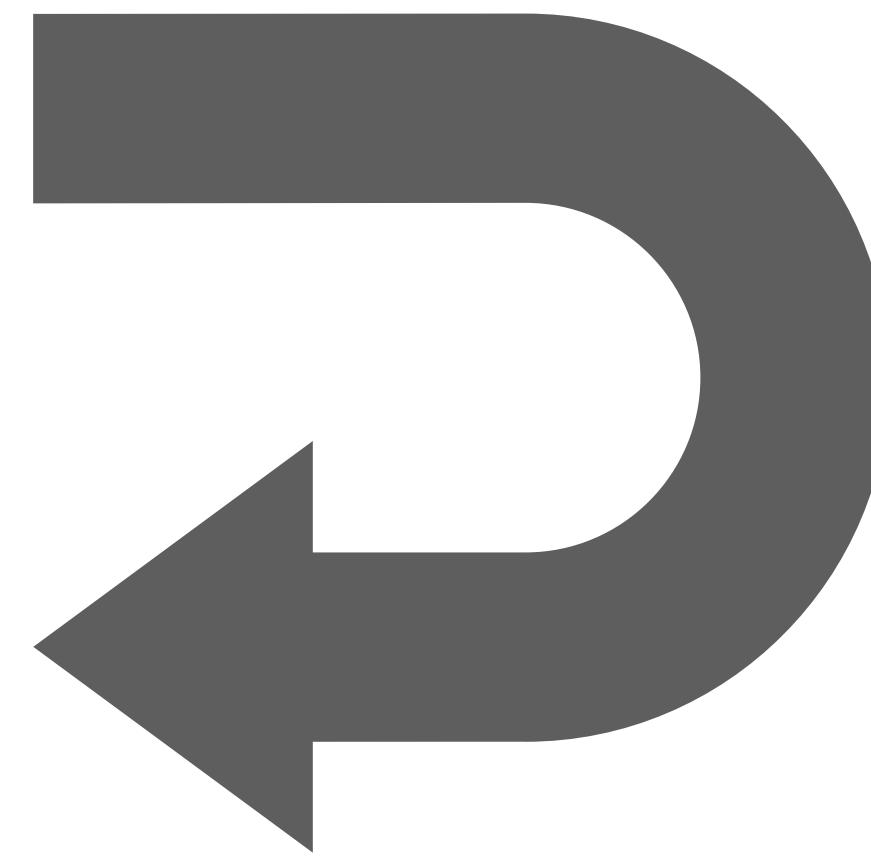


Supervised
Learning

But manual annotations are expensive:
e.g. 30min per image / requiring experts



Solving the problem of expensive annotations: self-supervision.

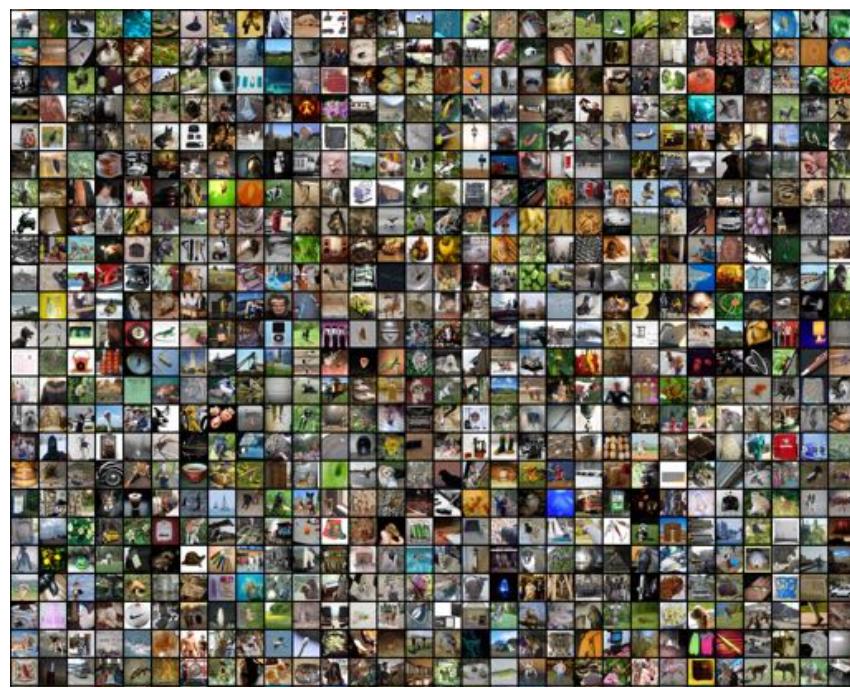


Self-supervision

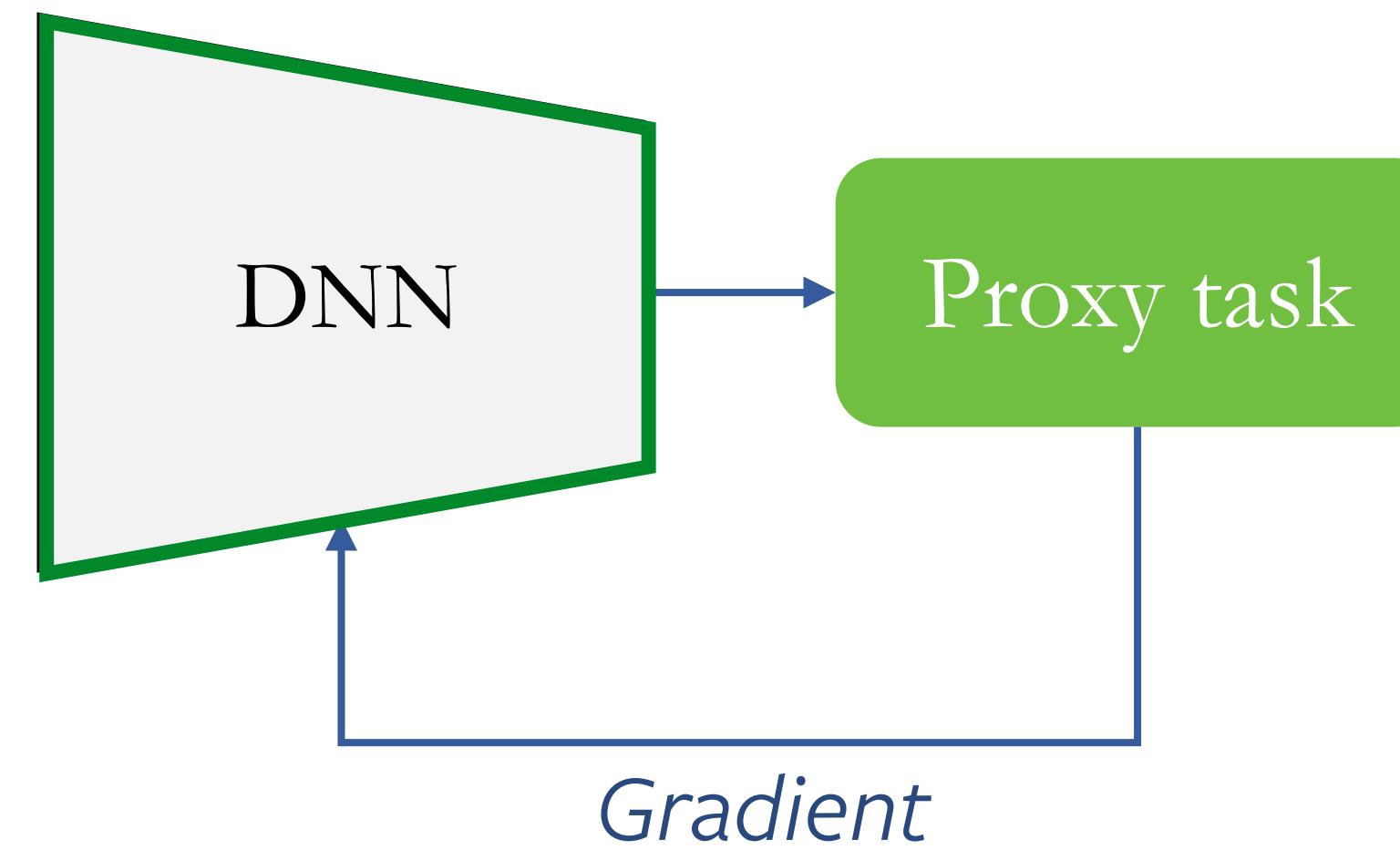
Extract a supervisory signal
from the raw data alone

General procedure of self-supervised learning.

Phase 1: Pretraining



Unlabelled data
+ transformations



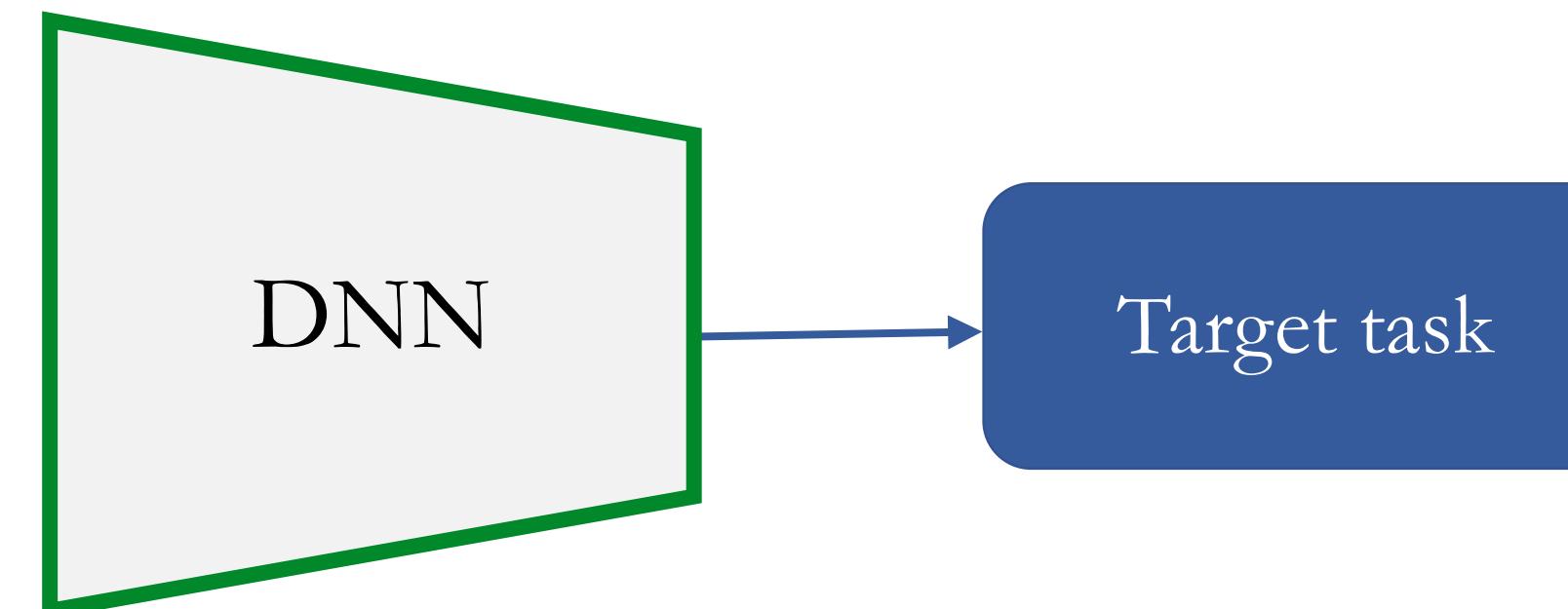
Types:

- Geometry based
- Clustering
- Contrastive
- Generative (partial/full)
- (more)

Phase 2: Downstream tasks



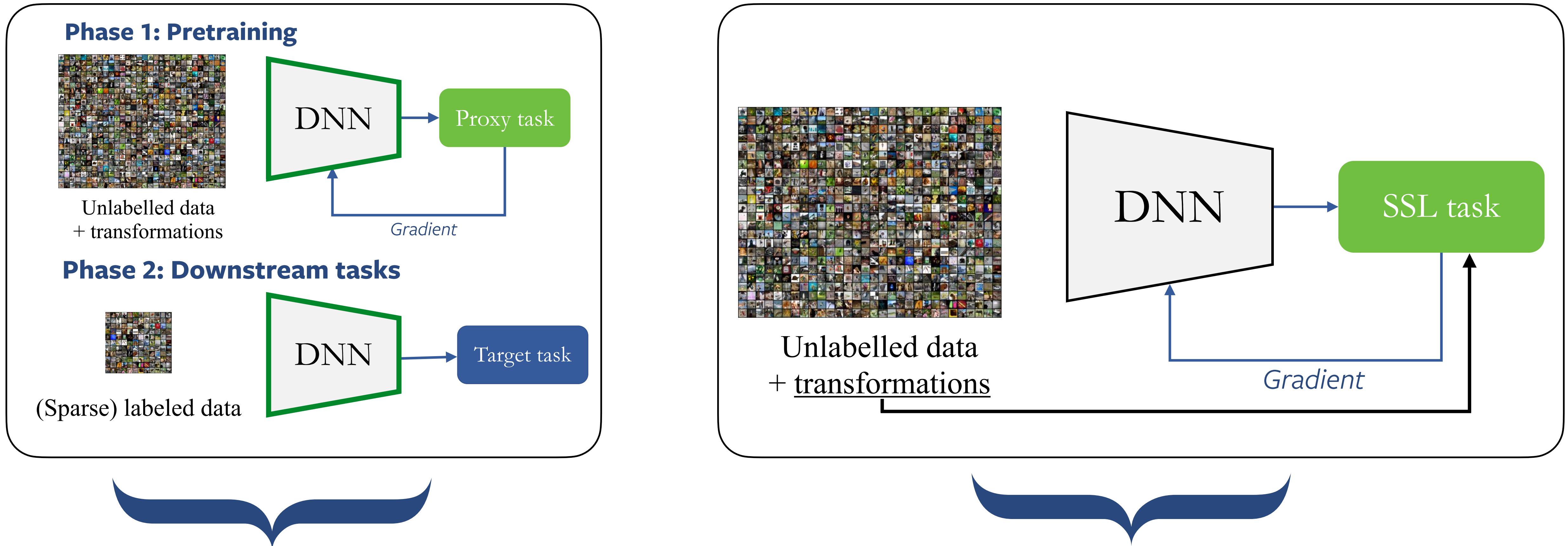
(Sparse) labeled data



Types:

- Limited fine-tuning (e.g. linear layer)
- Finetuning (w/ full or fraction of dataset)

General procedure of self-supervised learning.



Representation Learning

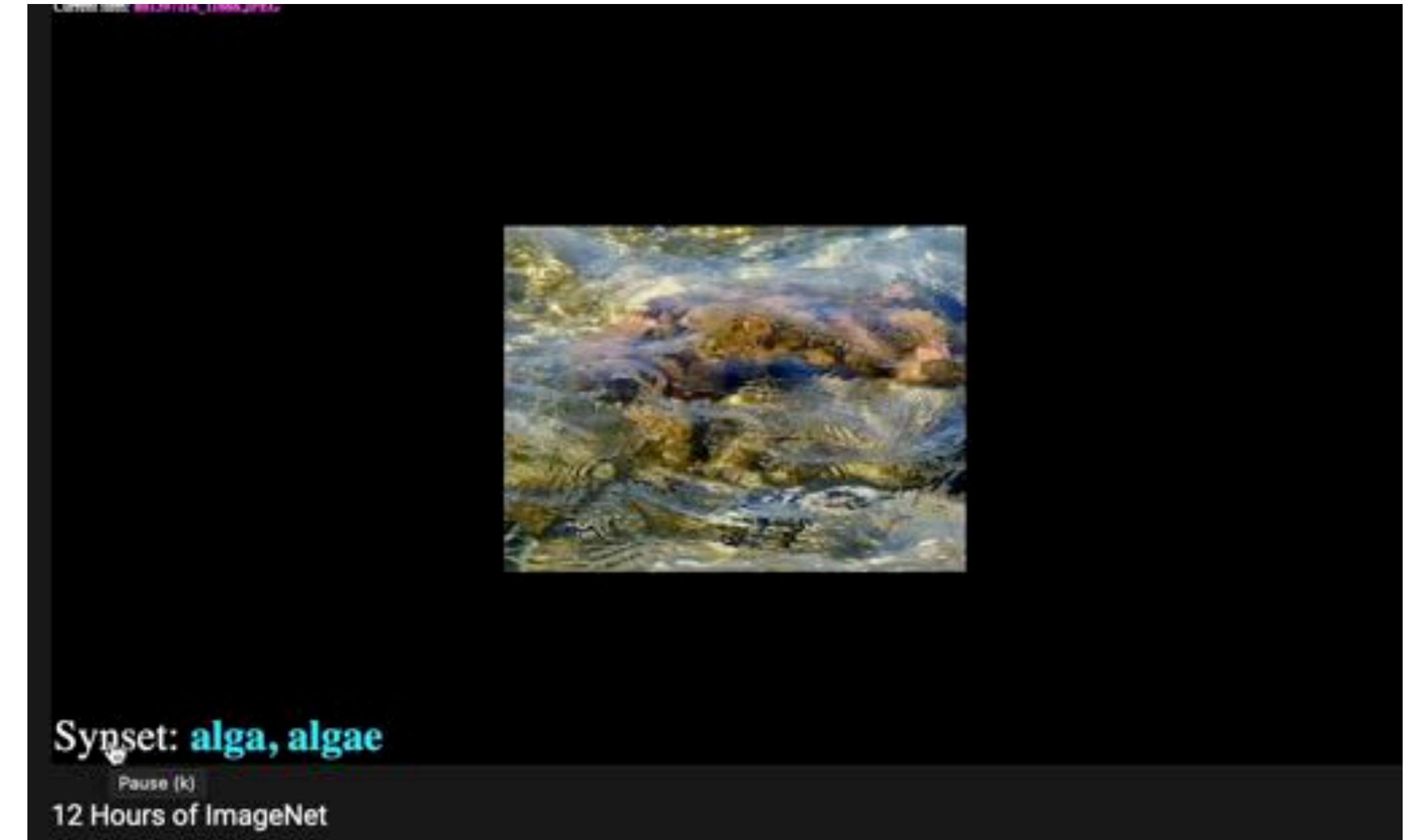
Introduction to self-supervised learning in computer vision

Part: “Why”?

Reason 1: Scalability



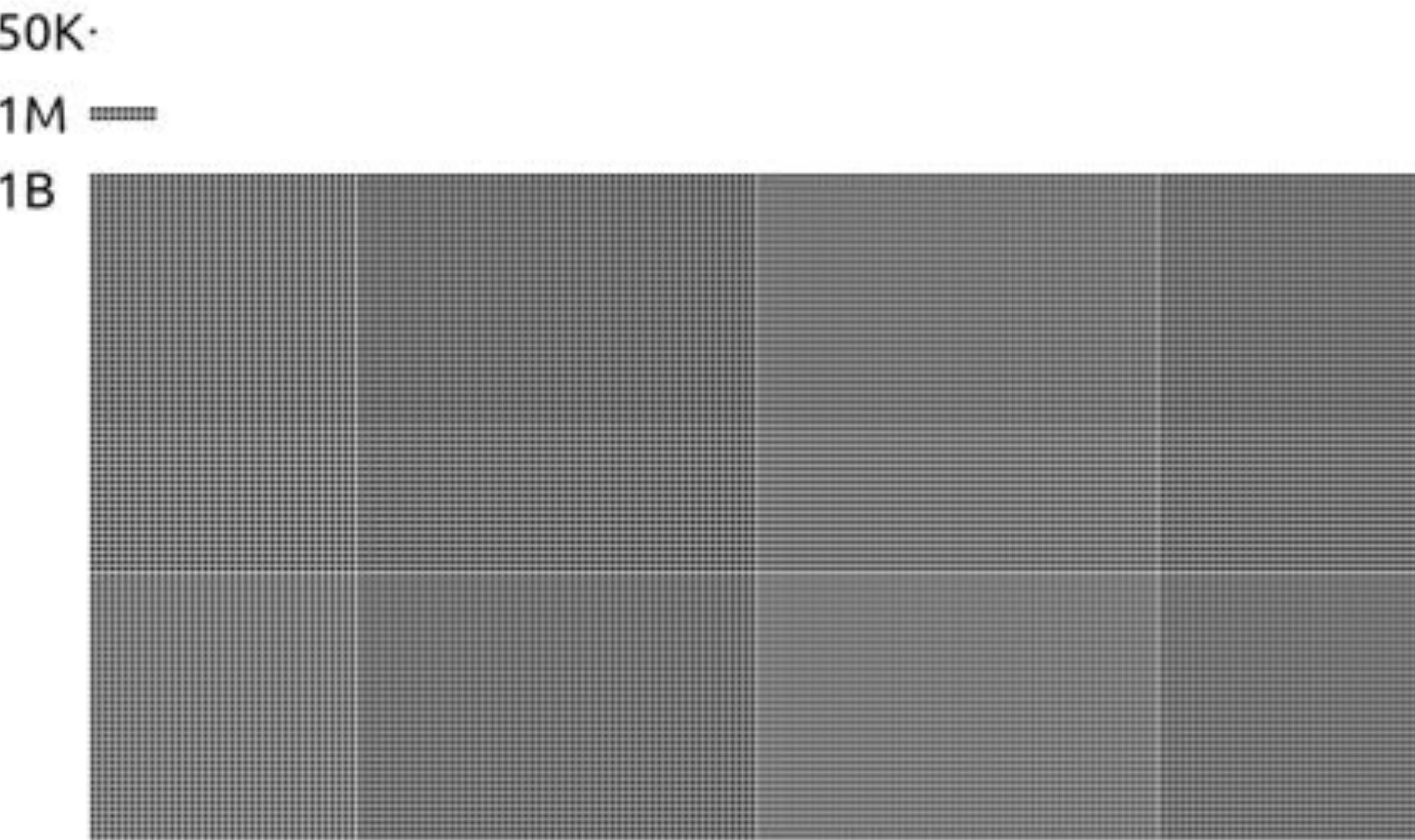
(above) $\times 50 = 1.2M$ images



$90\text{ms} * 1.2\text{M} = 30\text{h}$

Reason 1: Scalability

Instagram: >50B images



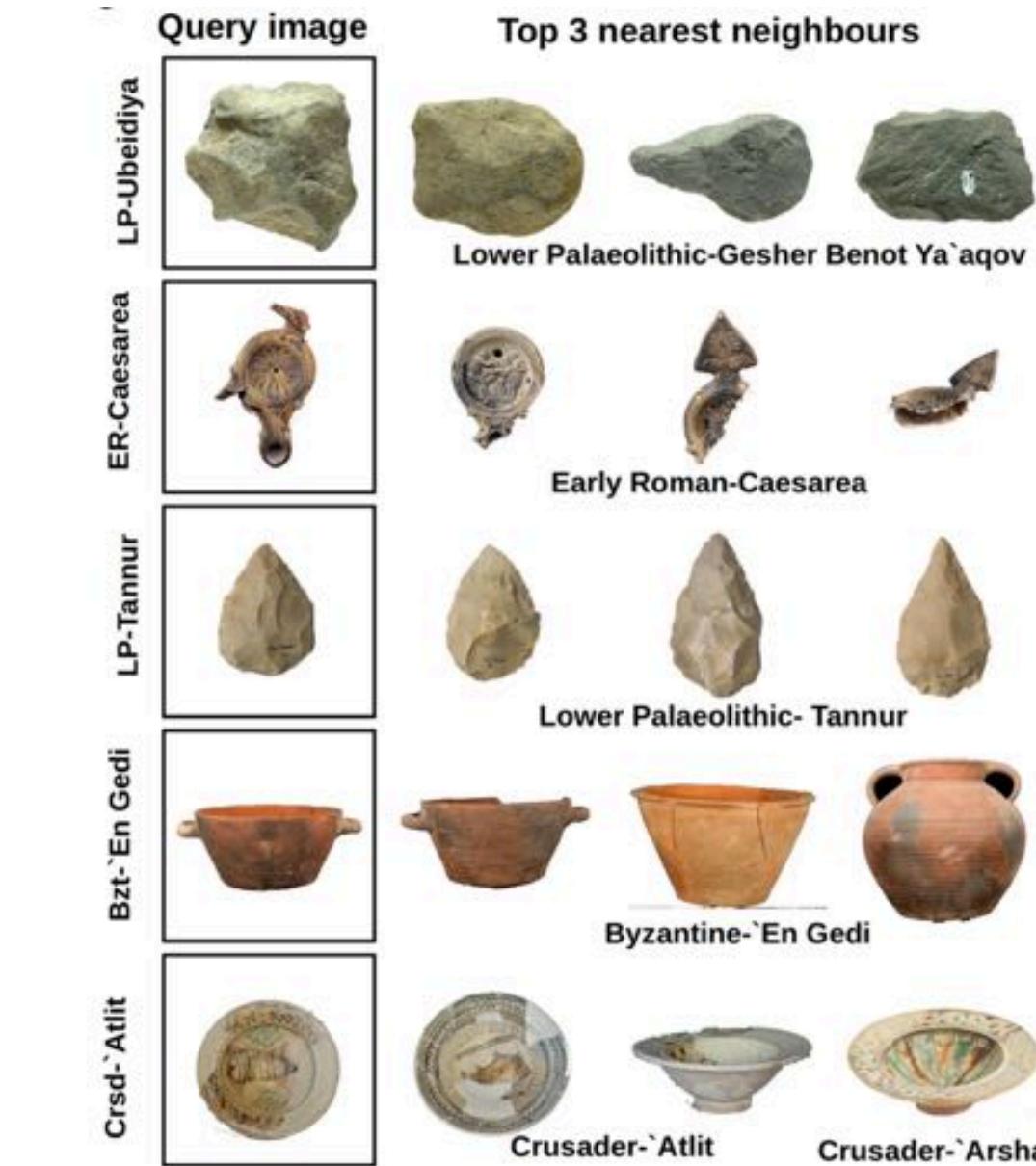
Annotation is expensive, yet datasets keep getting bigger.

Reason 2: Constantly changing domains



Unclear when & what to relabel. Again, large costs just to "keep up".

Reason 2: Accessibility & generalisability



Pretrained models are very useful for a variety of tasks.

Reason 4: Ambiguity of labels



"A house"?



"A boat"?

Bisexual, bisexual person

A person who is sexually attracted to both sexes

- supernumerary (1)
- inhabitant, habitant, dweller, denizen, indweller (4)
- debaser, degrader (1)
- achiever, winner, success, successor (5)
- contemplative (0)
- Cancer, Crab (0)
- national, subject (18)
- interpreter (0)
- numen (0)
- hoper (0)
- gainer (0)
- buster (0)
- biter (1)
- sensualist (12)
 - cocksucker (0)
 - erotic (0)
 - epicure, gourmet, gastronome, bon vivant, epicurean (0)
 - voluptuary, sybarite (0)
 - hedonist, pagan, pleasure seeker (1)
 - playboy, man-about-town, Corinthian (0)
 - bisexual, bisexual person (3)

Nonsensical
visual labels

Labels are ambiguous at best, discriminating and bias-propagating at worst.
Do we really wish to provide our models with these priors?

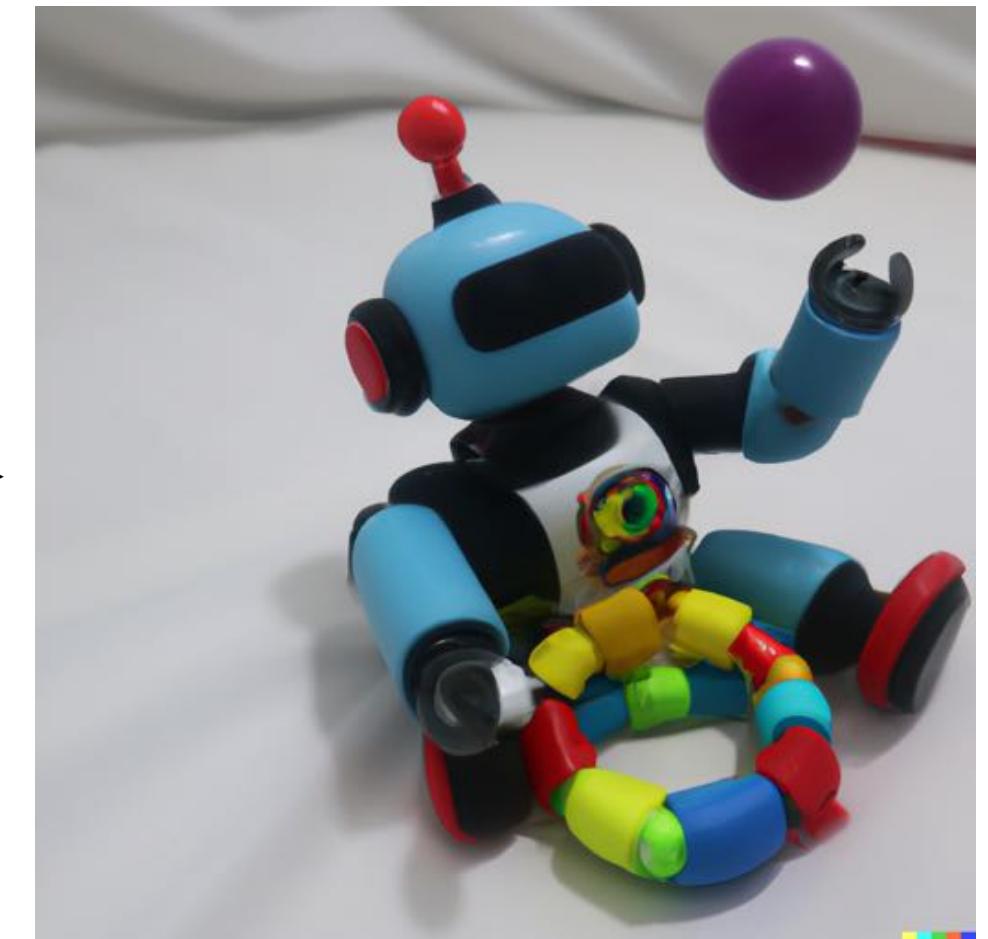
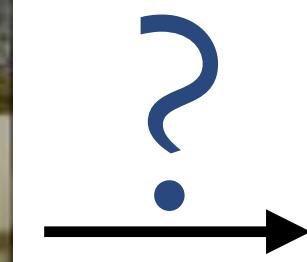
https://en.wikipedia.org/wiki/List_of_house_styles

<https://www.shutterstock.com/image-illustration/flat-ships-sailing-yachts-marine-sailboats-1903407259>

<https://excavating.ai/> Crawford & Paglen

Reason 5: Investigating the fundamentals of visual understanding

The image shows a screenshot of a Meta AI research page. At the top left is the Meta AI logo. To its right are navigation links: Research, Publications, and a partially visible 'P'. Below this, under the heading 'RESEARCH', is the title 'Self-supervised learning: The dark matter of intelligence'. Underneath the title is the date 'March 4, 2021'.



As babies, we learn how the world works largely by observation. We form generalized predictive models about objects in the world by learning concepts such as object permanence and gravity. Later in life, we observe the world, act on it, observe again, and build hypotheses to explain how our actions change our environment by trial and error.

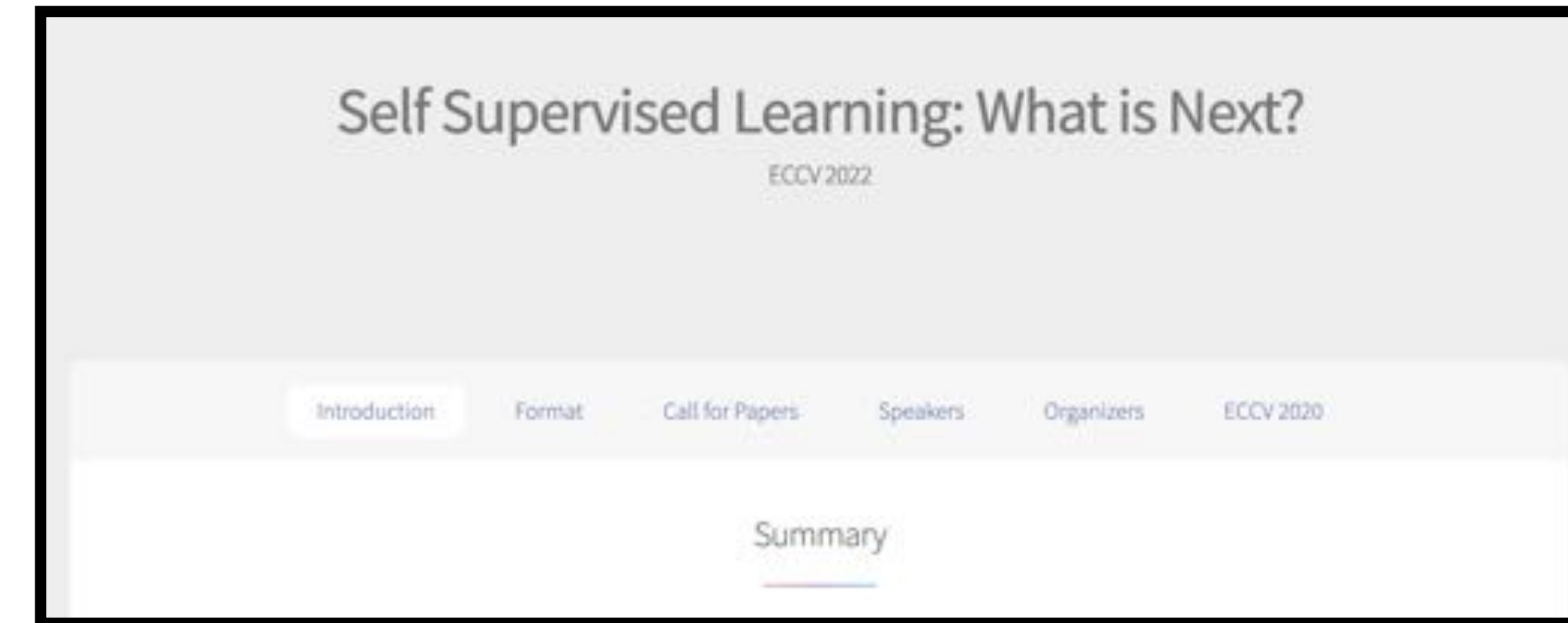
What, if there are, are the limits of learning without labels?

Overview of self-supervised learning methods (the “how”)

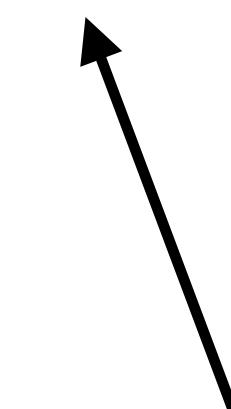
Here, we will only cover a fraction of the existing literature.
Further details and recent developments can be found here:



CVPR'21 Tutorial by Bursuc et al.



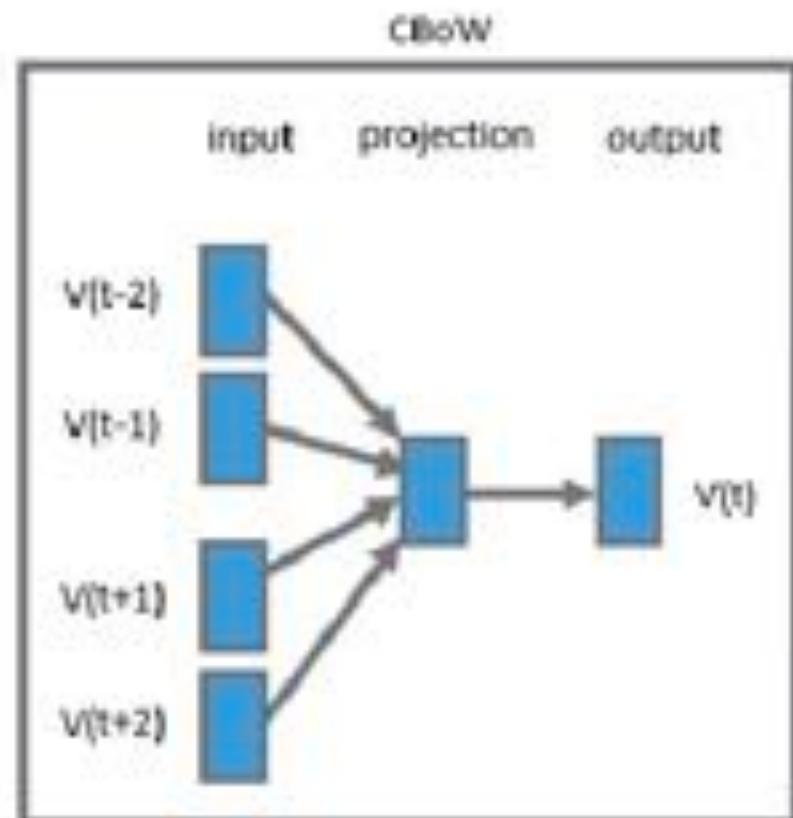
ECCV'20/22 workshop by Asano et al.



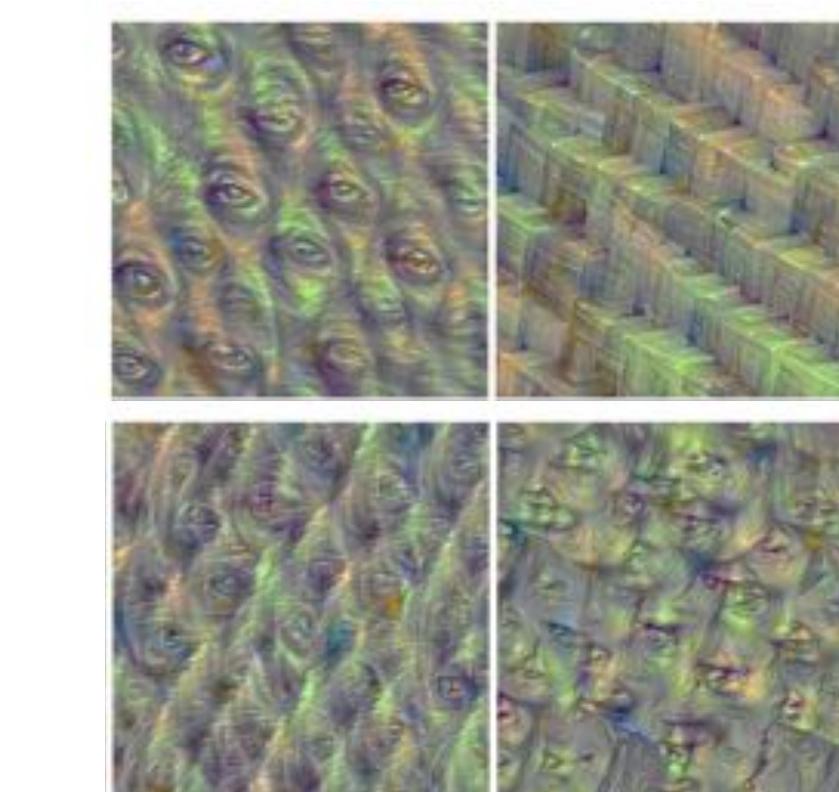
This October!

Early methods

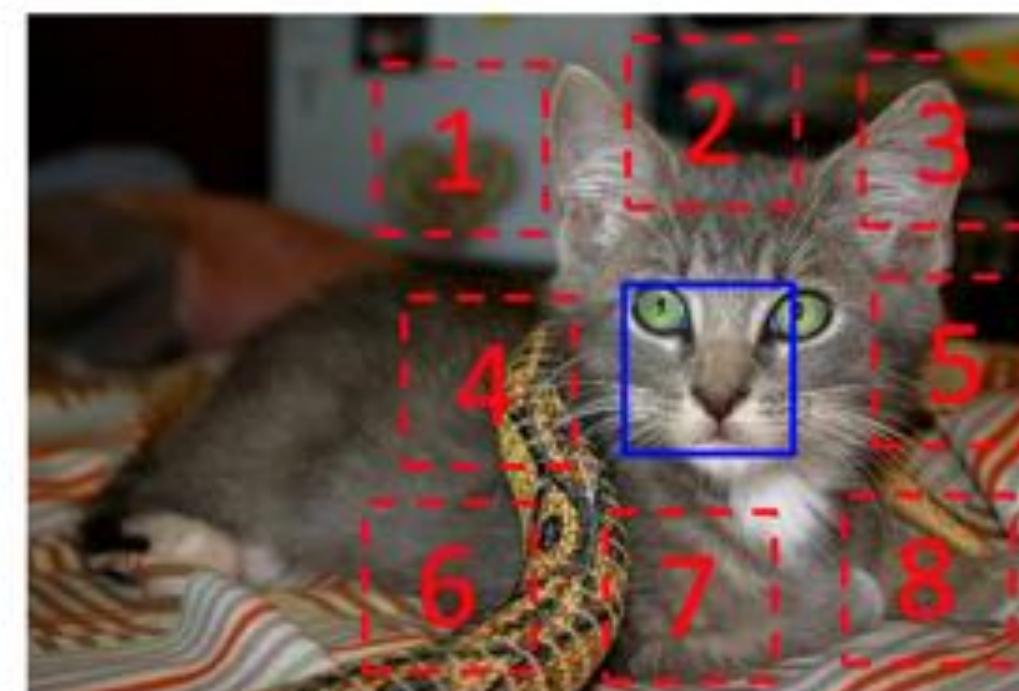
Word2Vec



Motivated from NLP

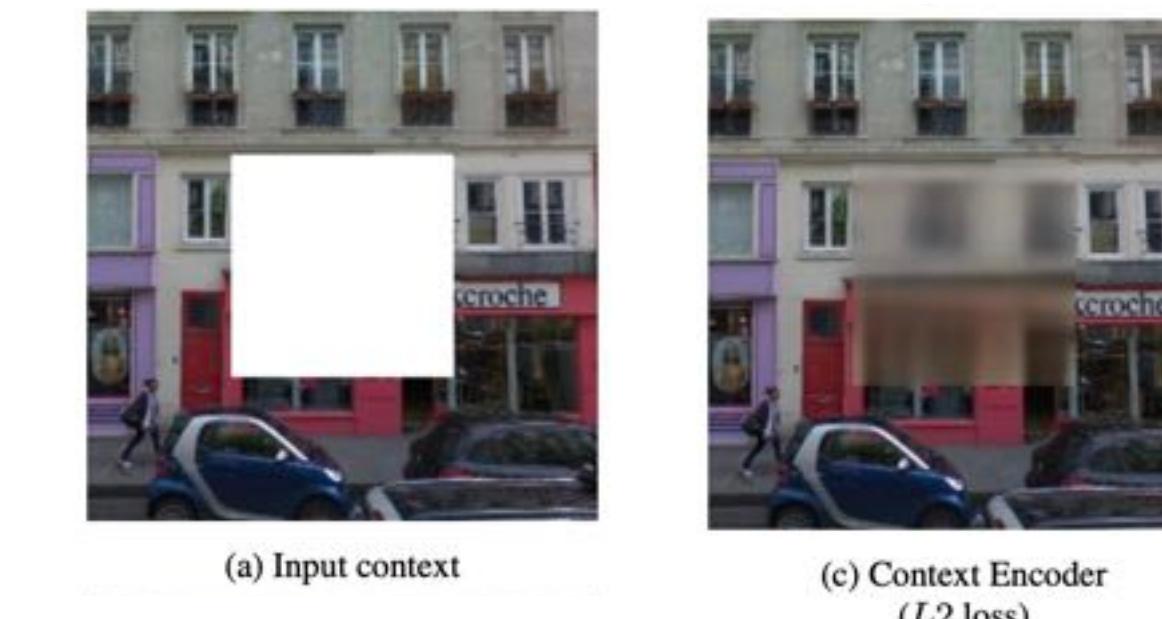


Context Prediction



$$X = (\text{cat face}, \text{background}); Y = 3$$

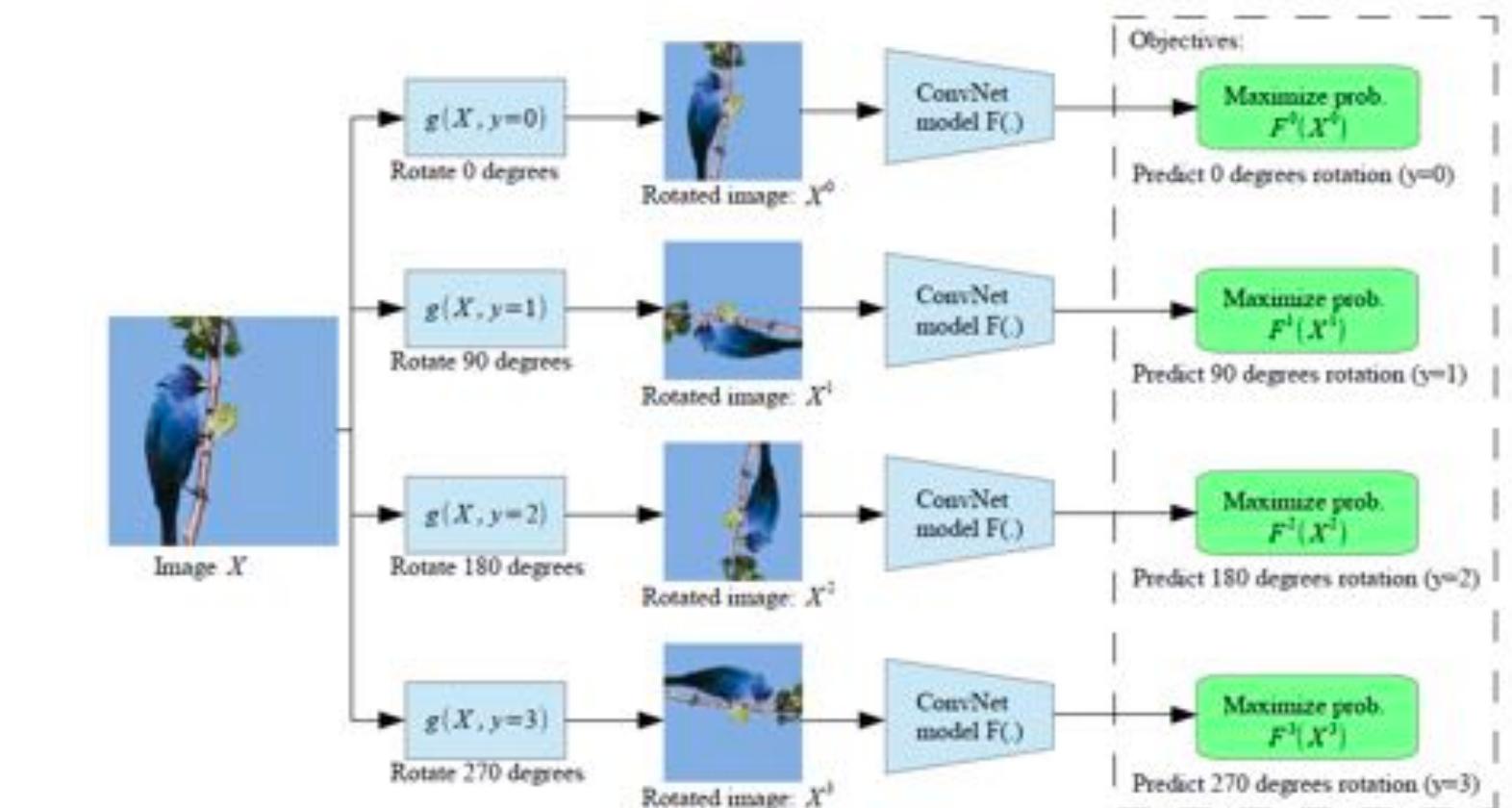
Context Encoders



(a) Input context

(c) Context Encoder
(L_2 loss)

RotNet



Learning without labels is meaningful and possible.

https://www.researchgate.net/figure/Word2Vec-CBOW-and-Skip-gram-There-are-two-different-methods-in-the-Word2Vec-algorithm_fig2_320829283

Doersch et al. *Unsupervised Visual Representation Learning by Context Prediction*. ICCV 2015.

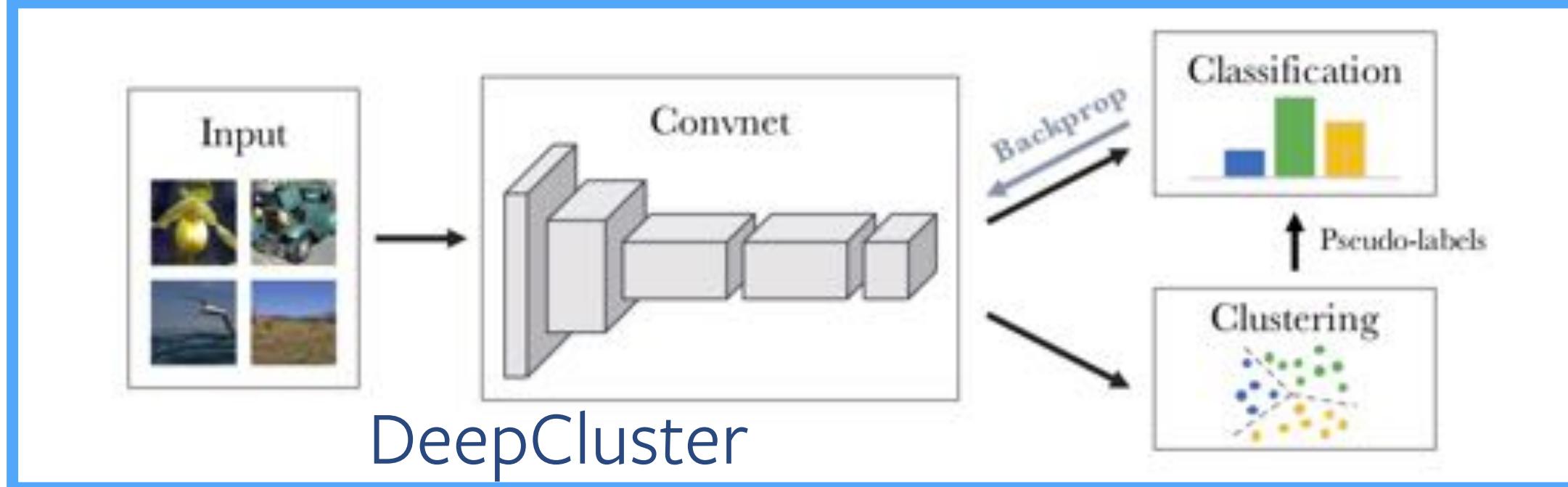
Pathak et al. *Context Encoders: Feature Learning by Inpainting*. CVPR 2016.

Gidaris et al. *RotNet: Unsupervised Representation Learning by Predicting Image Rotations*. ICLR 2018

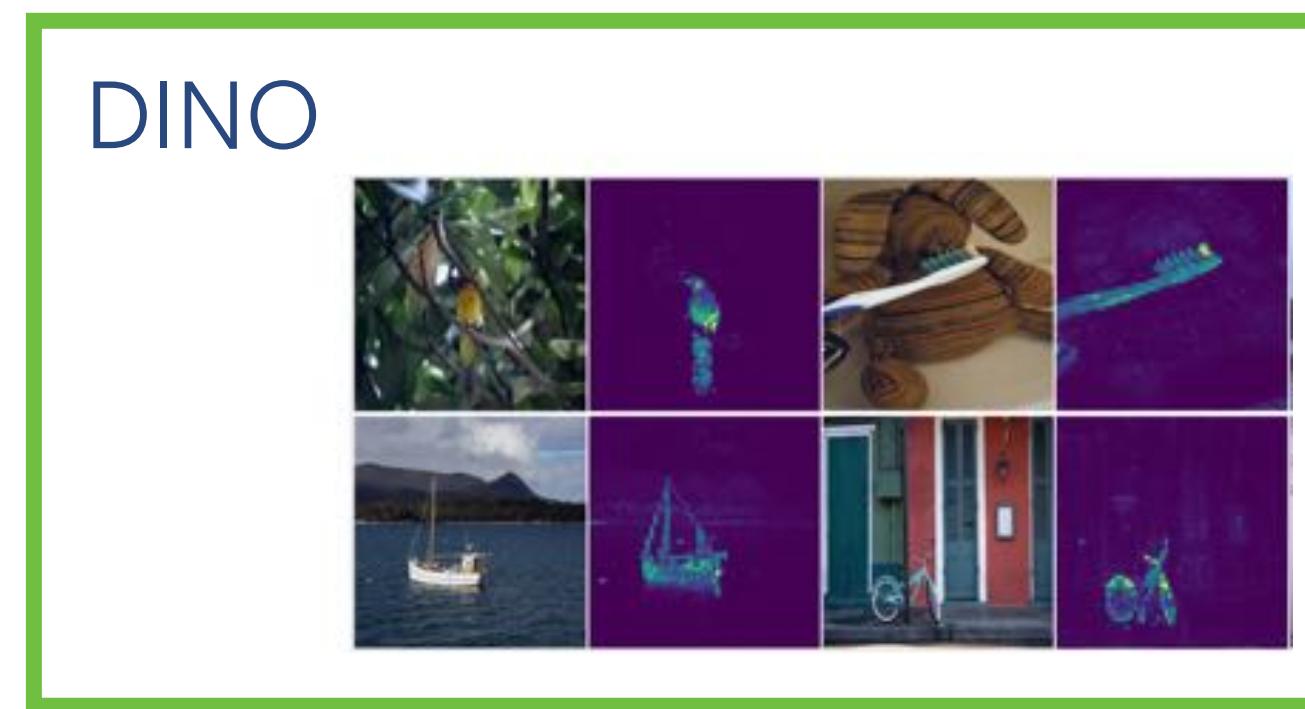
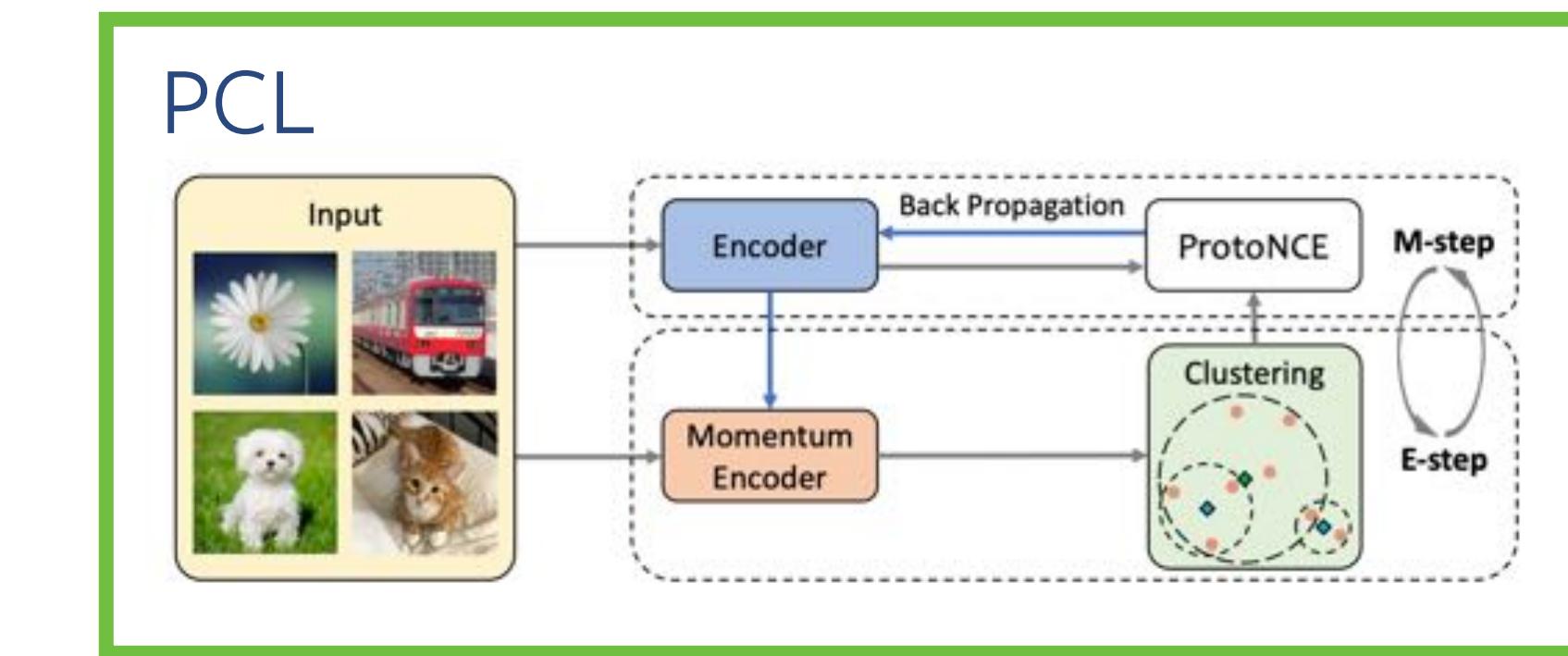
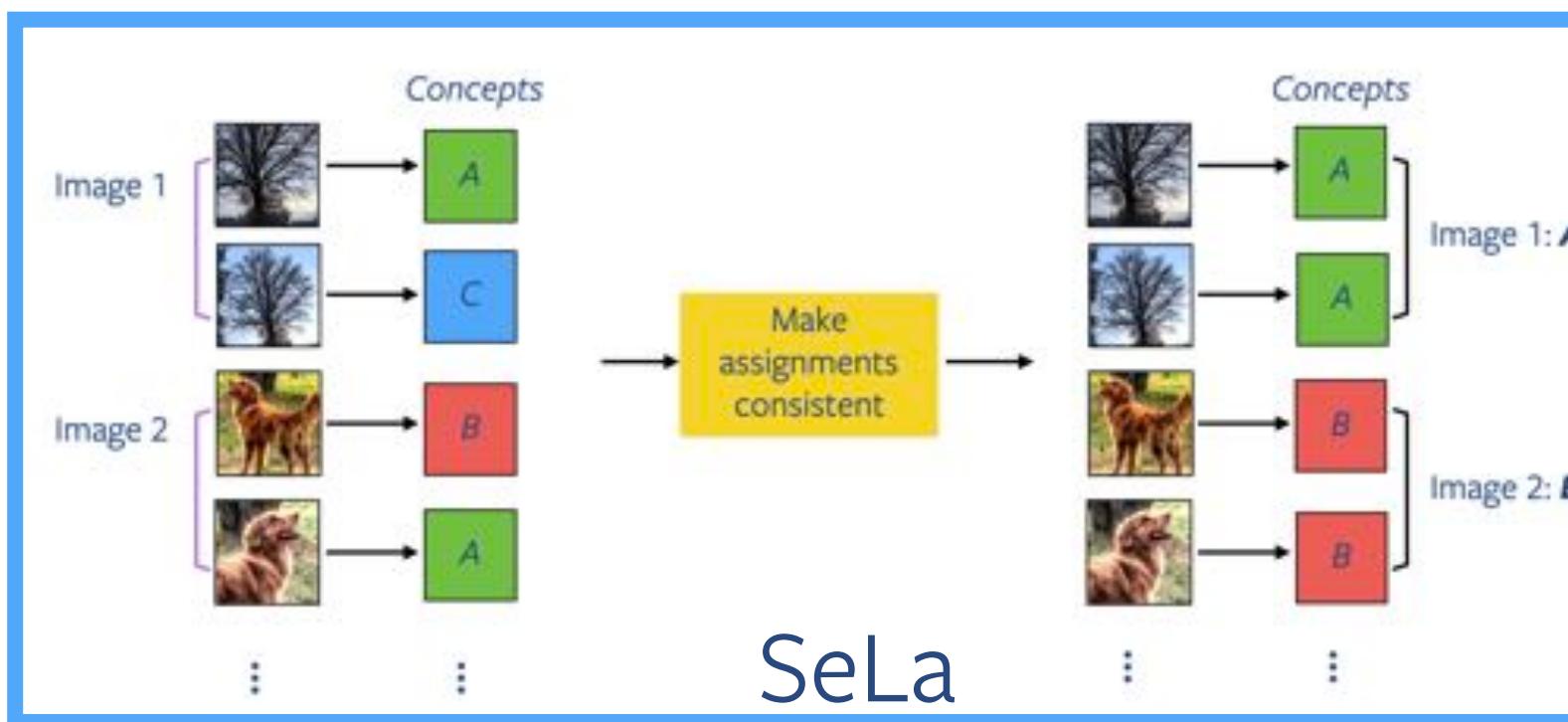
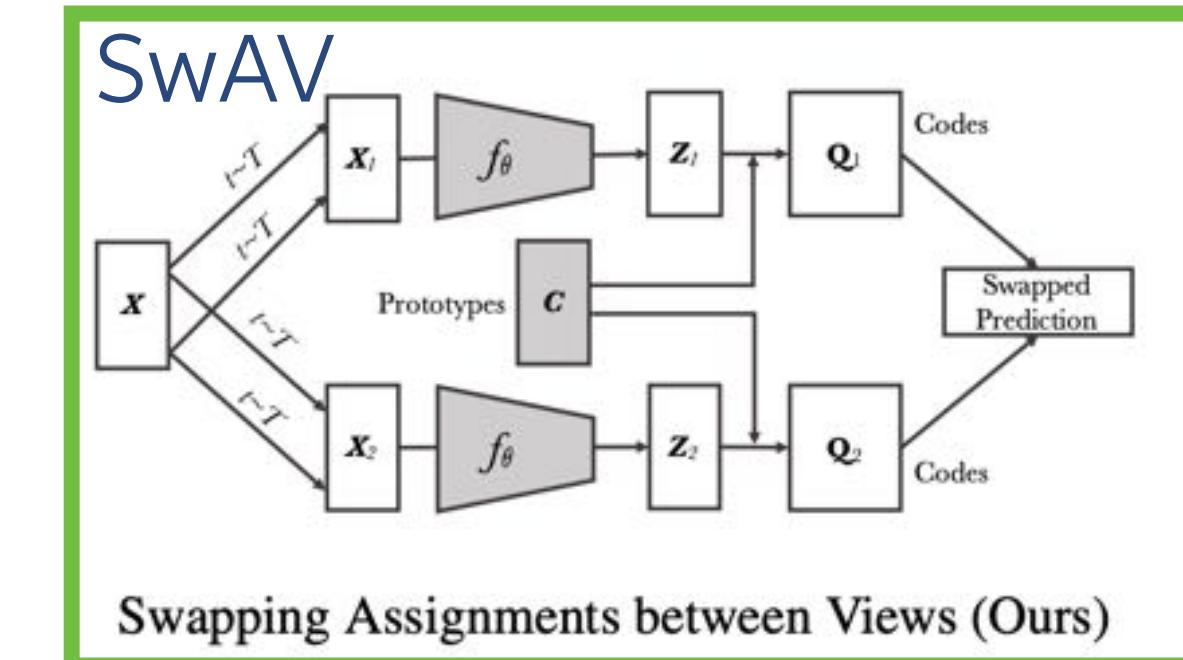


Clustering

Alternate between clustering and network learning



Online & implicit clustering



Clustering is a strong pretext task and serves a useful purpose (~labelling/categorizing).

Caron et al. Deep Clustering for Unsupervised Learning of Visual Features. ECCV'18

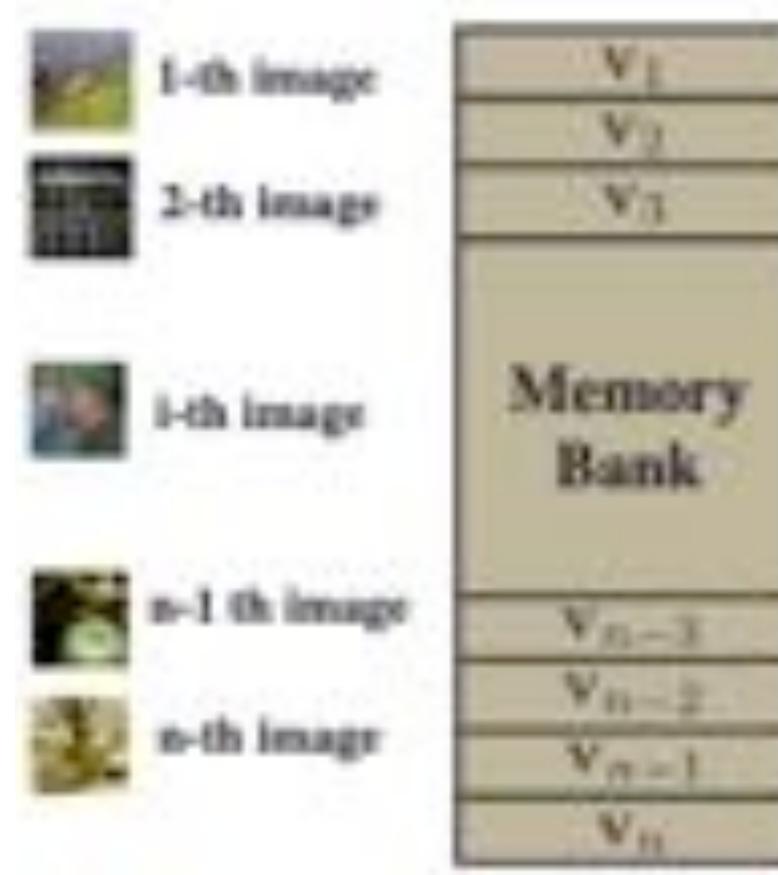
Asano et al. Self-labelling via simultaneous clustering and representation learning. ICLR'19

Caron et al. Unsupervised Learning of Visual Features by Contrasting Cluster Assignments. NeurIPS'20

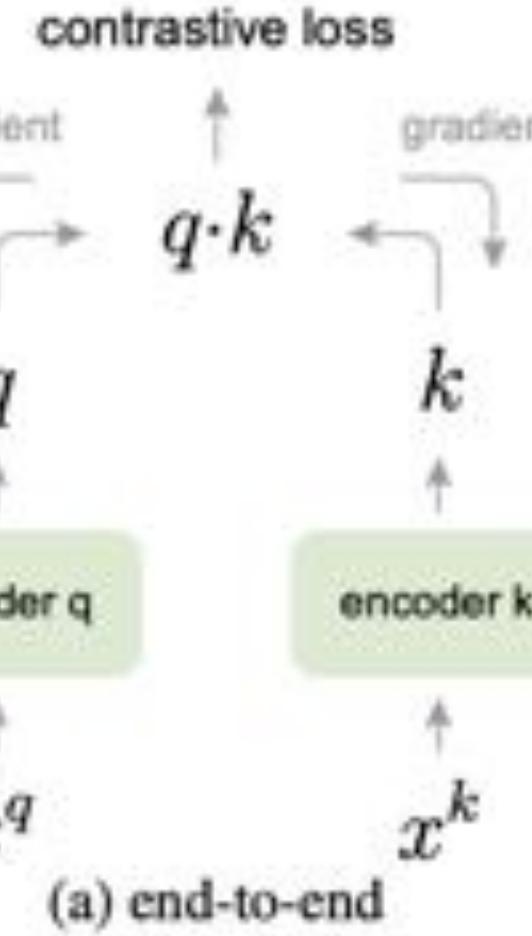
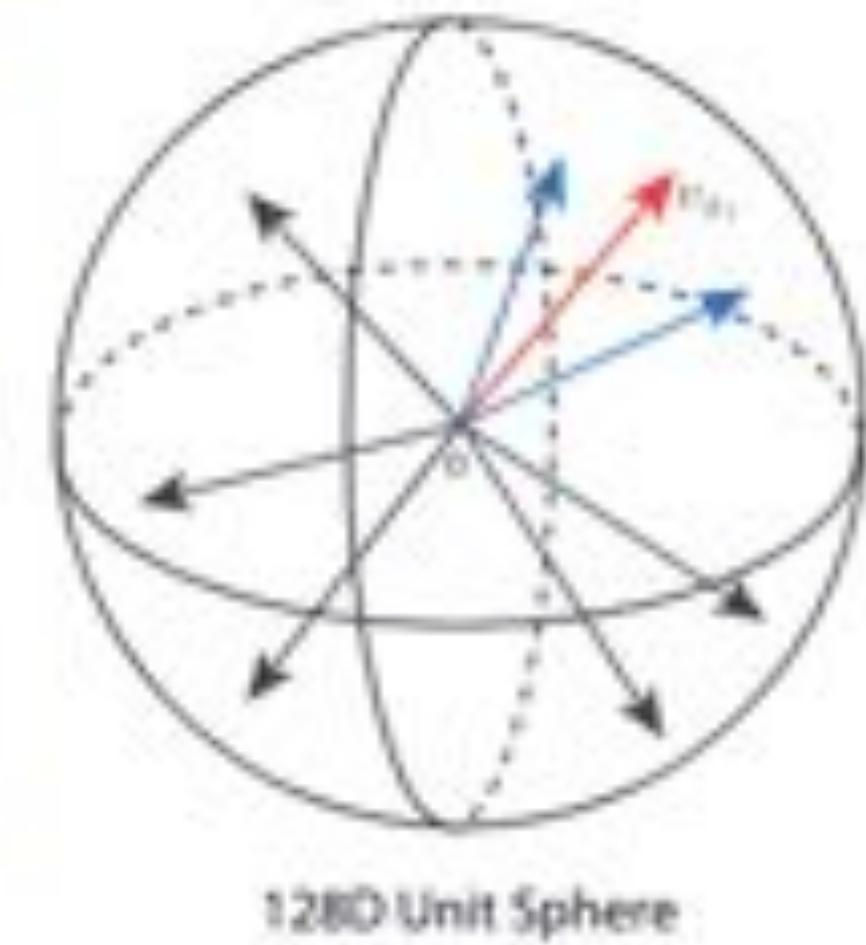
Li et al. Prototypical Contrastive Learning of Unsupervised Representations. ICLR'21

Caron et al. Emerging Properties in Self-Supervised Vision Transformers. ICCV'21

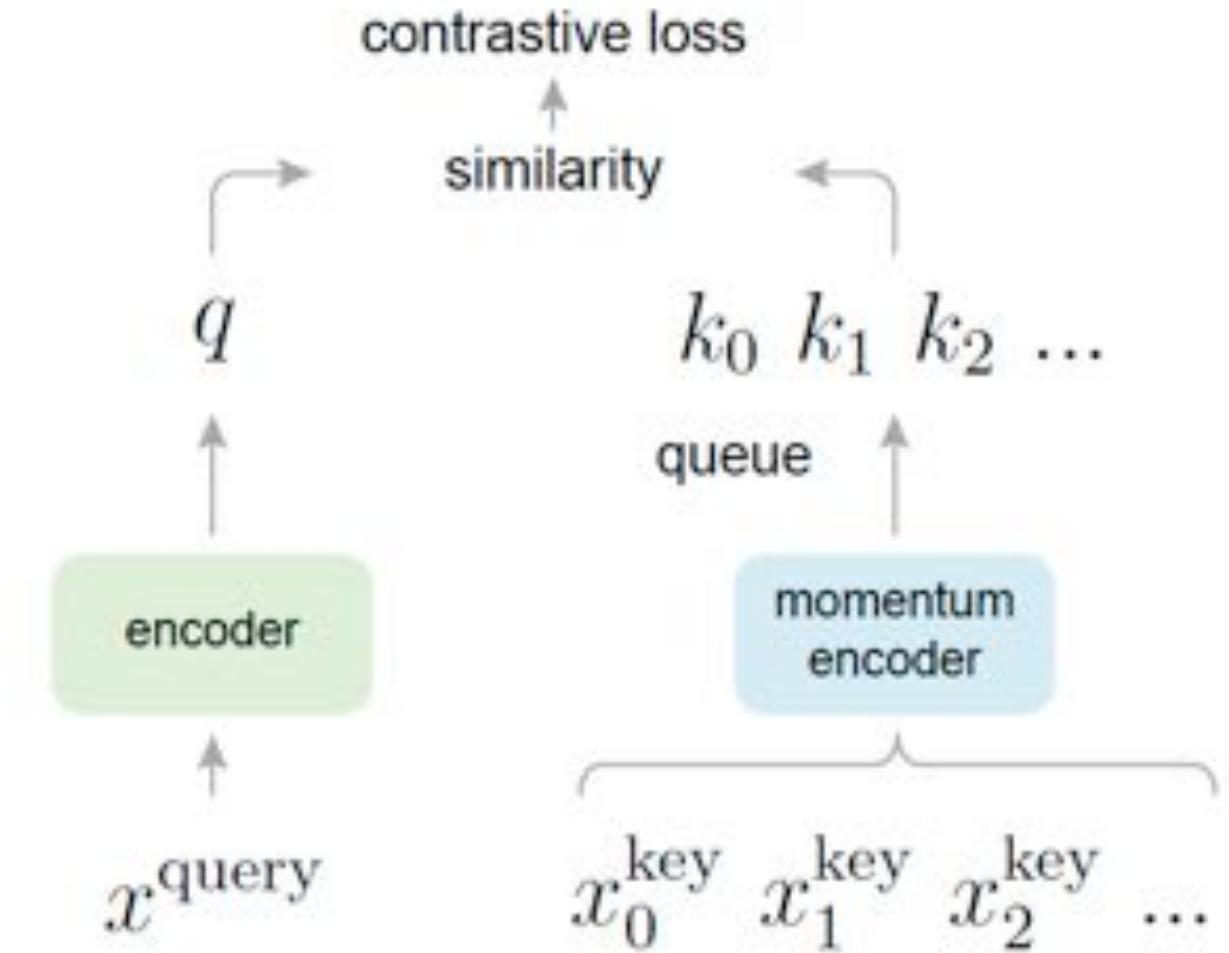
Noise-contrastive self-supervised learning



NPID



SimCLR



MoCo

Momentum encoder:

```
# momentum update: key network  
f_k.params = m*f_k.params + (1-m)*f_q.params
```

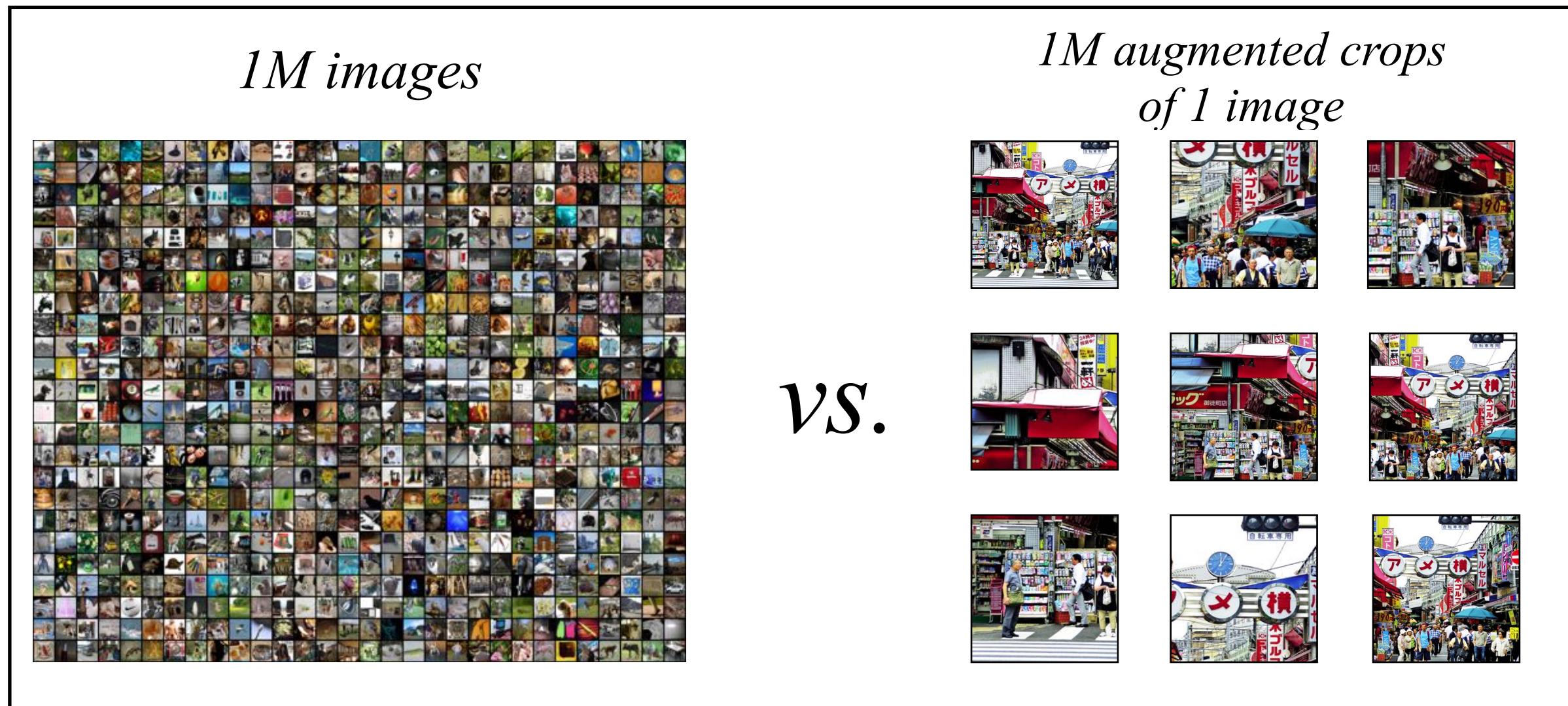
The start of large-scale & industrial self-supervised learning.
These works heavily rely on image augmentations.

How does one implement contrastive learning?

Learn it at the tutorial!



About image-augmentations for self-supervised learning.



Downstream task: object detection,
COCO R50-C4 finetuning, 1x

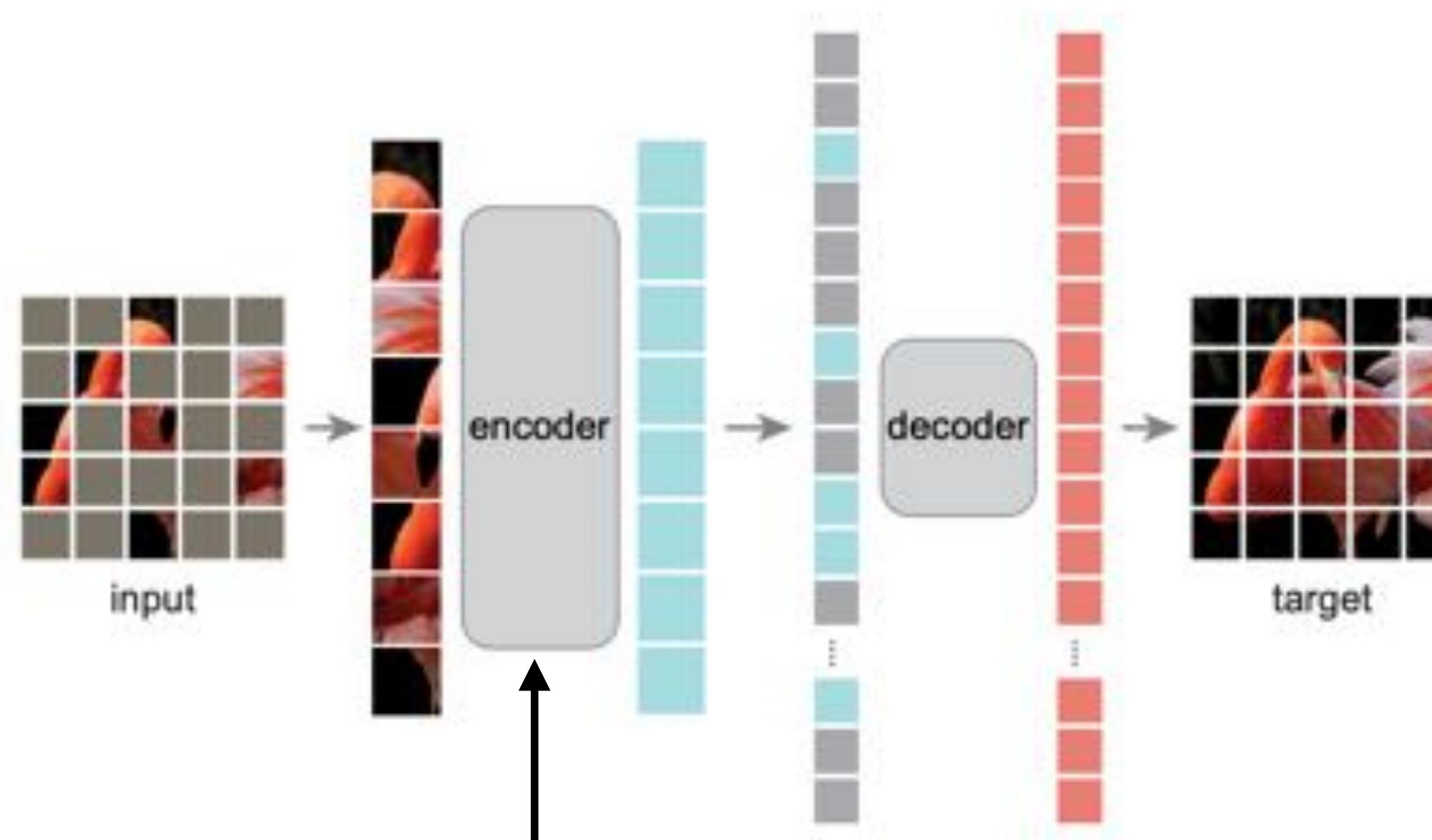
pre-train	Bounding-box			Segmentation		
	AP ^{bb}	AP ₅₀ ^{bb}	AP ₇₅	AP ^{mk}	AP ₅₀ ^{mk}	AP ₇₅ ^{mk}
Random	26.4	44.0	27.8	29.3	46.9	30.8
Supervised	38.2	58.2	41.2	33.3	54.7	35.2
ours 1-image A	36.5	55.2	39.2	32.1	52.2	34.0
MoCo-v1	38.5	58.3	41.6	33.6	54.8	35.6
MoCo-v2	39.0	58.6	41.9	34.2	55.4	36.2

Method, Image A		
BiGAN	RotNet	DeepCluster
20.4	19.9	20.7

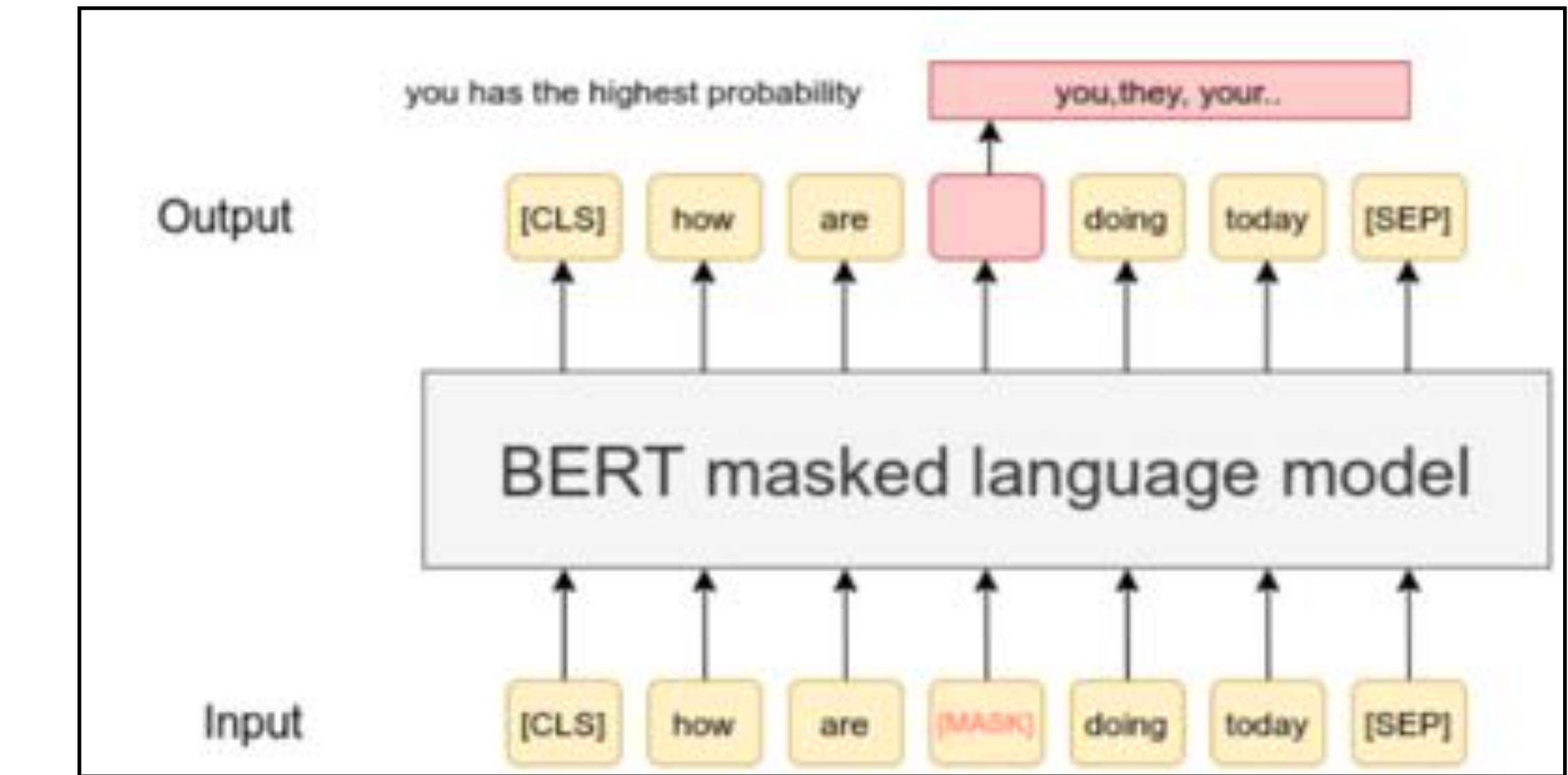
Self-supervised learning's "learning": mostly from augmentations.

+10% mAP from a single image and augmentations
Within 3% of MoCo-v2 on full ImageNet

Masked Image Modelling



Vision Transformer



Back to NLP

Datasets: Pretraining and downstream



- Class-balanced dataset, via search engine
- Unclear image licences
- Particular choice of classes, e.g. 120 classes of dogs
- Object-centric, stereotypical images
- Many problematic images (see Prabhu & Birhane)



- Random images from YFCC-100M
- All images with complete CC-BY licences
- No people, nor identifiable information
- Natural images as humans take them
- Likely a better indicator for billions-level pretraining

Clustering (❄️)

IN-1k
ObjectNet
Places205
Flowers

SVM low-shot (❄️+🔥)

Places205
Pascal VOC
Herbarium-19

Linear probing (❄️+🔥)

IN-1k
Places205
CIFAR-100
Flowers
...

Finetuning (🔥)

MS-COCO:
detection, segmentation, key point detection, dense pose estimation

Pascal VOC:
detection

LVIS v1.0:
detection



Generalizable representations and ‘free labelling’ of images

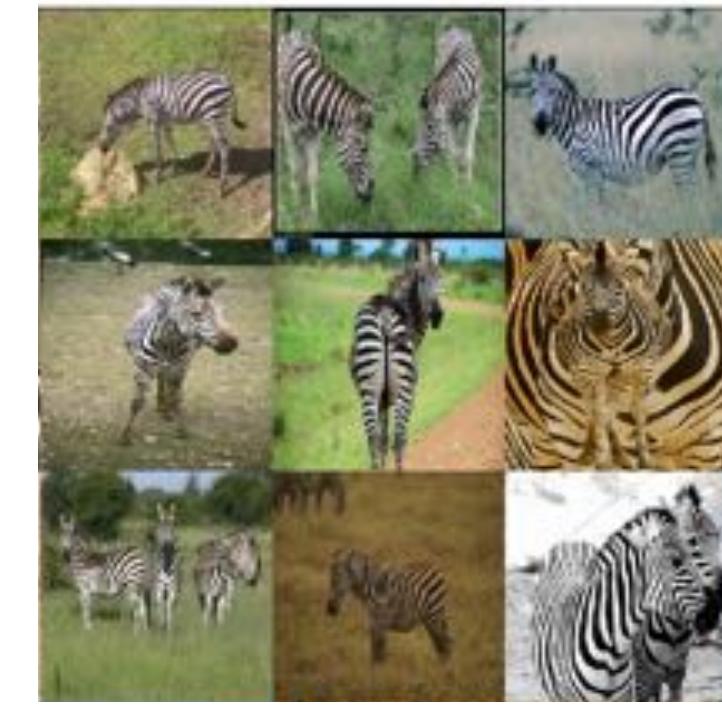
Self-labelling via simultaneous clustering and representation learning (ICLR'20 spotlight)

YUKI M. ASANO, CHRISTIAN RUPPRECHT, ANDREA VEDALDI

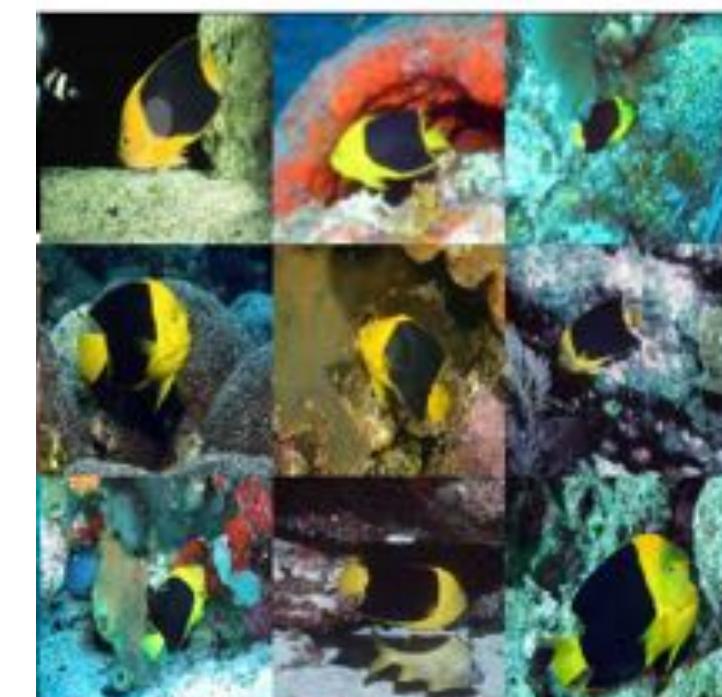
Goal: Discover visual concepts without annotations.



(above) $\times 50 = 1.2\text{M}$ images



concept "A"



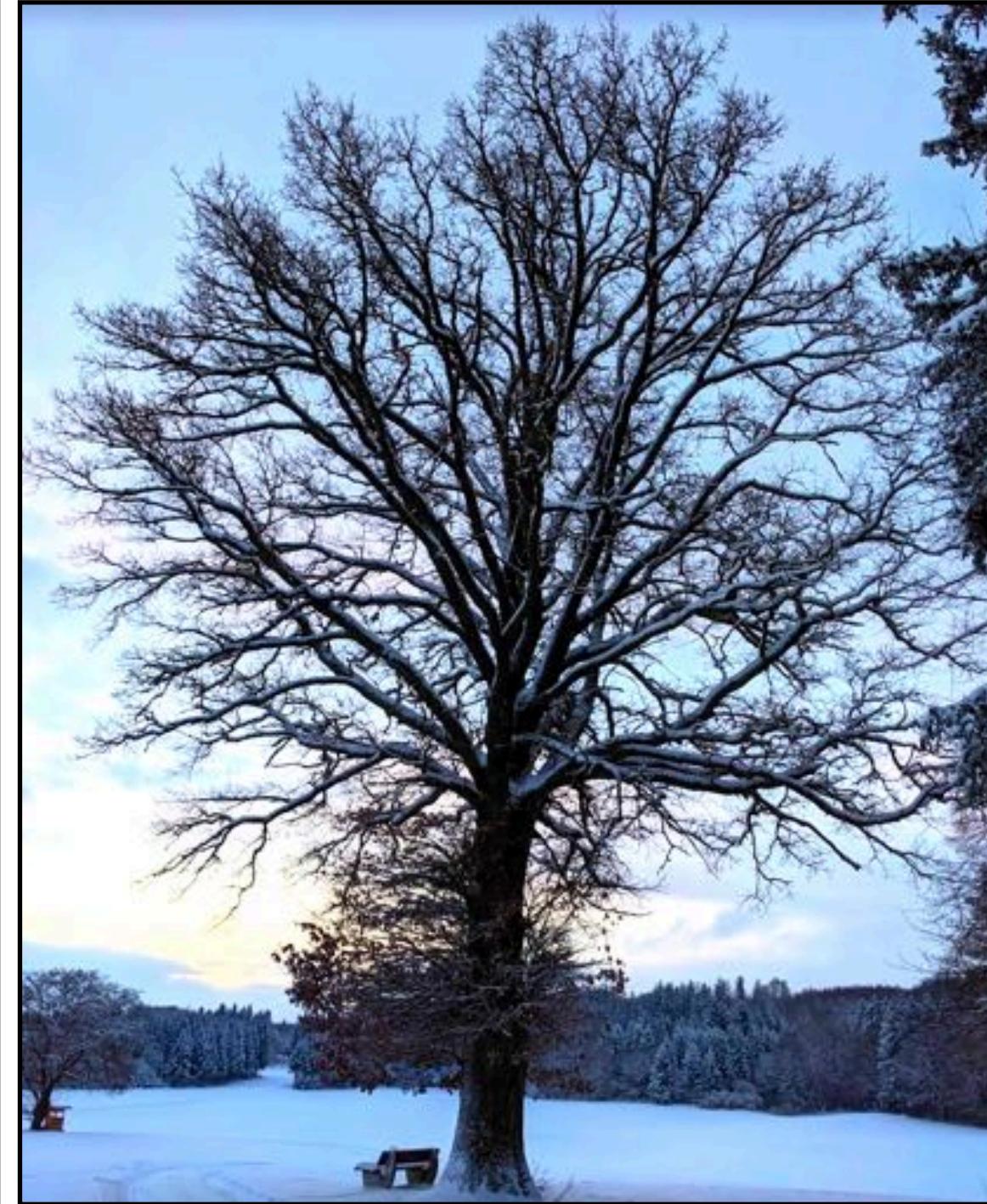
concept "Z"



The key to image understanding is separating meaning from appearance.



Original



Different lighting



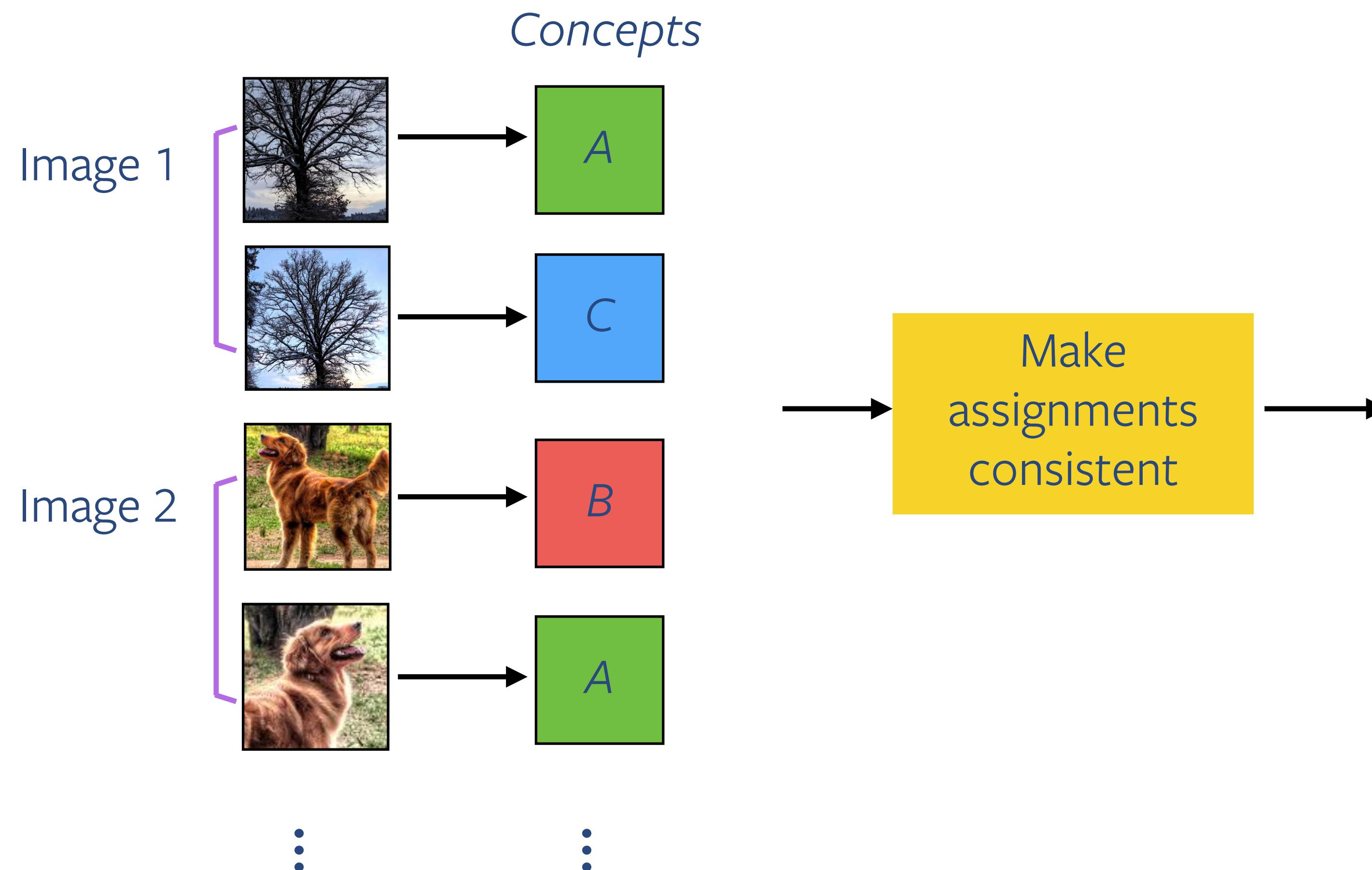
Mirrored



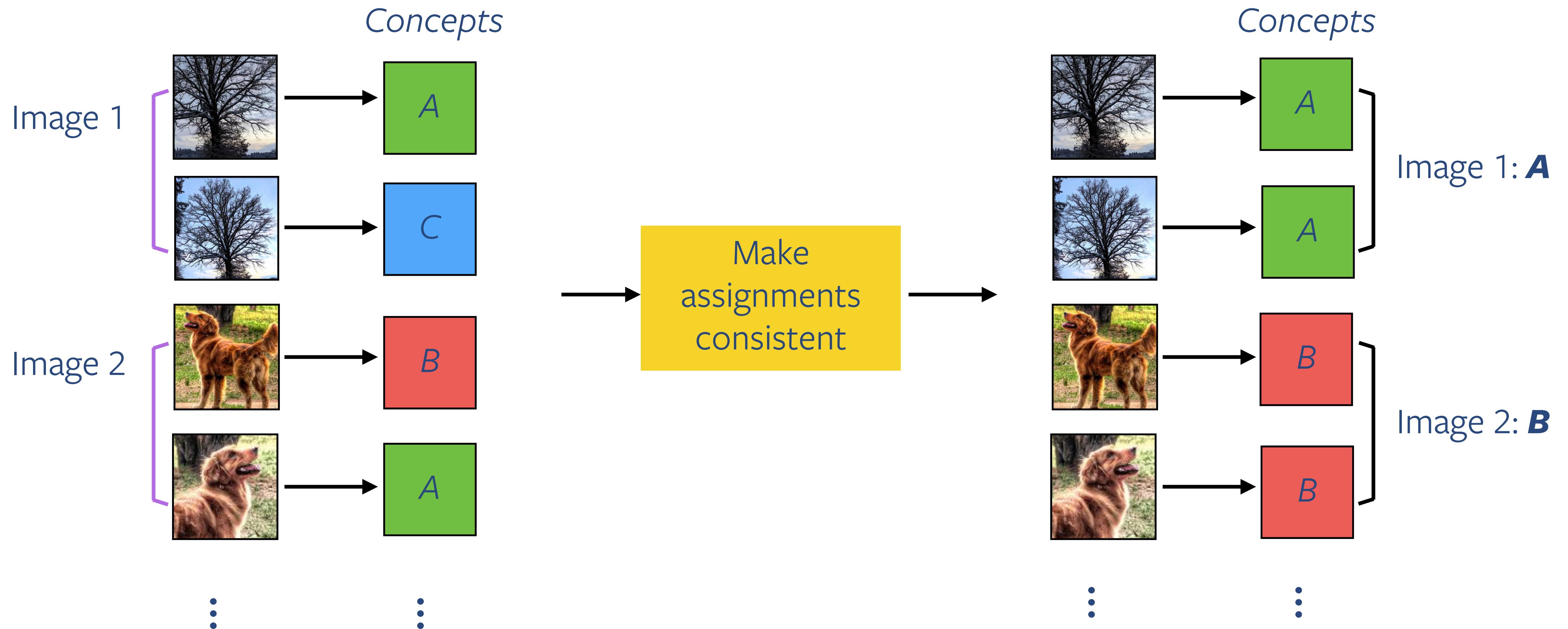
Different zoom

Augmentations

Our work applies the idea of augmentation invariance to assign concepts.



Our work applies the idea of transformation invariance to assign concepts.



How can we optimize the labels and make assignments consistent?

If we had ground-truth labels

$$\min_{\Phi} L(\mathbf{y}, \Phi),$$

$$\text{where } L(\mathbf{y}, \Phi) = \frac{1}{N} \sum_{i=1}^N \log p(y_i | \mathbf{x}_i, \Phi)$$

- L is the loss (cost) function
- Φ is the deep neural network model
- y are the labels

Our novel contribution *without* ground-truth

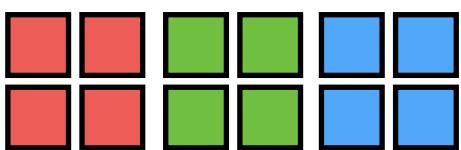
Solution sketch:

1. Represent via an assignment table q and optimize:

$$L(q, \Phi) = \frac{1}{N} \sum_{i=1}^N \sum_y q(y | \mathbf{x}_i) \log p(y | \mathbf{x}_i, \Phi)$$

But: The trivial solution for q is to set all labels to be the same

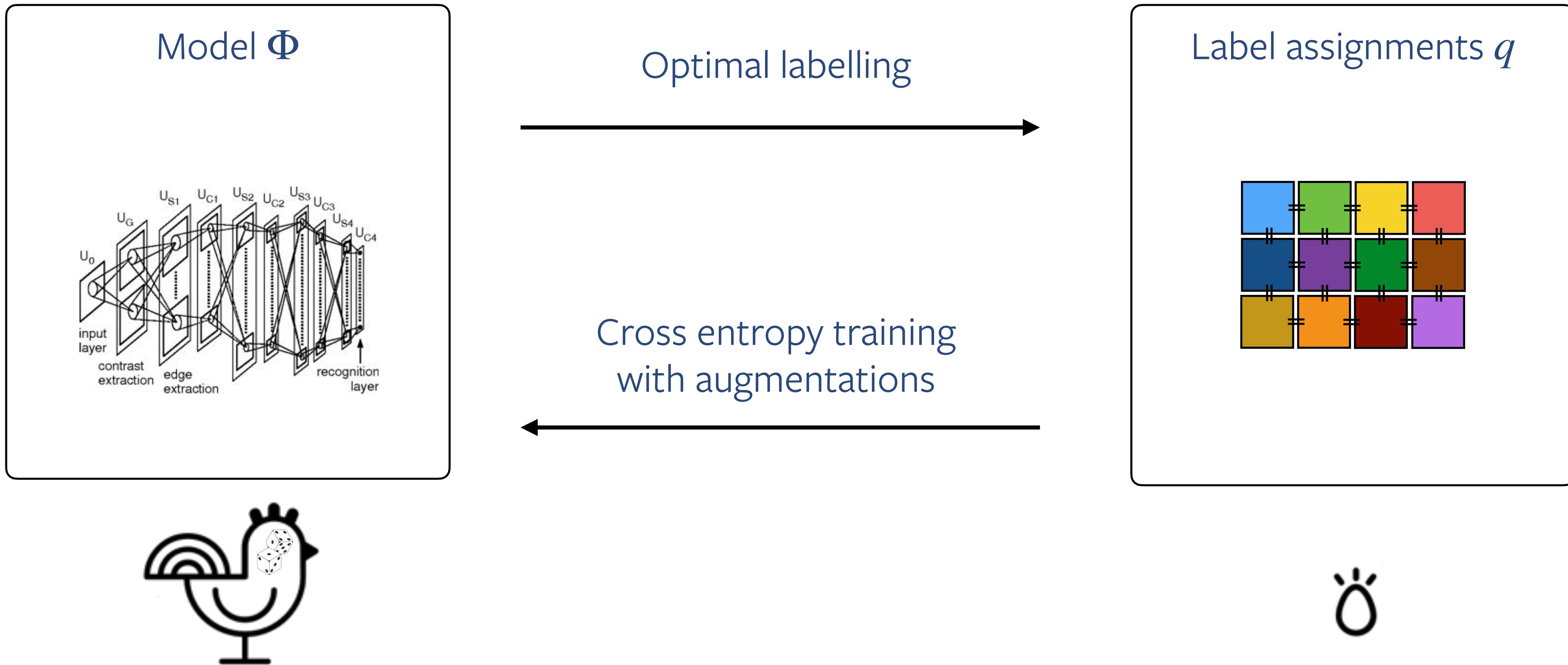
2. Use pseudolabels an equal number of times:



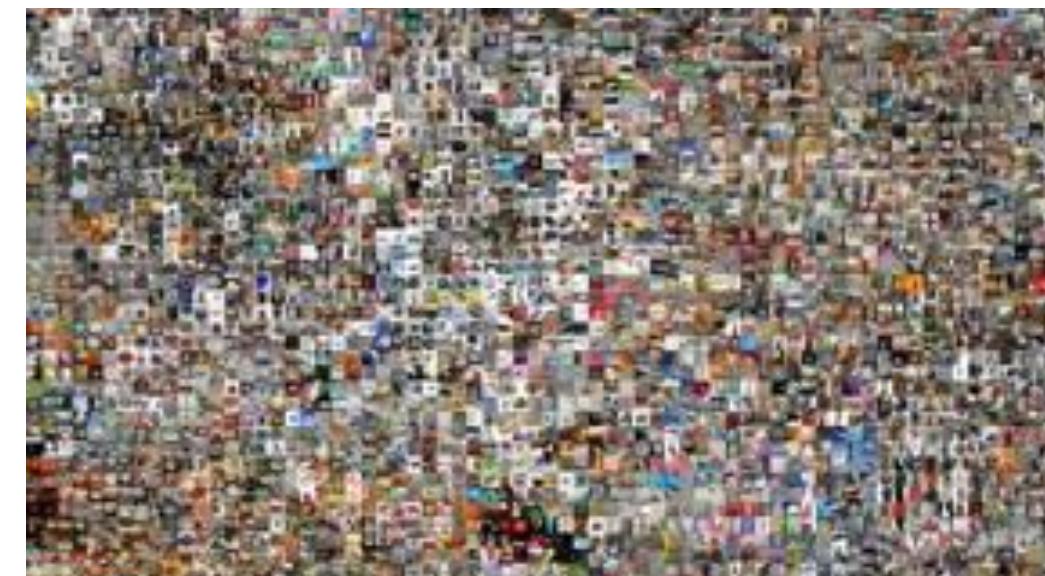
3. Pose as approximate optimal transport:

$$\min_{q, \Phi} L(q, \Phi) \text{ s.t. } \sum_{i=1}^N q(y | \mathbf{x}_i) = \frac{N}{K},$$

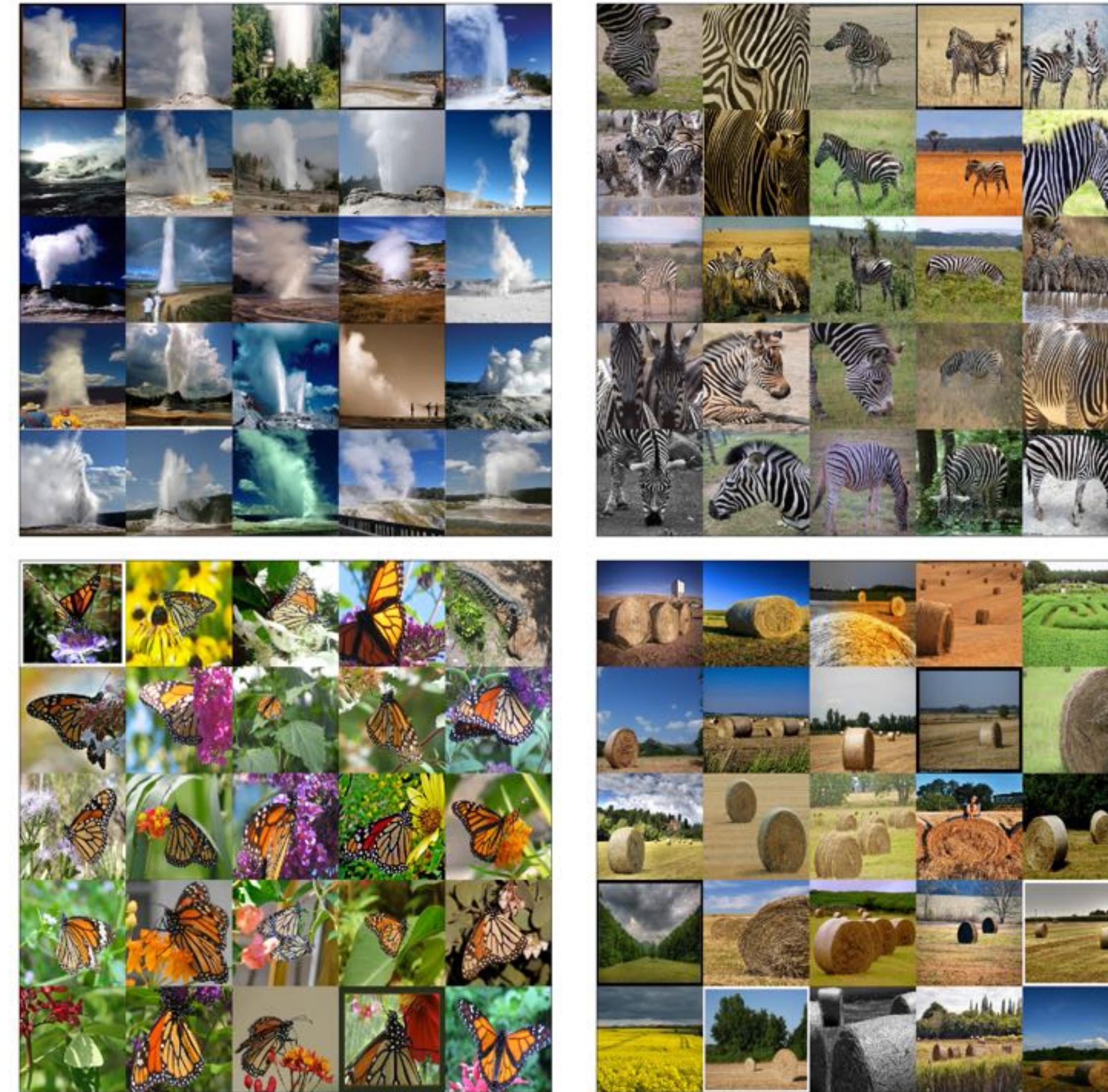
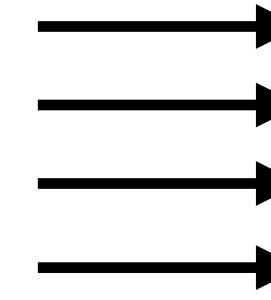
Algorithm



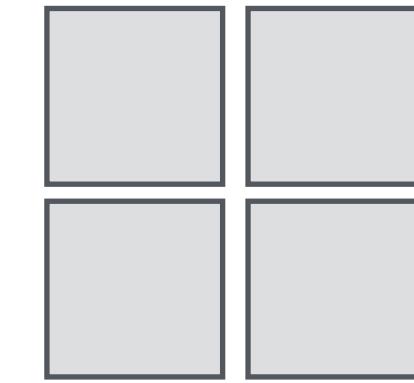
Our method applied on 1.2 million images: Examples



1.2M images



Legend:

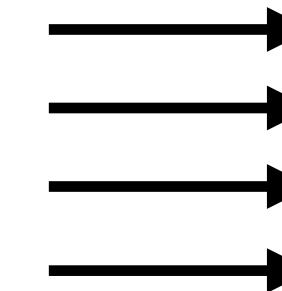


Concept

Automatically discovered concepts match manual annotation.



1.2M images



Legend:



Concept



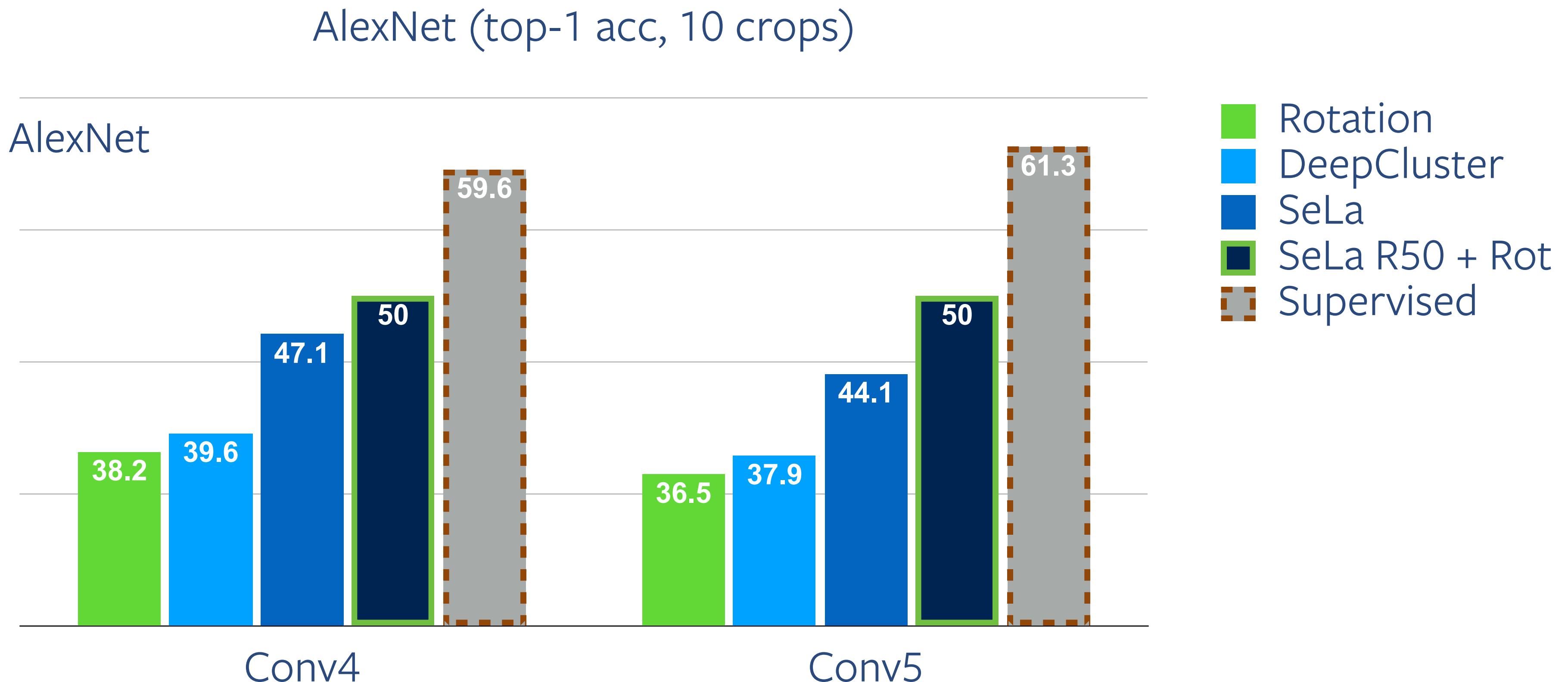
Manually
annotated label

Explore all clusters:



AlexNet, ImageNet linear probes

- Big jump on DeepCluster
- SoTA or close to SoTA for AlexNet



More recently...

Method	Top-1		Δ
	2x224	2x160+4x96	
Supervised	76.5	76.0	-0.5
<i>Contrastive-instance approaches</i>			
SimCLR	68.2	70.6	+2.4
<i>Clustering-based approaches</i>			
SeLa-v2	67.2	71.8	+4.6
DeepCluster-v2	70.2	74.3	+4.1
SwAV	70.1	74.1	+4.0

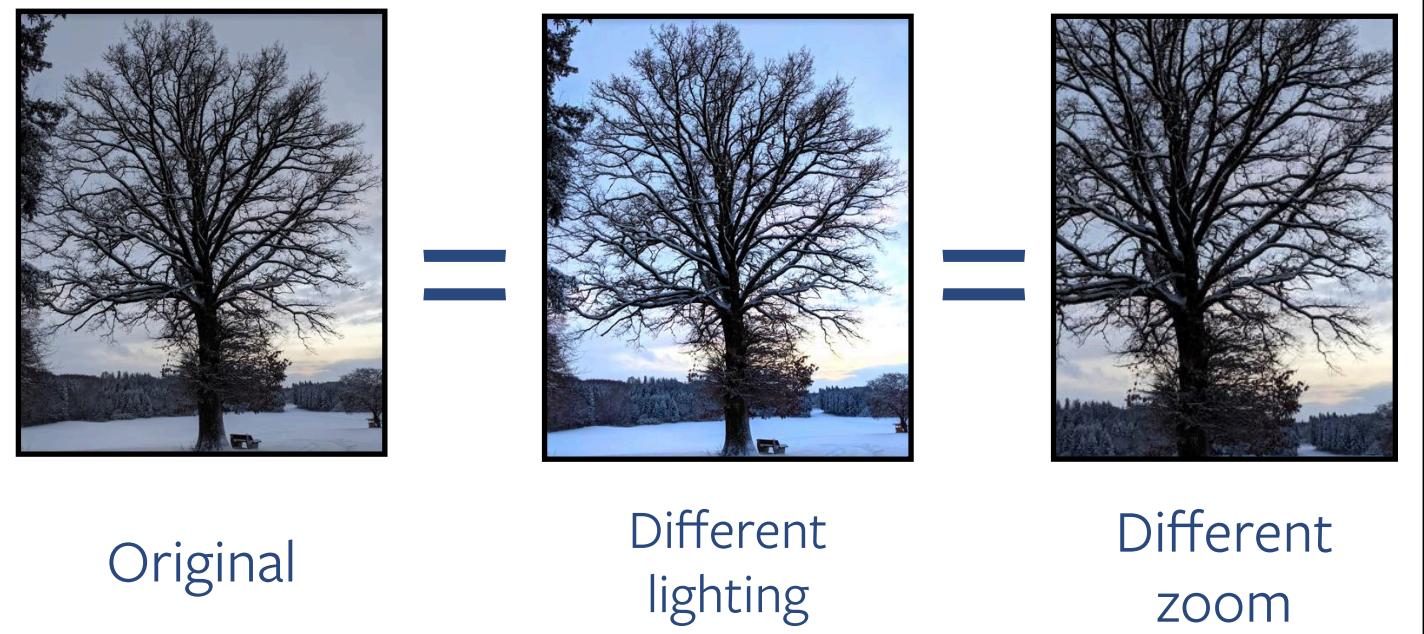
- SwAV uses SeLa's SK algo
- SeLa-v2 better than SimCLR

	Method	Momentum	Operation	Top-1
1	DINO	✓	Centering	76.1
2	-	✓	Softmax (batch)	75.8
3	-	✓	Sinkhorn-Knopp	76.0
4	-		Centering	0.1
5	-		Softmax (batch)	72.2
6	SwAV		Sinkhorn-Knopp	71.8

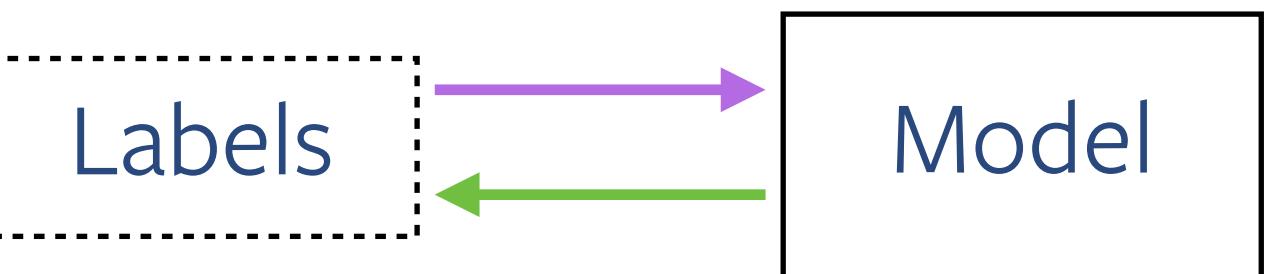
- DINO with SeLa's SK: same performance.

Self-supervised labelling from three core ideas

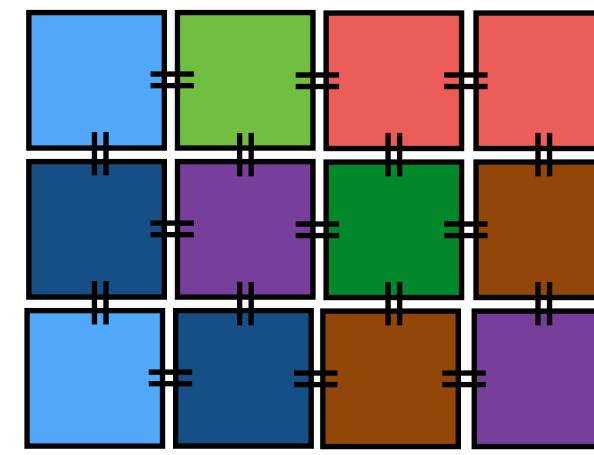
Invariance to augmentations



Virtuous cycle of labelling and representation learning



Balanced labelling



(1) Transformations

- Data augmentations “infuse knowledge”

(2) Useful labels

- Labels discovered are similar to ground-truth
- Can be used to analyze how the network “sees” the data

(3) Balanced pseudo-labelling

- Well defined, fast objective
- No trivial solutions

Multi-modal SSL: leveraging video & audio for labelling

The key to image understanding is separating meaning from appearance.



Original



Different lighting



Mirrored



Different zoom

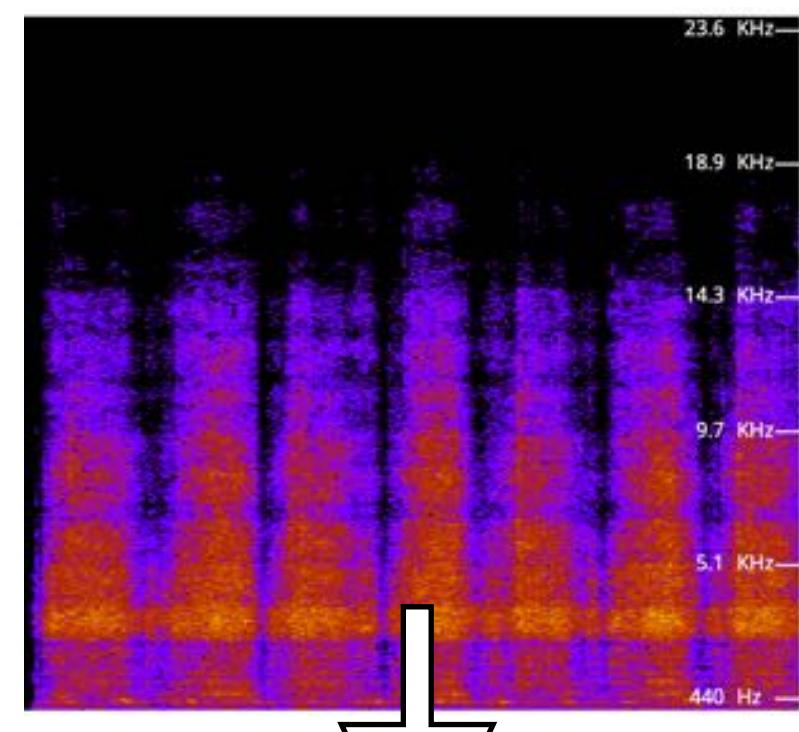
Augmentations

Multiple modalities can yield useful semantic information.

Video



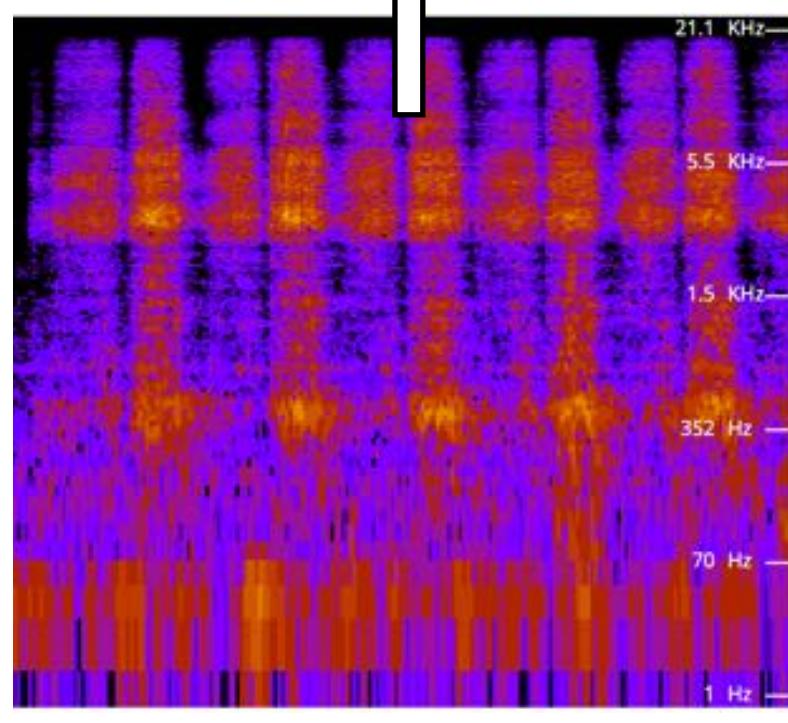
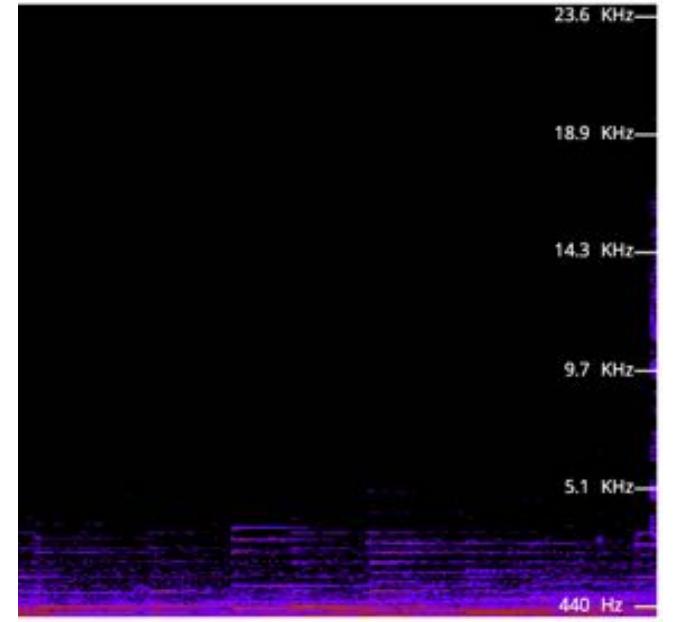
Audio



Video



Audio

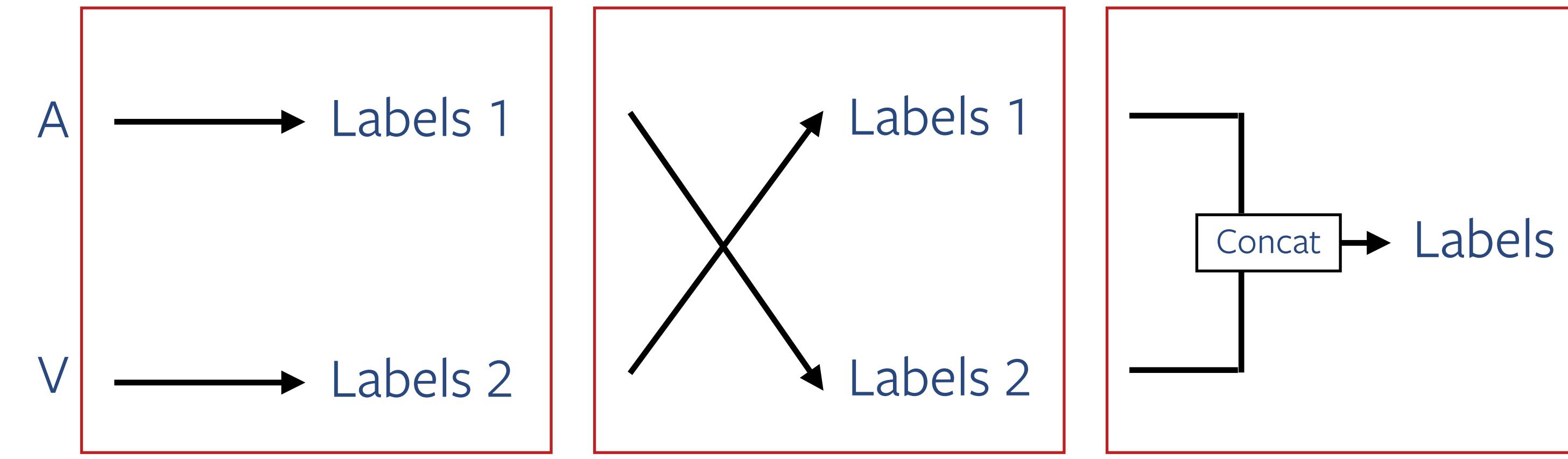
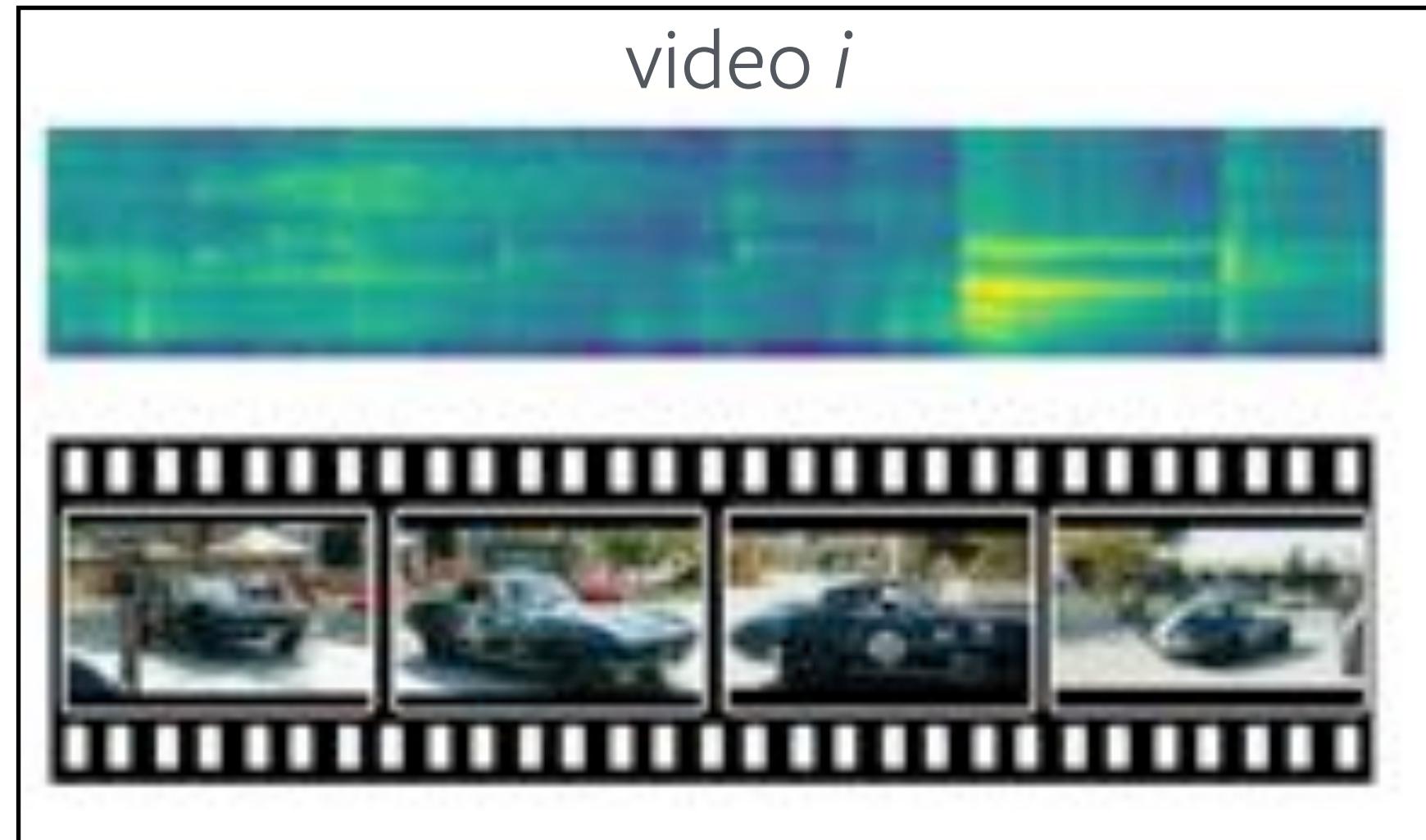


Why clustering for videos?

- 🤔 Clustering works well for images
- 💰 Videos are expensive to annotate.
- 📈 Video content is rapidly increasing.



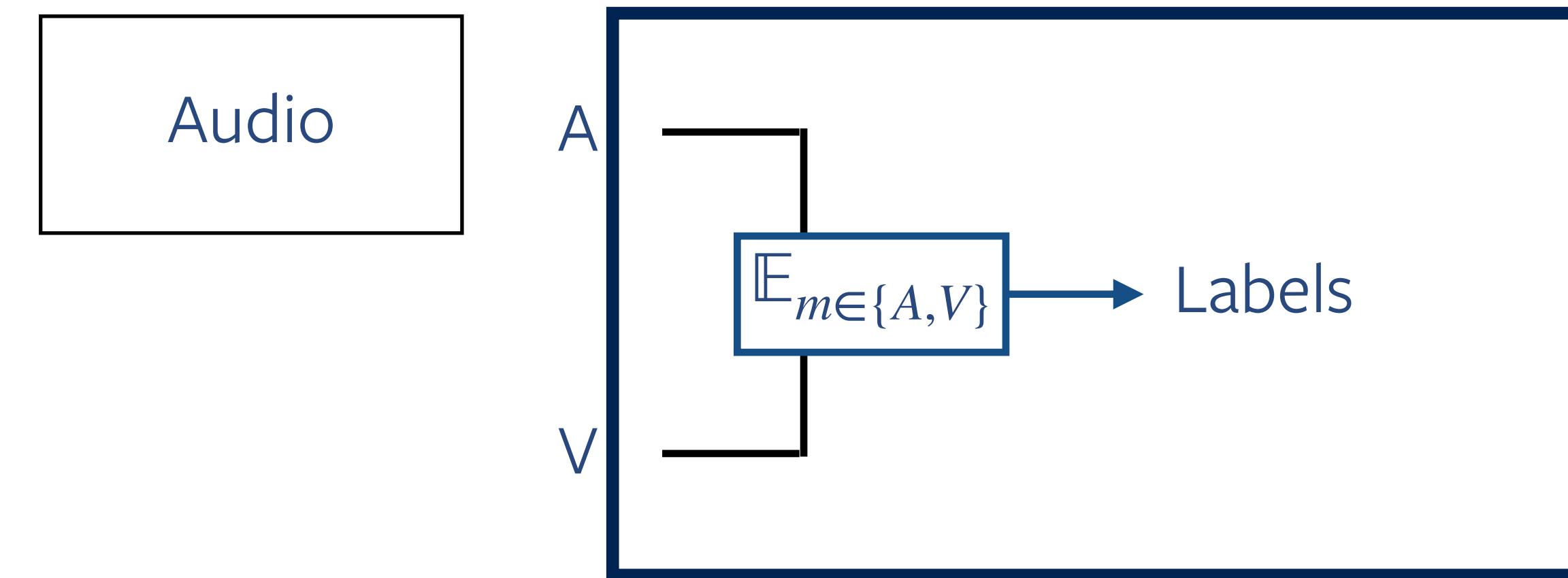
Clustering multi-modal data



- ✗ does not use same-source information
- ✗ two different sets of clusters

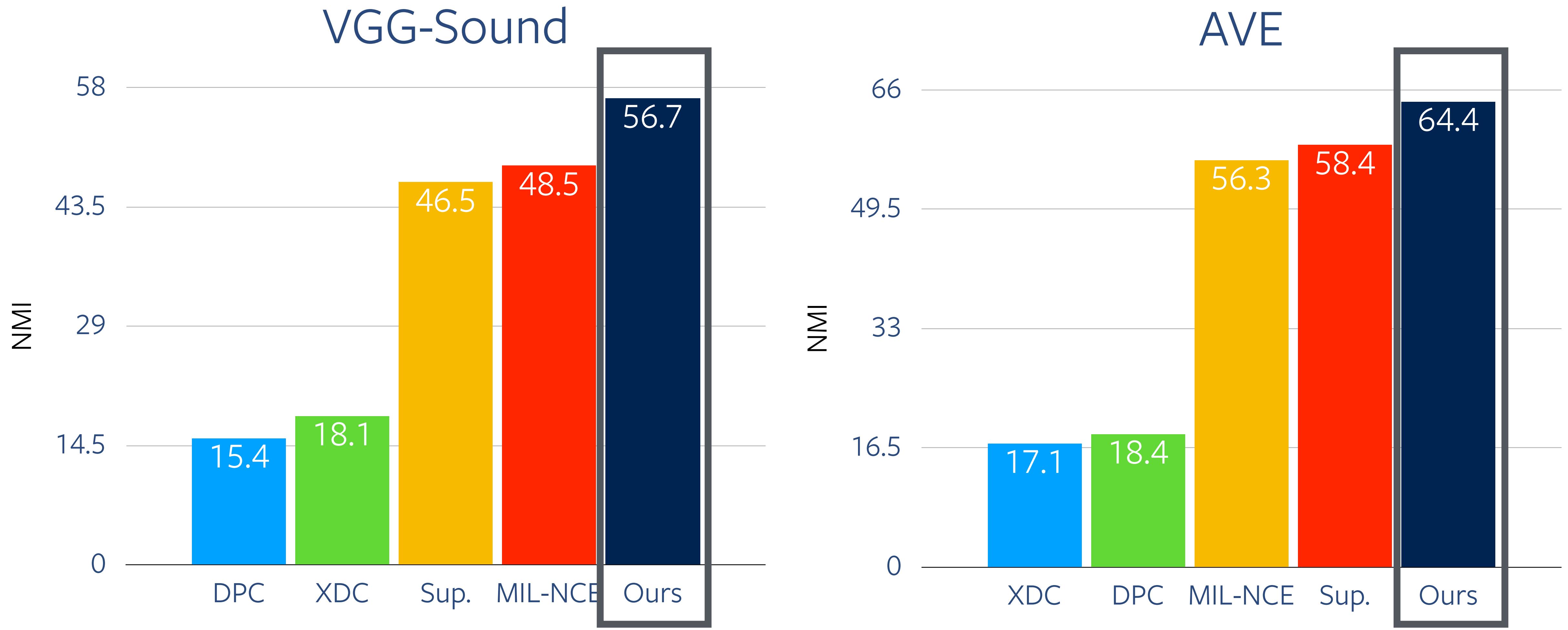
- ✗ two different sets of clusters
- ✗ hard to interpret

- ✗ concatenation can just rely on stronger modality and ignore the other



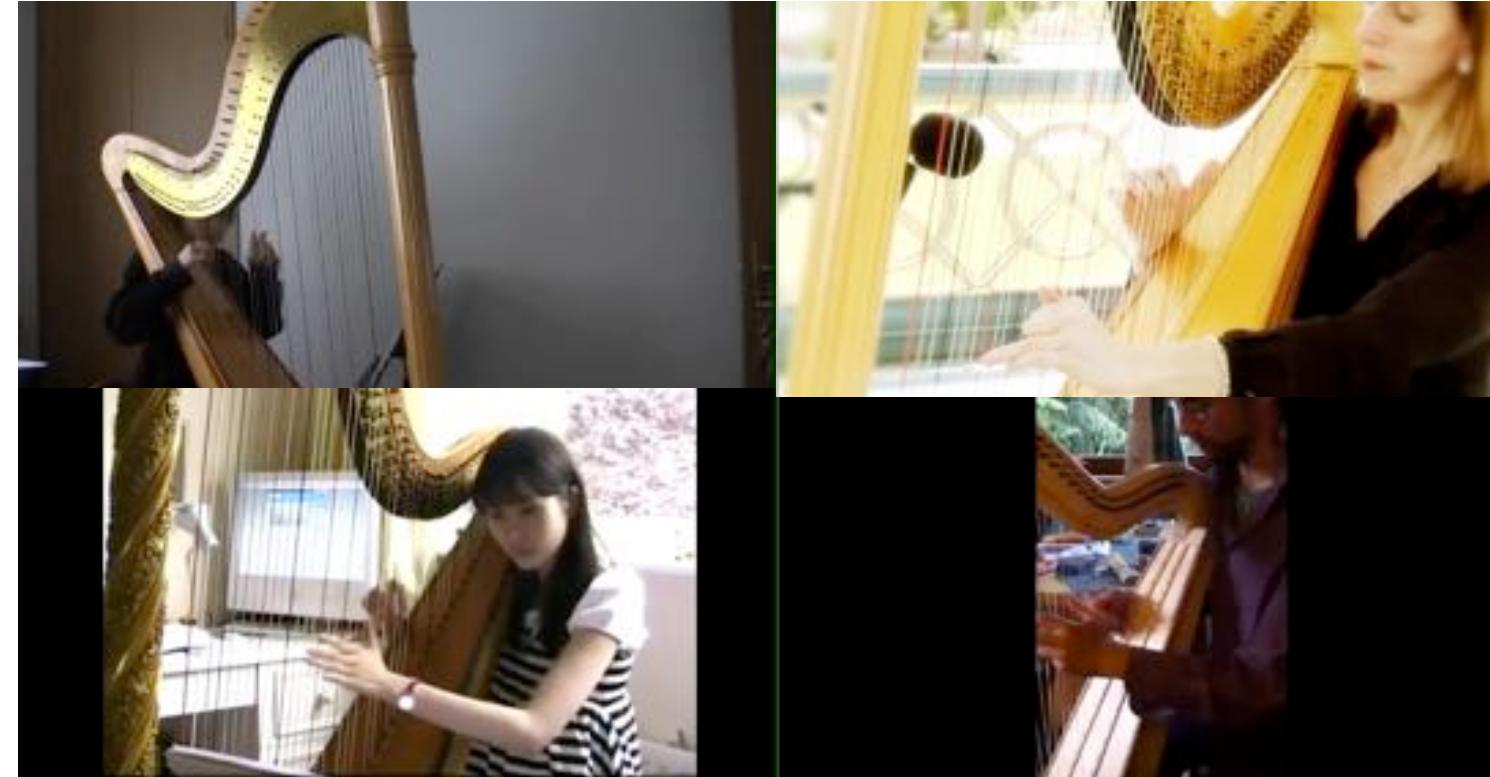
Our idea:
view each modality as an augmentation

Good feature representations $\not\rightarrow$ good clustering



[Han et al., arXiv; Alwassel et al., NeurIPS 2020; Miech et al., CVPR 2020]

Discovering concepts without manual annotations from 230K videos



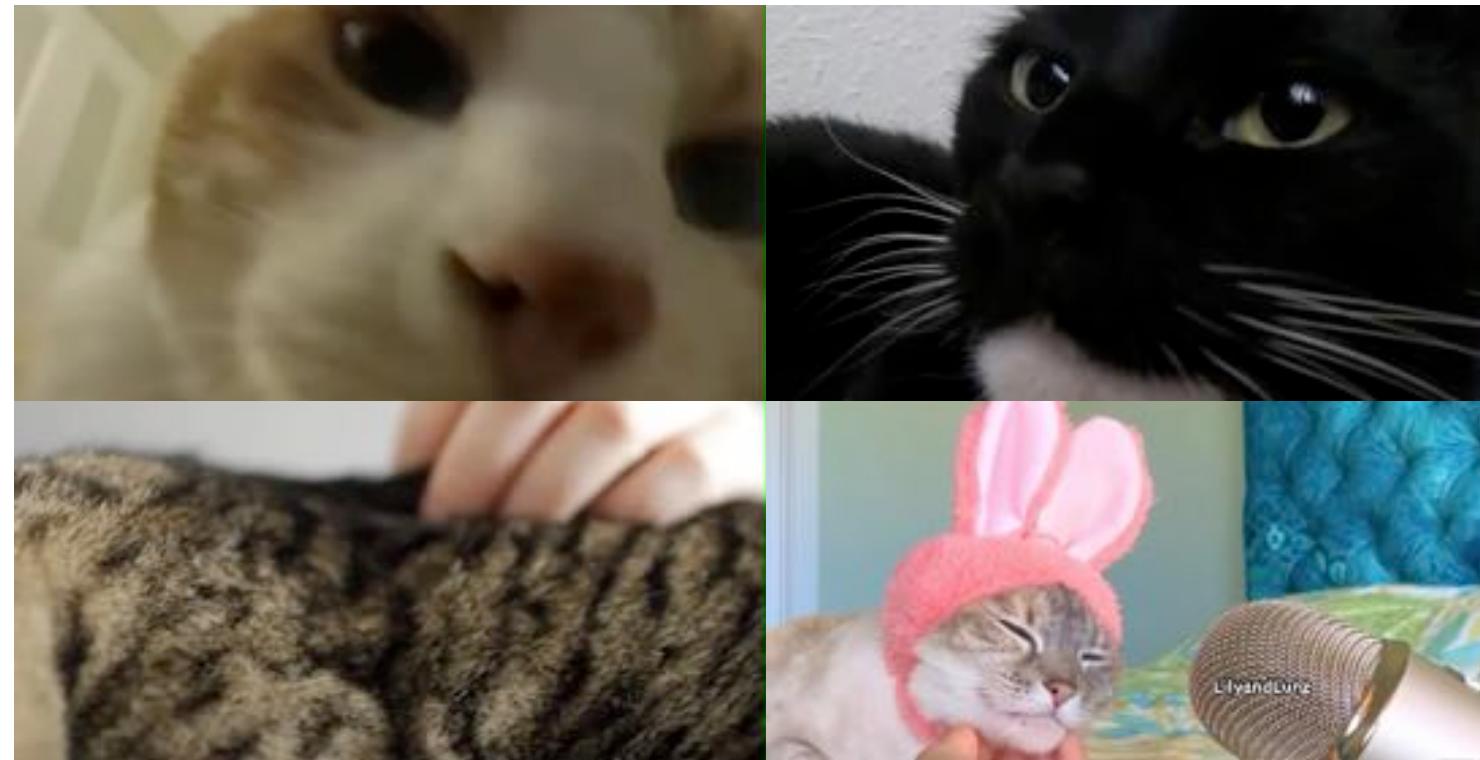
"Playing harp"



"Fireworks"



"Hockey game"



"Cat growling"



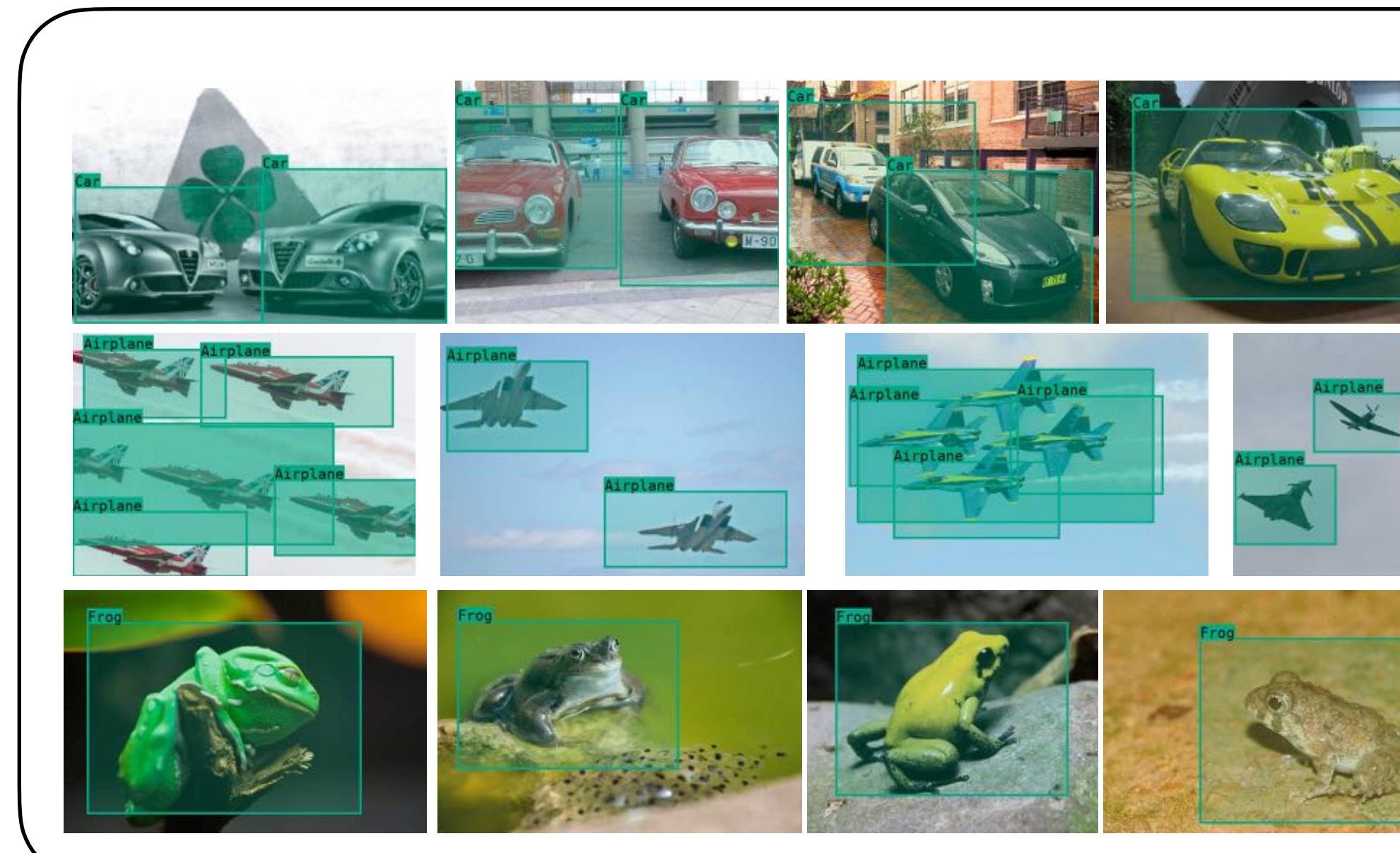
"Electric guitar"



"Vehicle driving"

“Self-Labelling” videos: three main findings

- Clustering framework of SeLa is well-suited
- Generalizable to model any distribution (Zipf/exponential etc.)
- Good feature representations do not imply good clustering



Follow-up work:
using this + contrastive
learning, we can detect
objects self-supervisedly.

On Compositions of Transformations in Contrastive Self-Supervised Learning

Mandela Patrick*, Yuki M. Asano*, Polina Kuznetsova, Ruth Fong, João F. Henriques, Geoffrey Zweig, Andrea Vedaldi
ICCV'21

Invariance vs distinctiveness

In contrastive learning, we define positives and negatives.

Should the representations enforce invariance or distinctiveness?

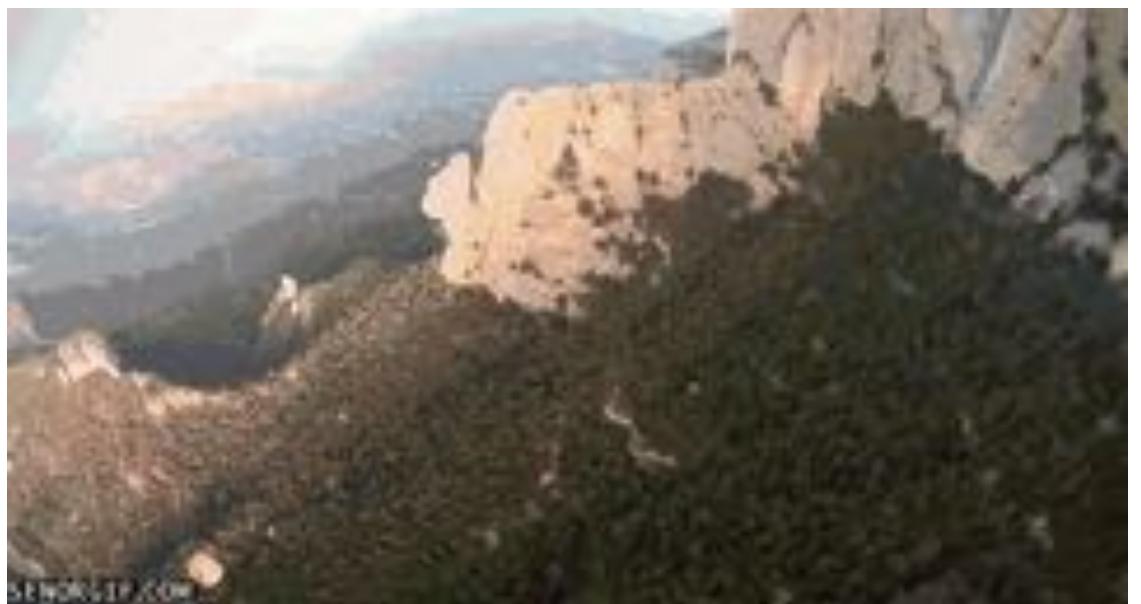


?
=



Learning hypotheses we test:

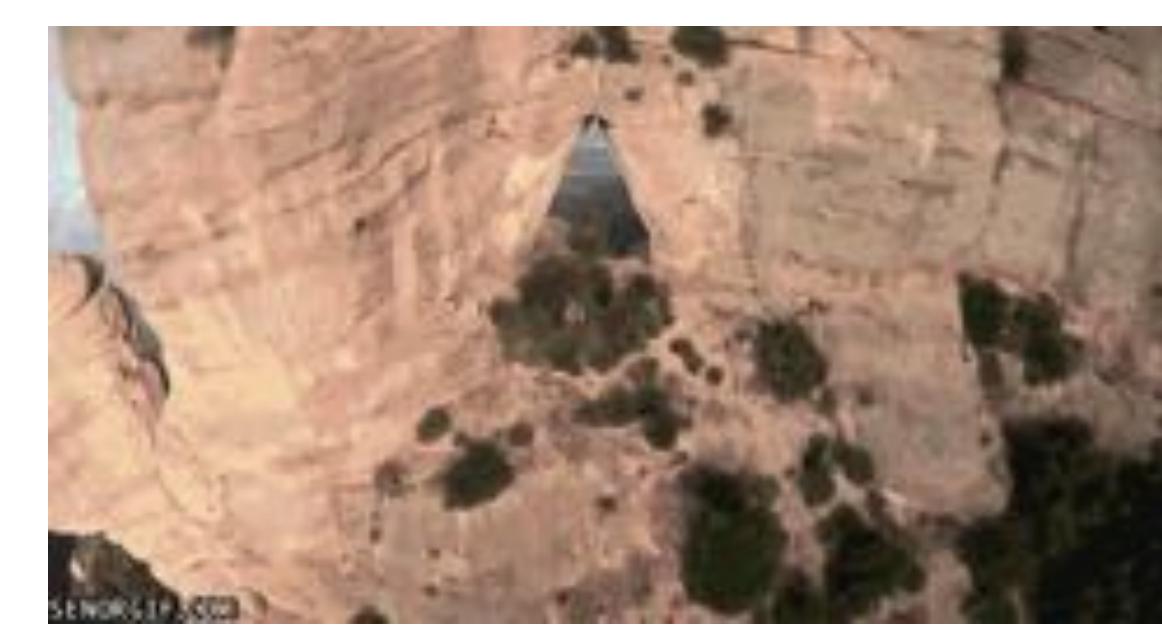
1. Sample Distinctiveness



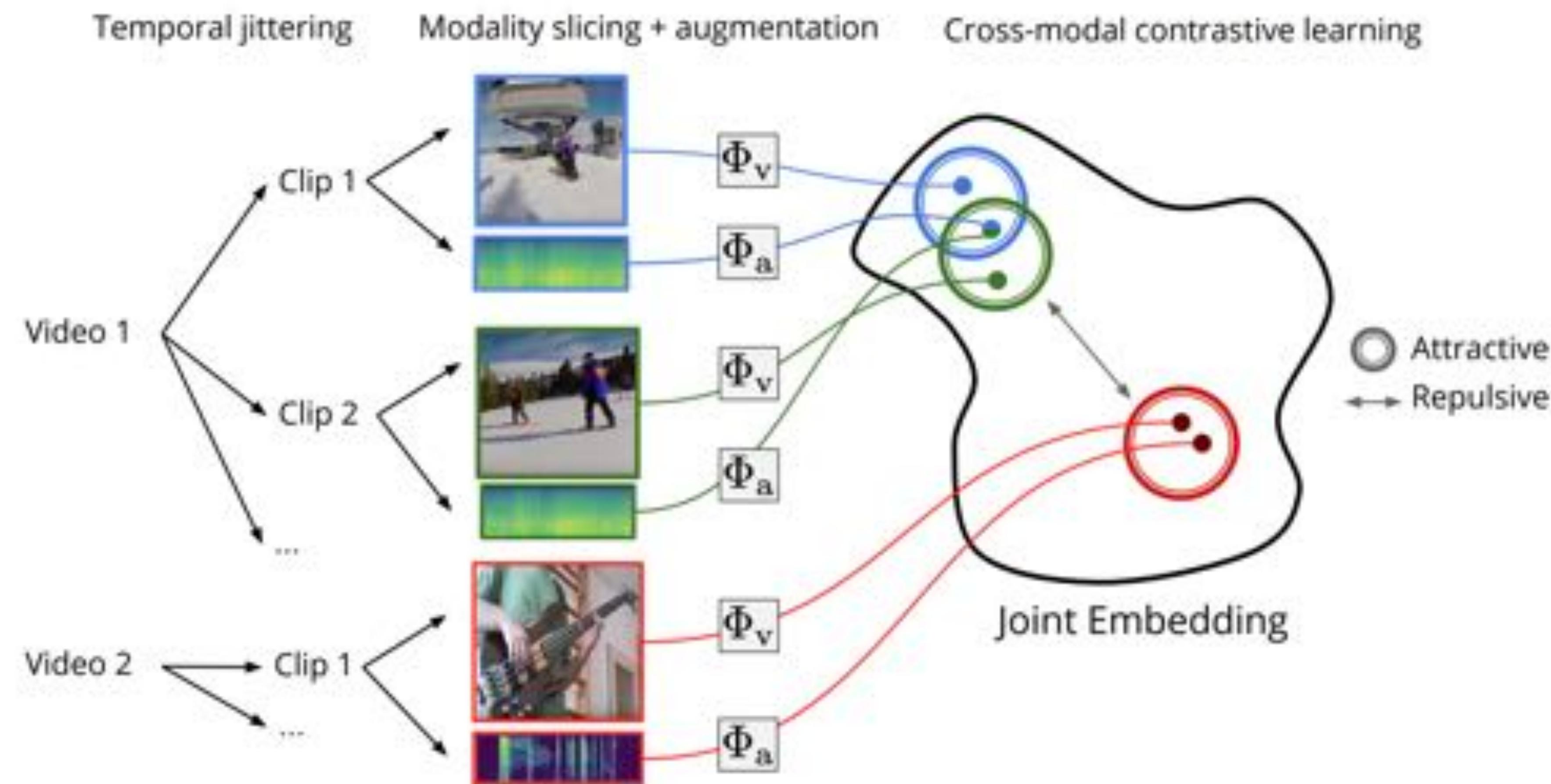
2. Time Reversal



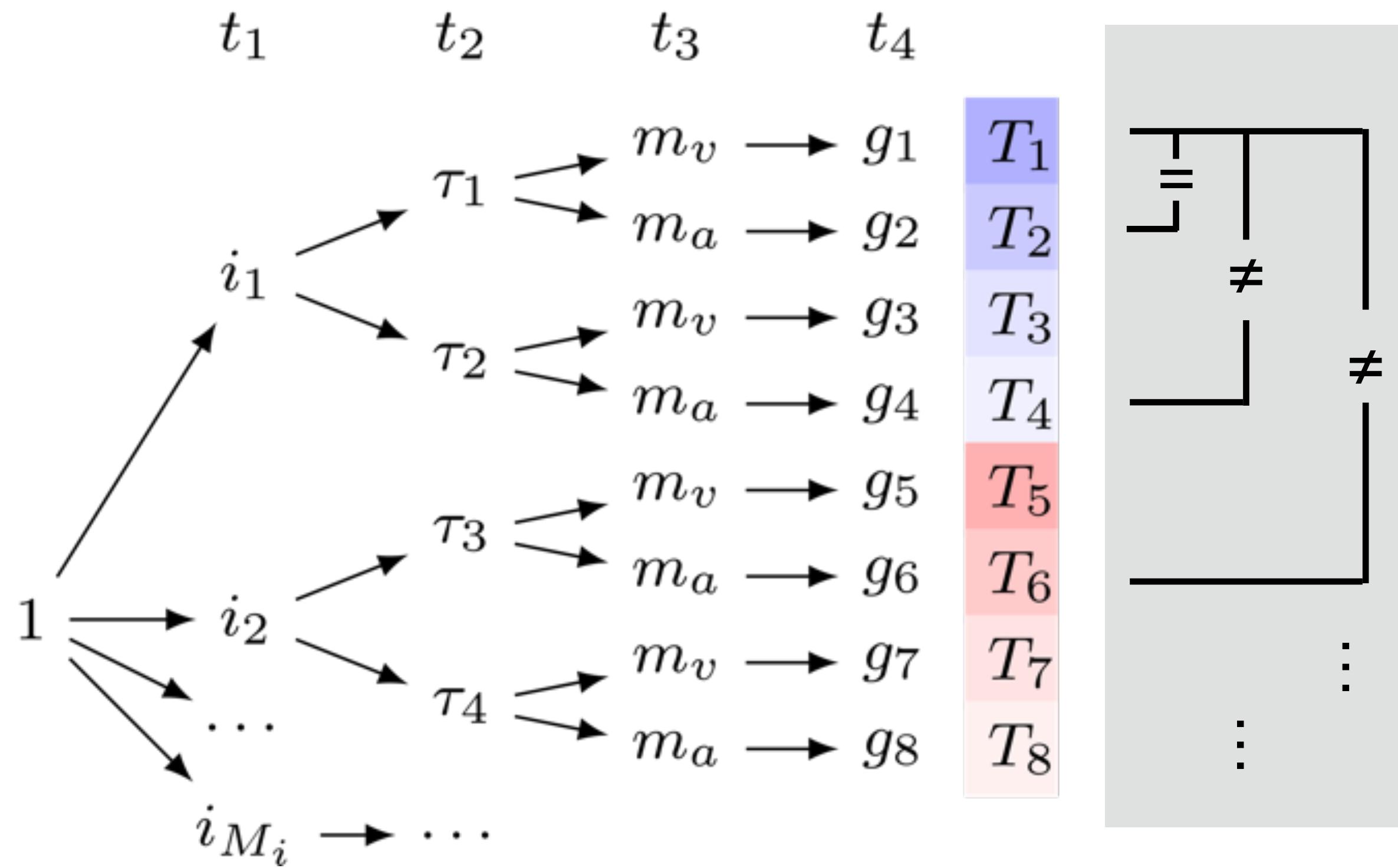
3. Time Shift



Framework



Example: distinctive to sample & time shift, invariant to modality



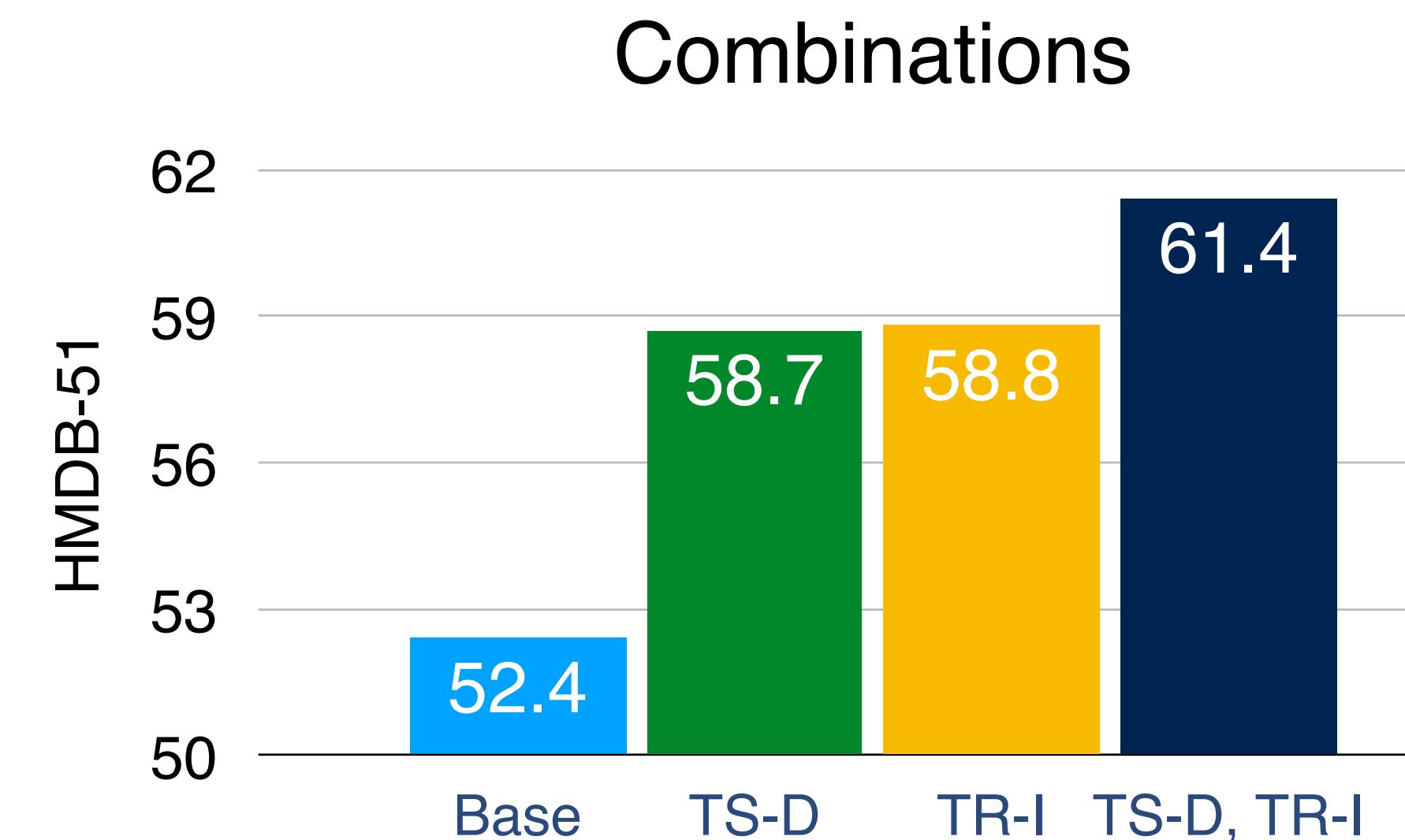
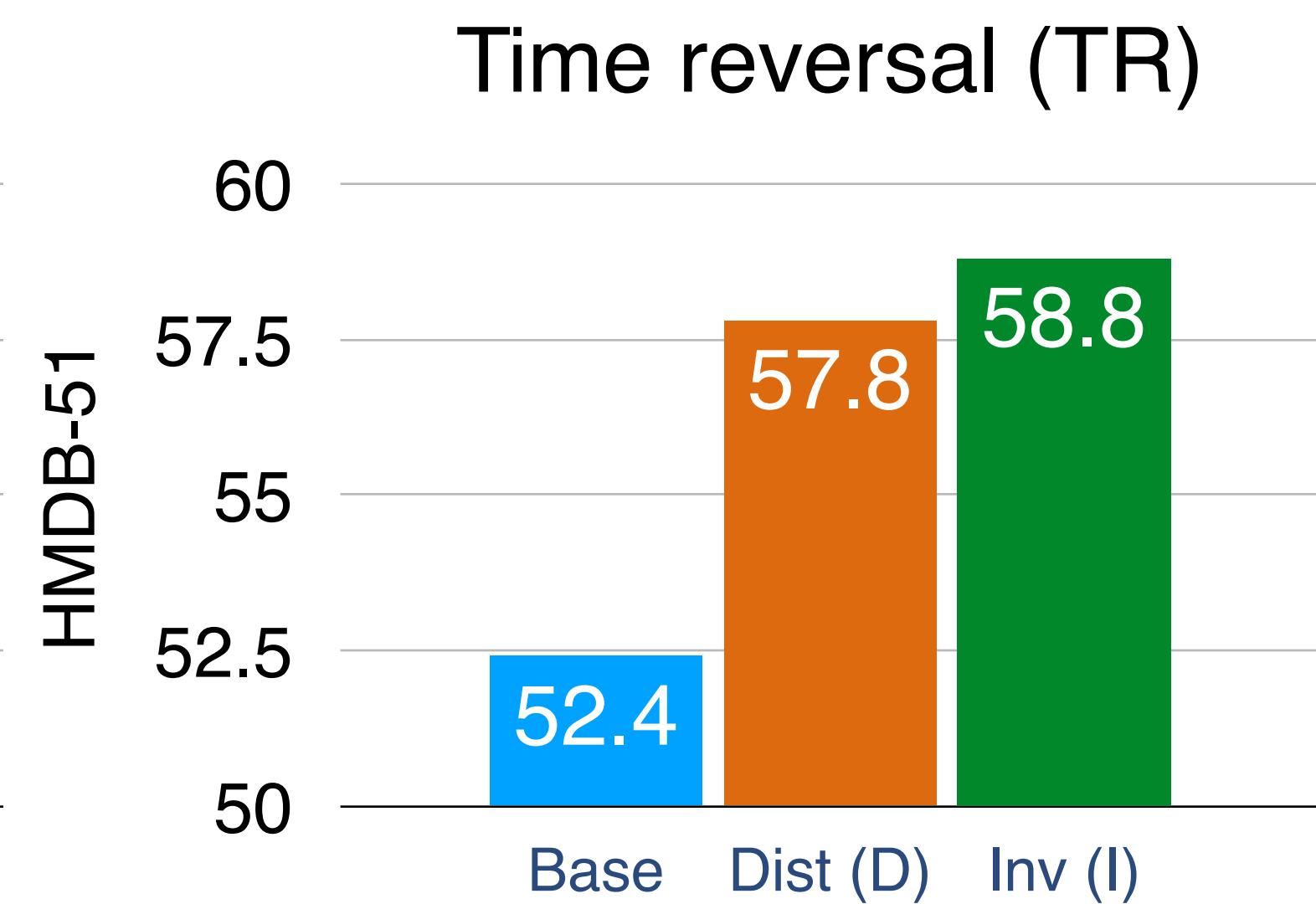
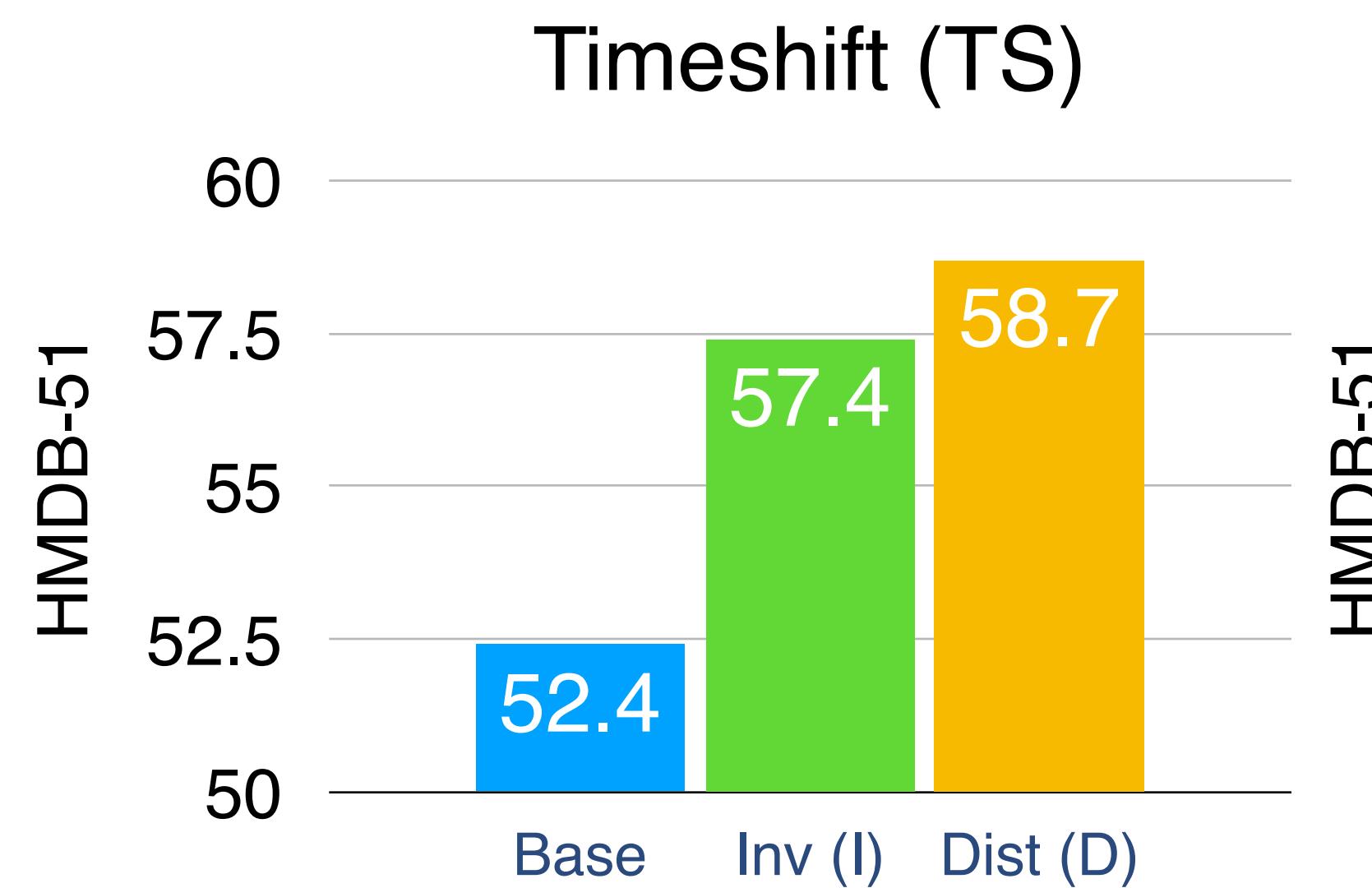
Same embedding for video and audio at *same time* from *same video*

Different embedding for *different time-indices* from the *same video*

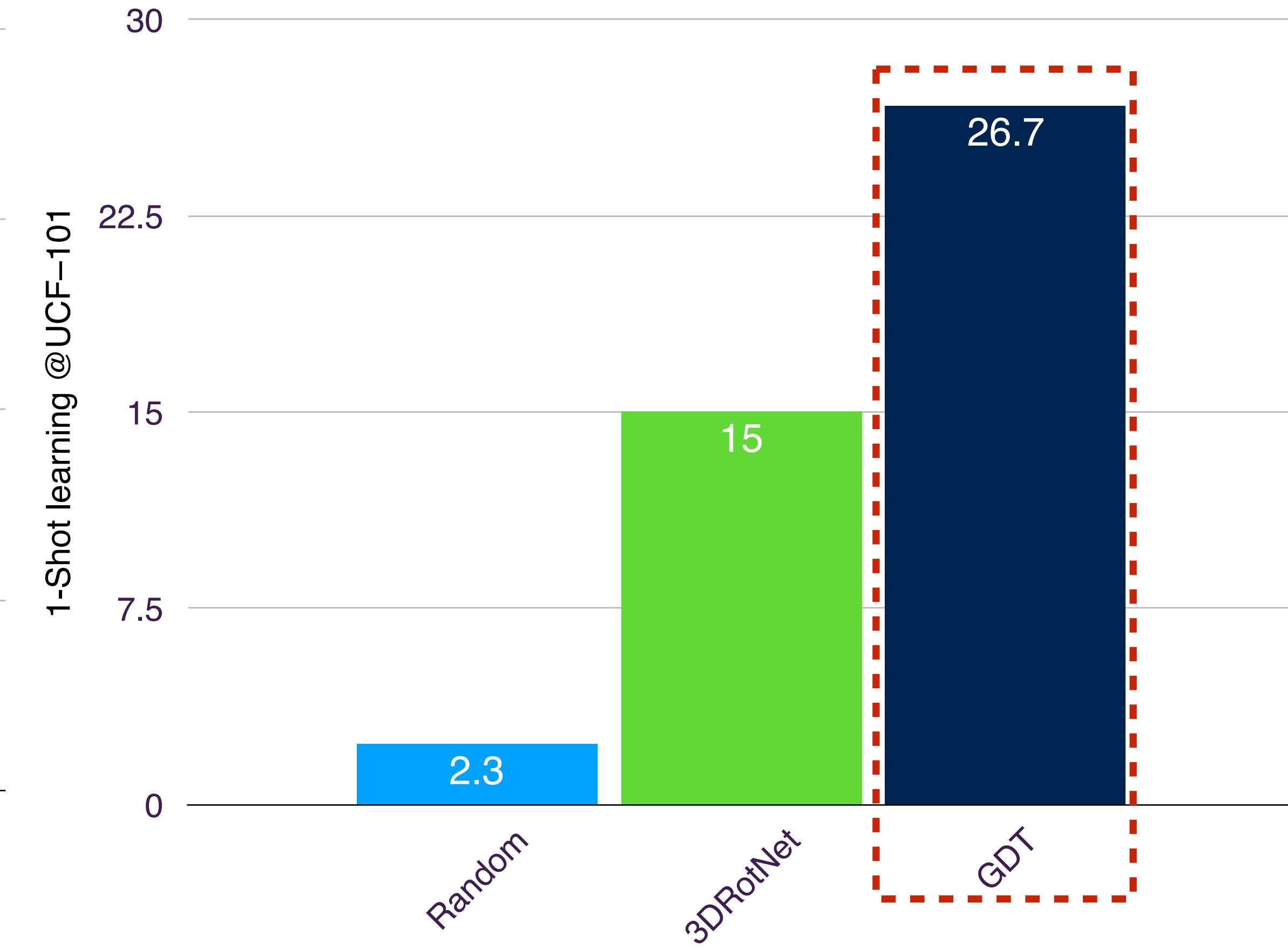
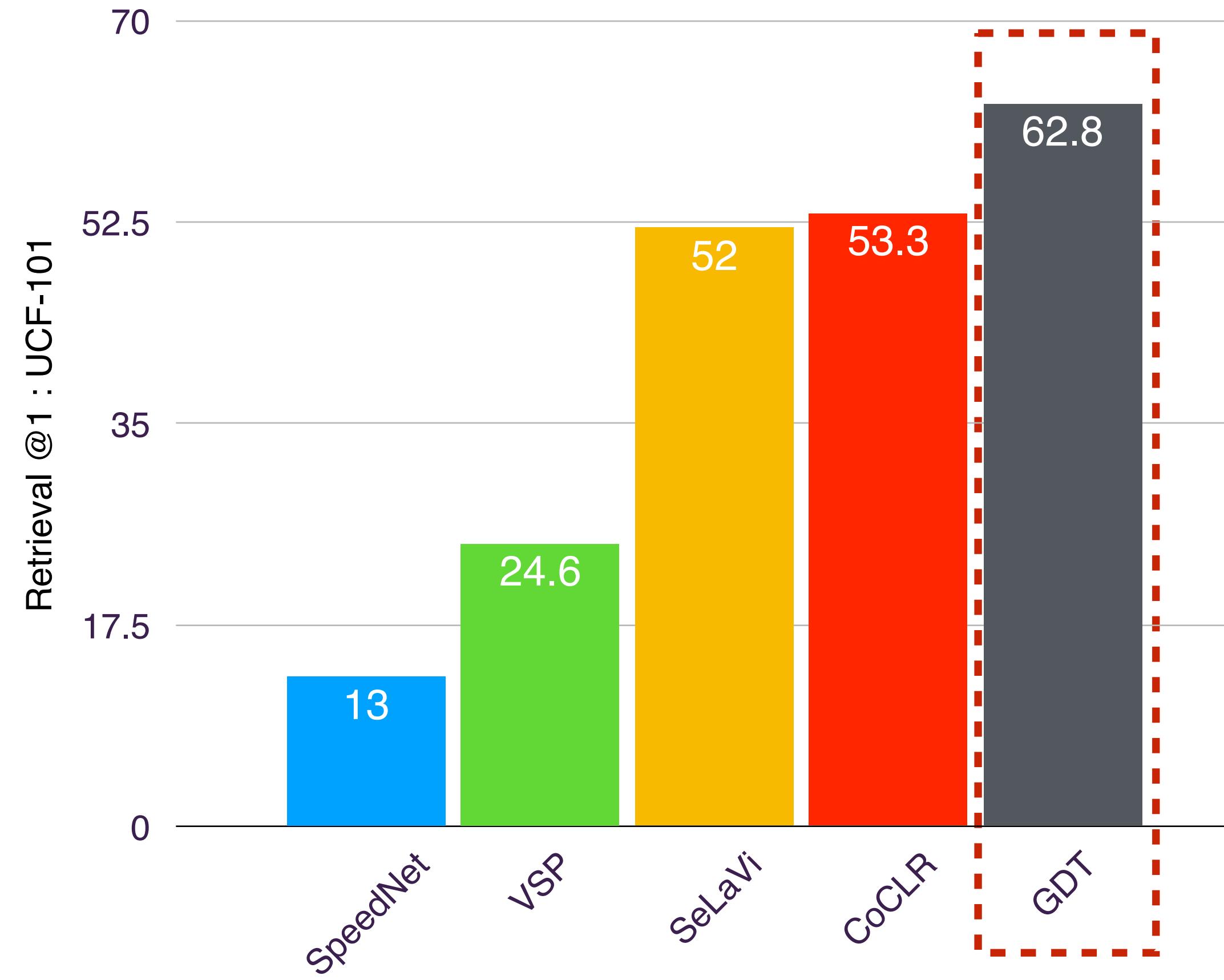
Different embedding for *different video*

Cross-modally

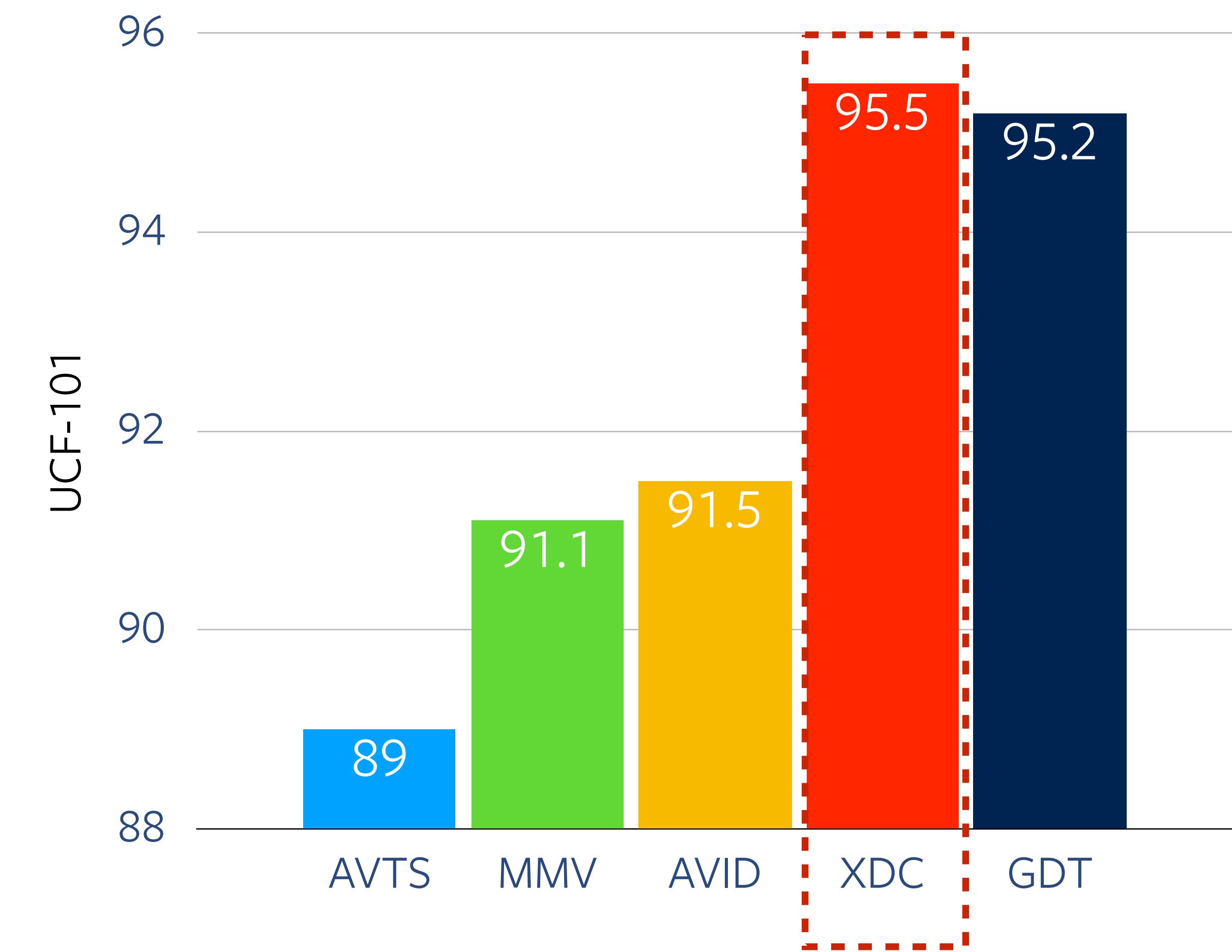
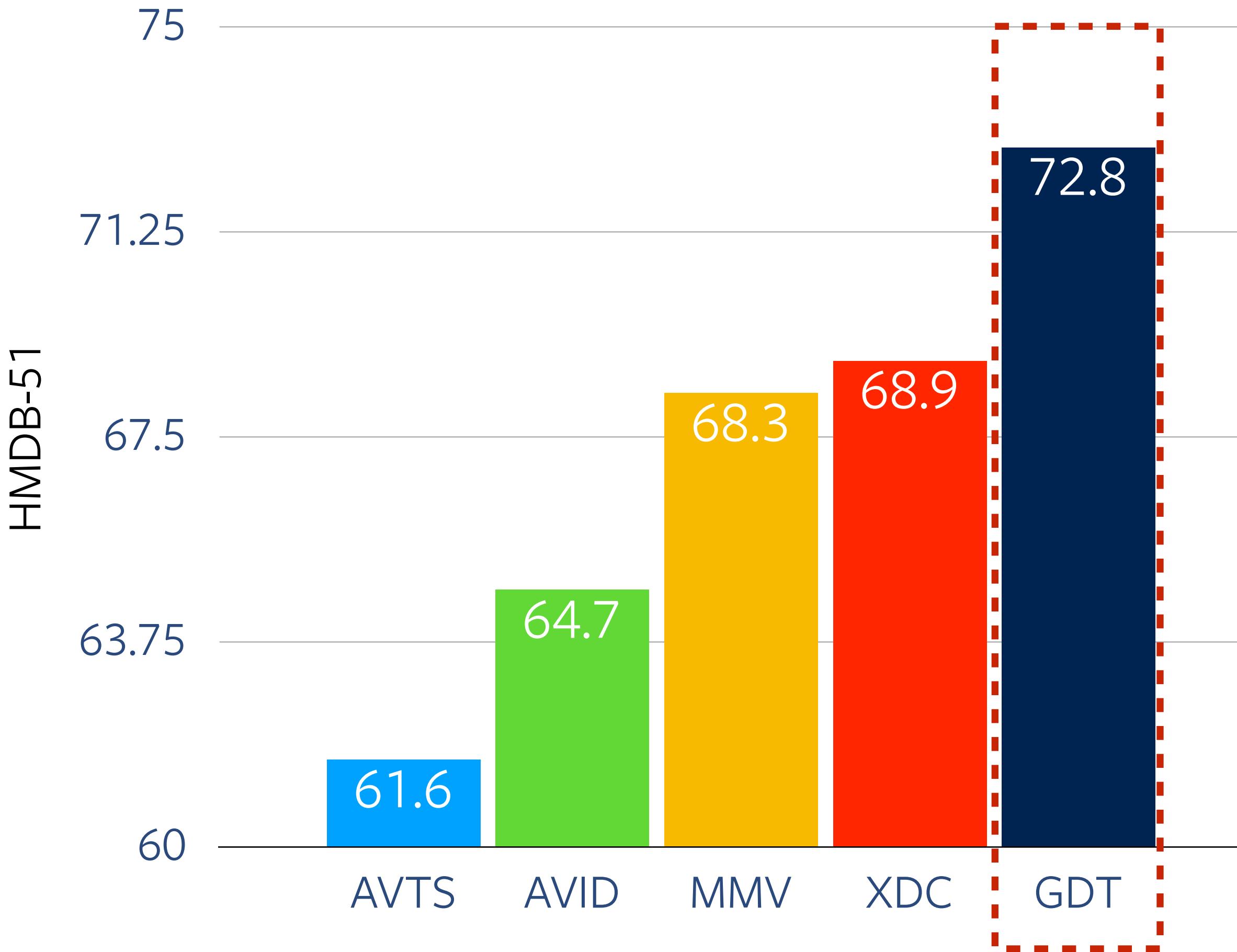
Gains from hypotheses



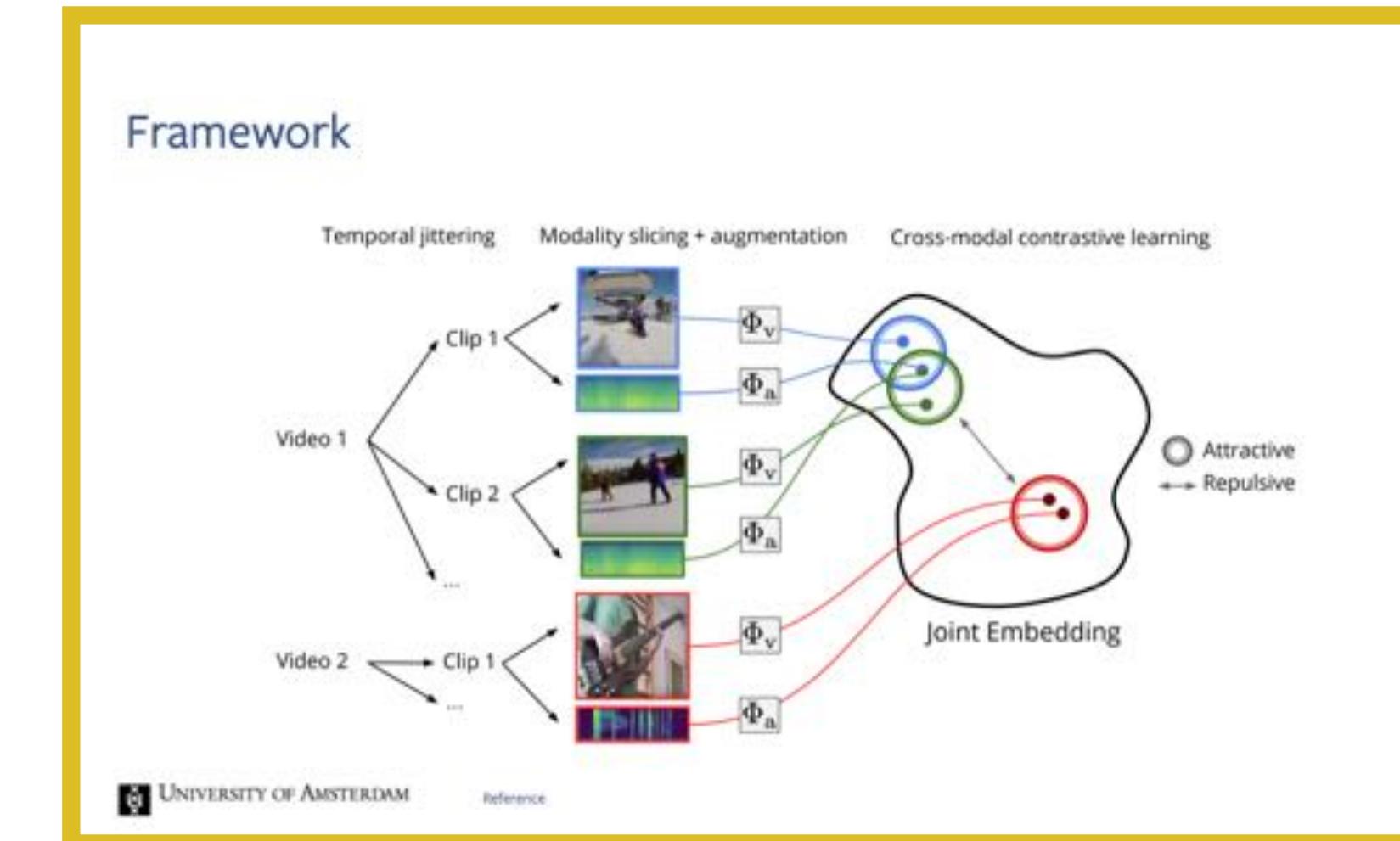
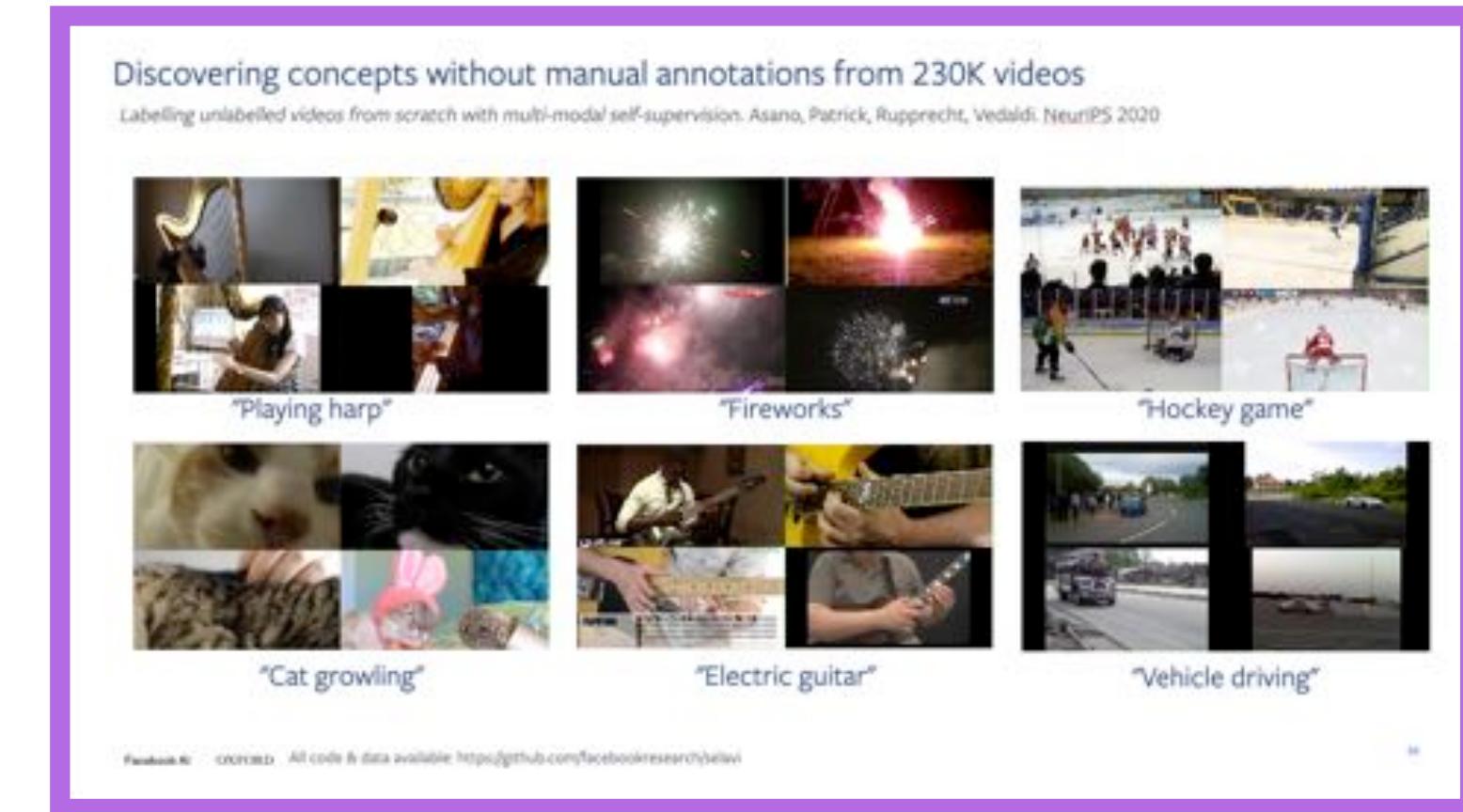
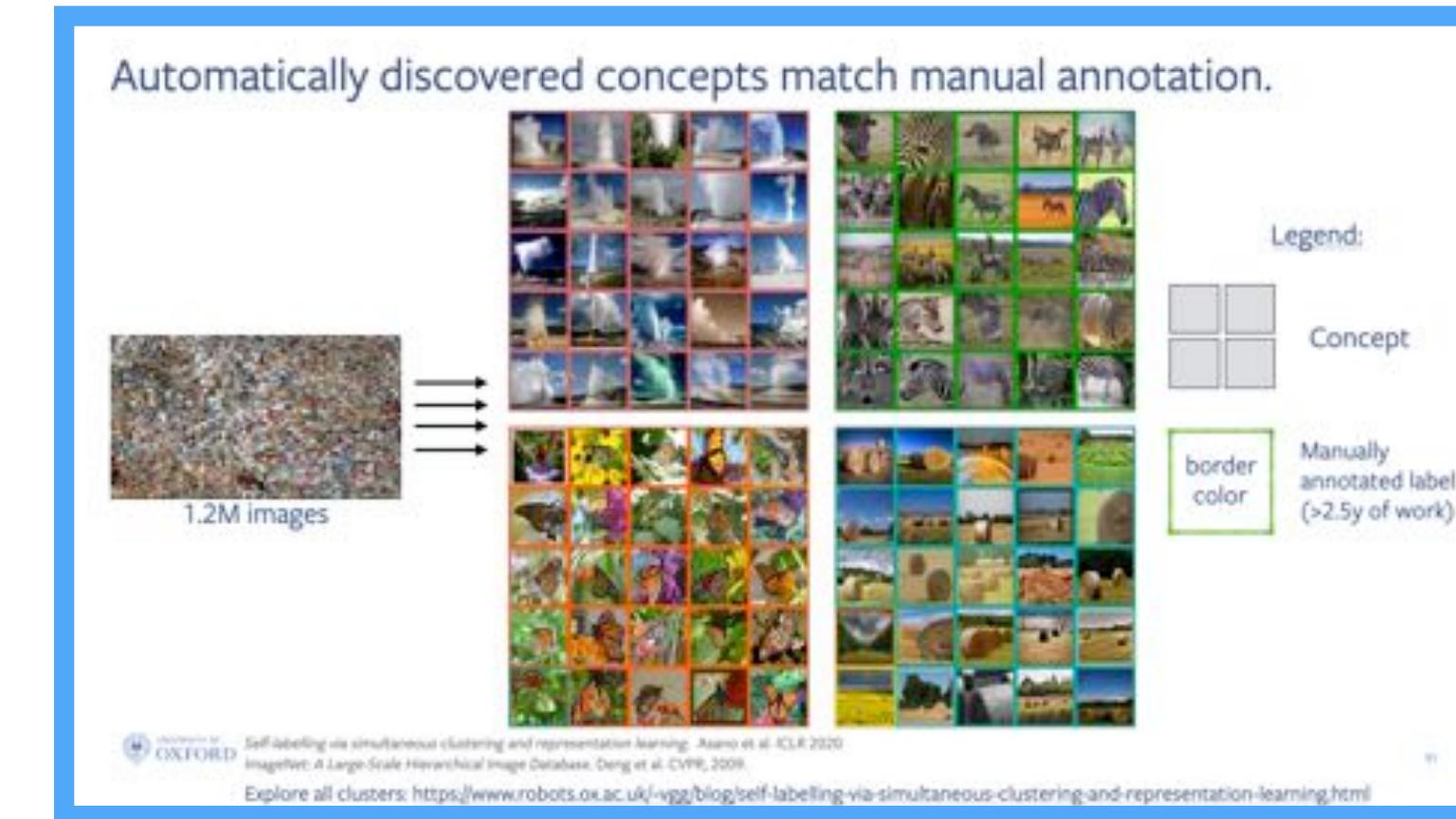
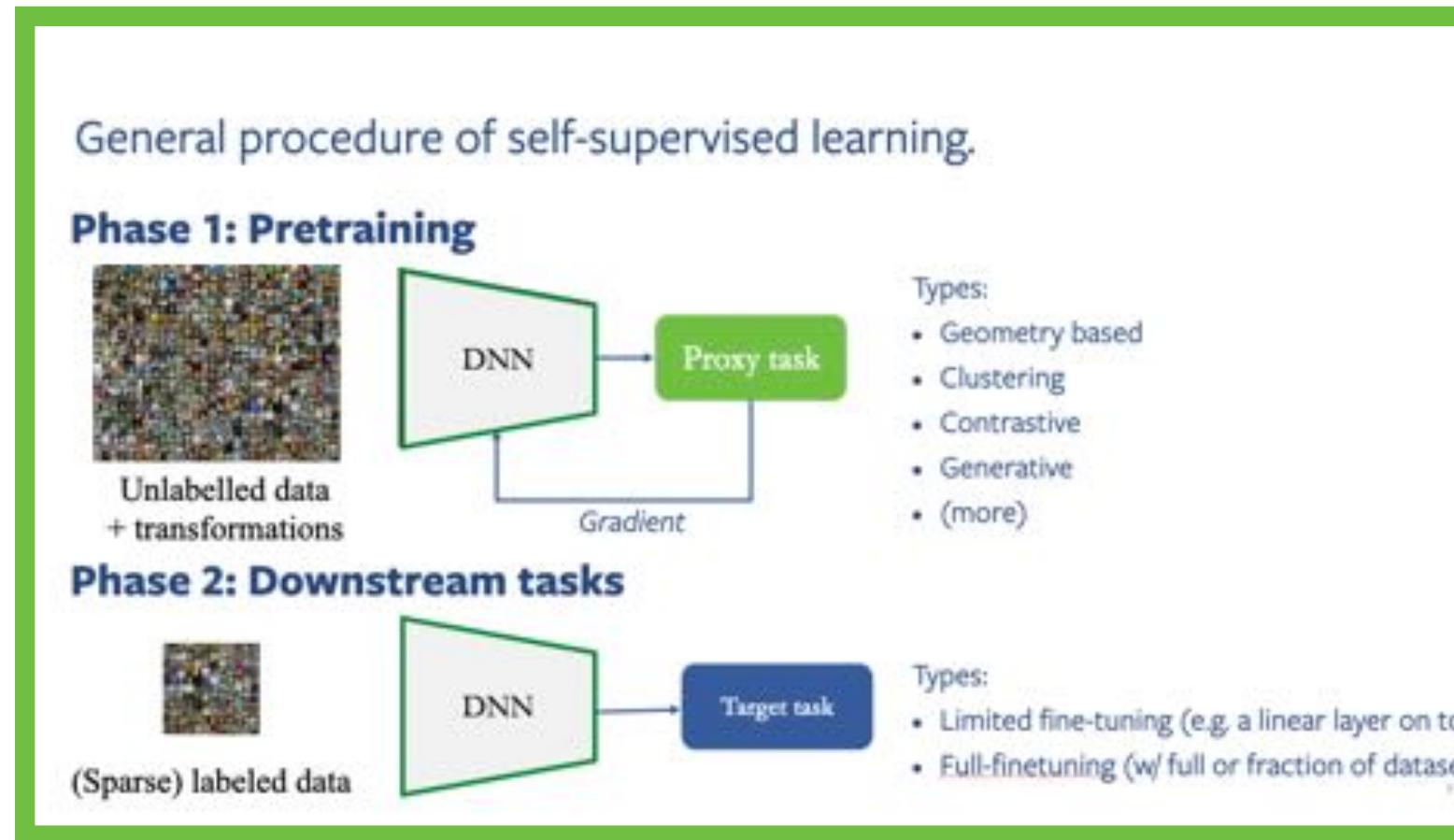
SOTA video action retrieval and few-shot learning results



SOTA finetuning video-action recognition results



Visual recap of this talk



Tutorial:

<https://urlis.net/5mj36>