

# Self-supervised Learning from Images and Videos

YUKI M ASANO

AT

SBA WORKSHOP, UNIVERSITY OF BIRMINGHAM



UNIVERSITY  
OF AMSTERDAM

# Hi, I'm Yuki

- Assistant Professor with Video & Image Sense (VIS) Lab
  - Self-supervised Learning
  - Multi-modal Learning
  - Large Model Adaptation
- More info: <https://yukimasano.github.io/>

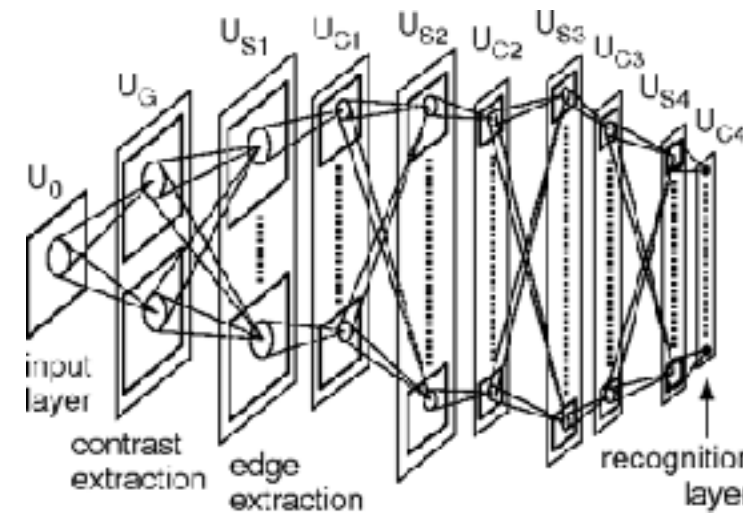
focus of today



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

## Algorithms

*Deep neural networks*



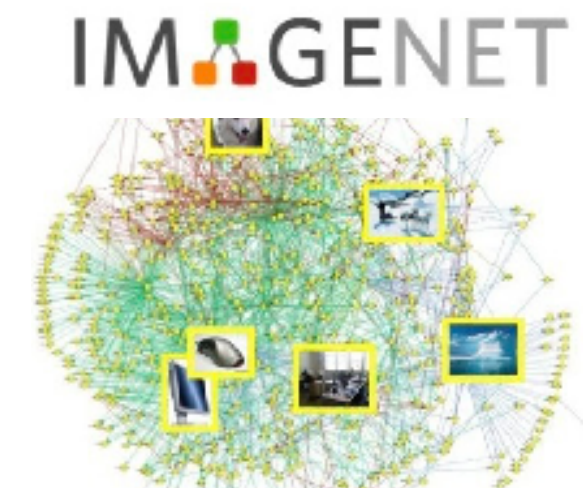
## Hardware

*GPUs*



## Data

*Large scale datasets*

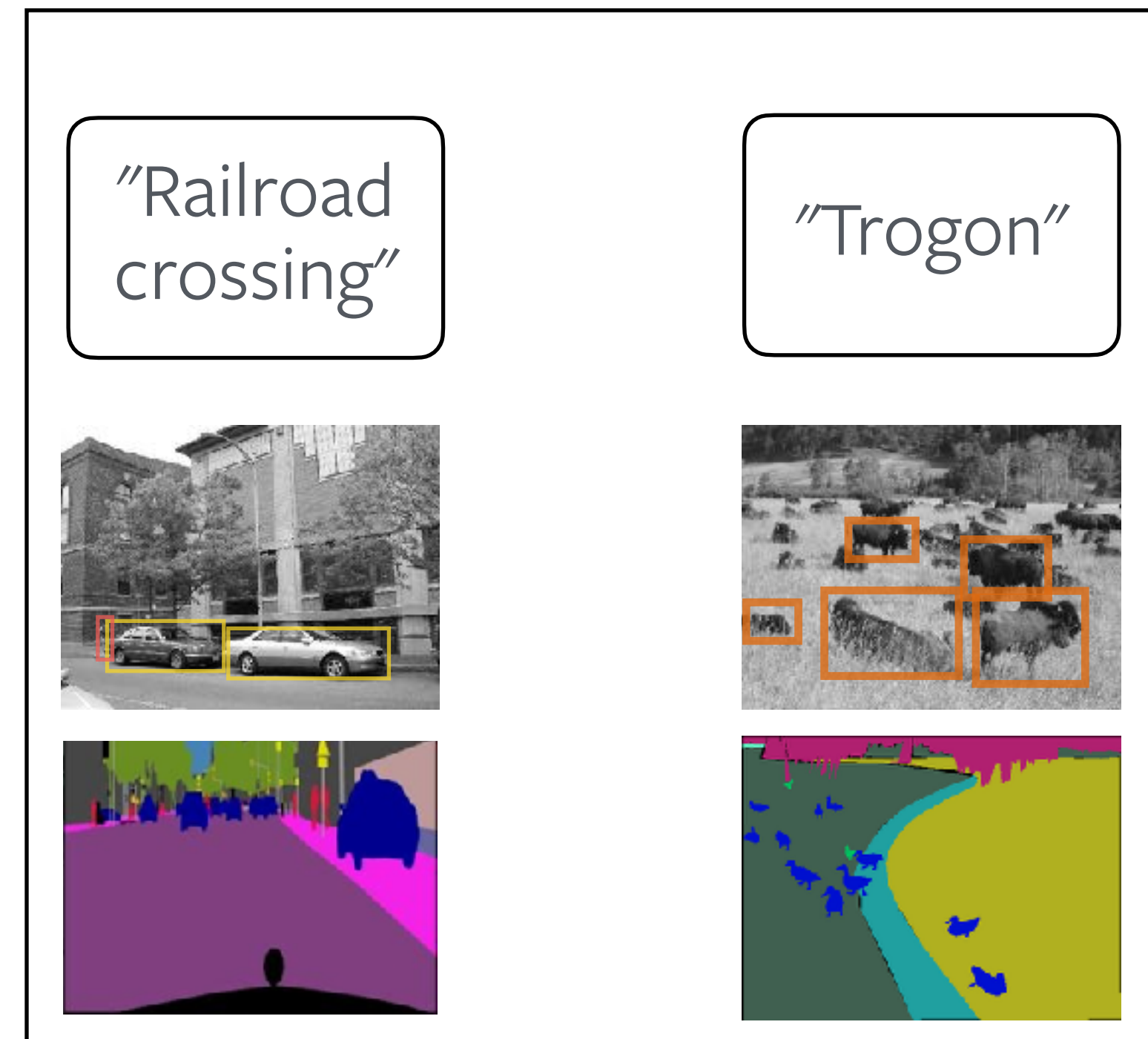


# Manual annotations for the data are limiting.

## Supervised Learning



Images are often cheap



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

# Self-supervised Learning replaces the need for labels & annotations.

## Self-supervised Learning

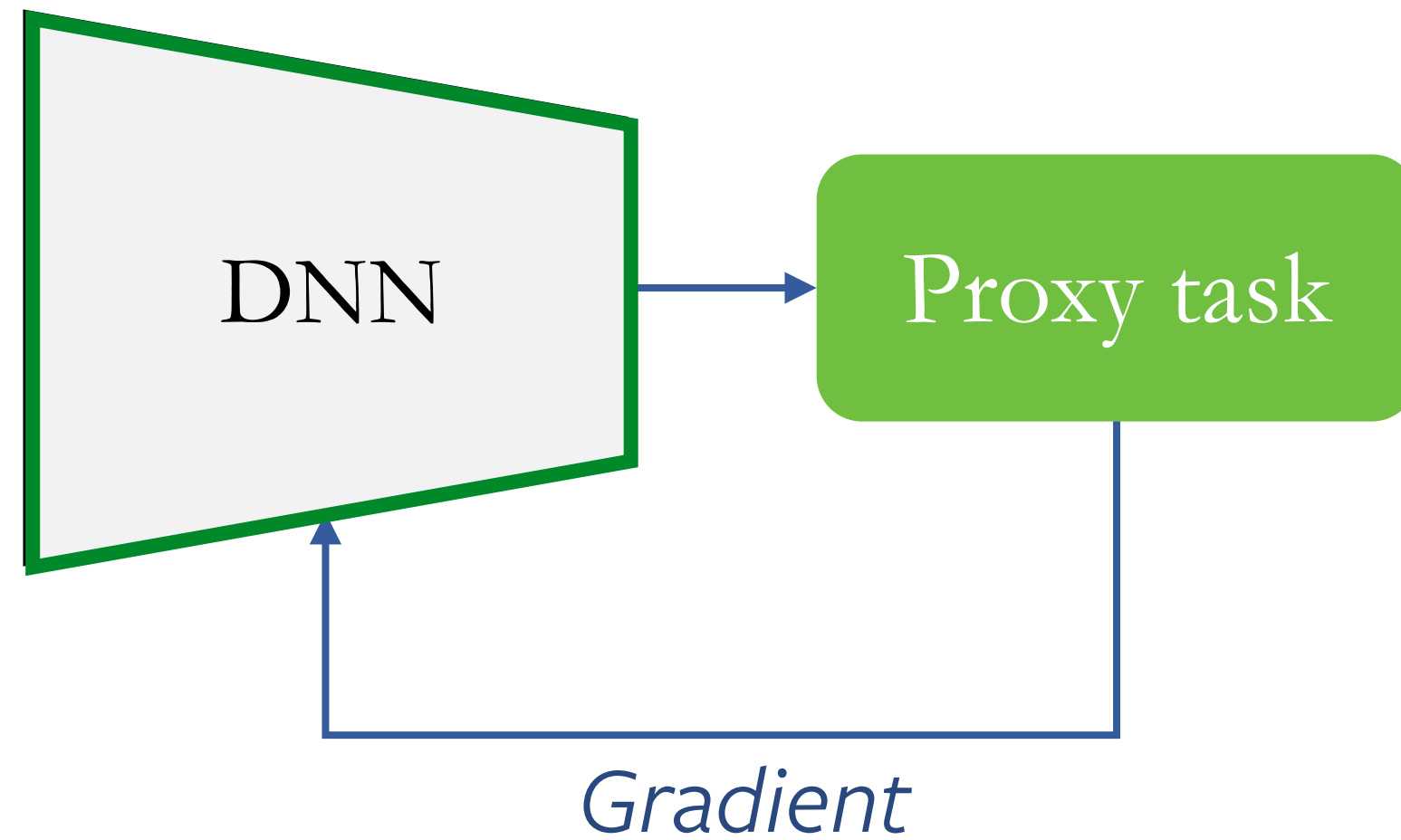


# General procedure of self-supervised learning.

## Phase 1: Pretraining



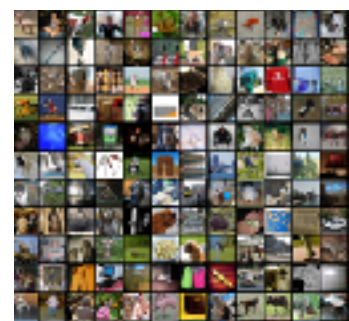
Unlabelled data  
+ transformations



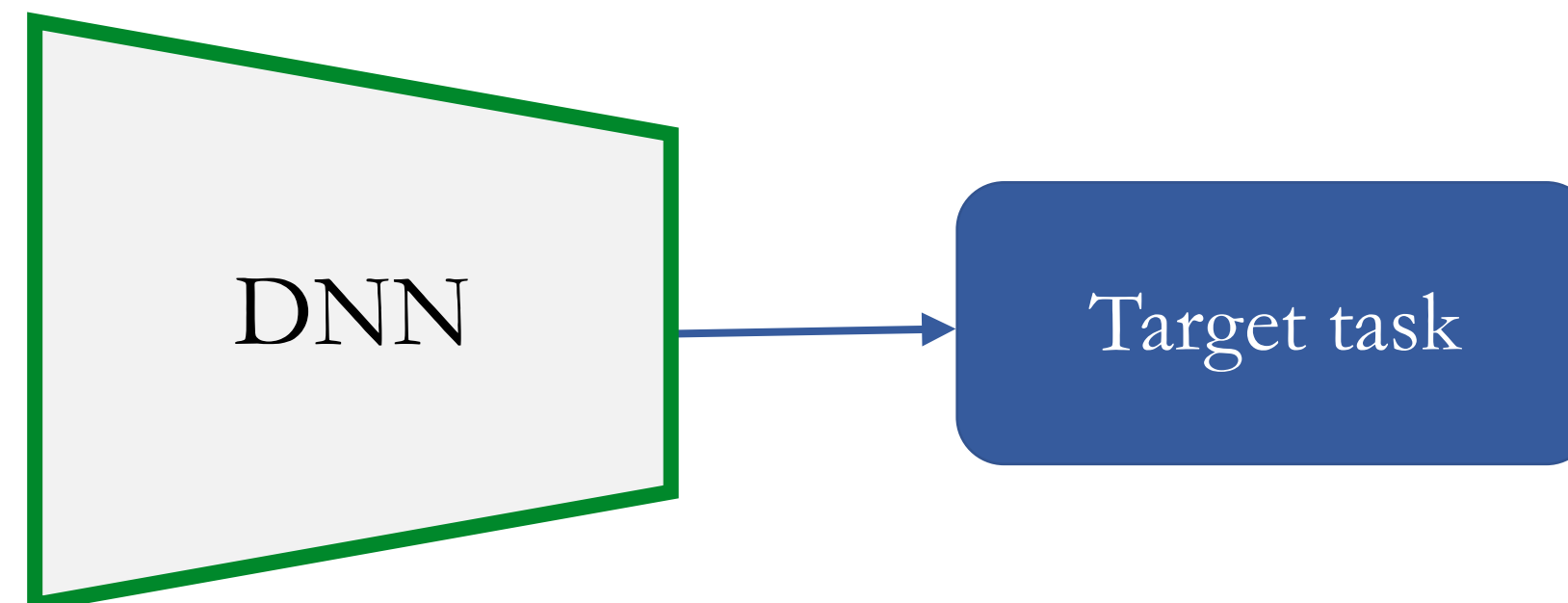
For example:

- 1) rotate image by multiple of 90 degrees
- 2) have network predict the rotation

## Phase 2: Downstream tasks



(Sparse) labeled data



Typically has less data

Why do we want to do self-supervised learning?



# Reason 1: Scalability and GPT as proof-of-concept

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 3: Ambiguity of labels and captions



"A house"?



"A boat"?

Example from Flickr30k



*A hot, blond girl getting criticized by her boss.*

Problematic captions

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

# Especially videos open exciting new directions



Visual development for AI



"Get" physics



Embodied AI

Bonus: *insane* scale:



*Is ImageNet worth 1 video? Learning  
strong image encoders from 1 long  
unlabelled video.*

VENKATARAMANAN, RIZVE, CARREIRA, AVRITHIS\*, ASANO\*.

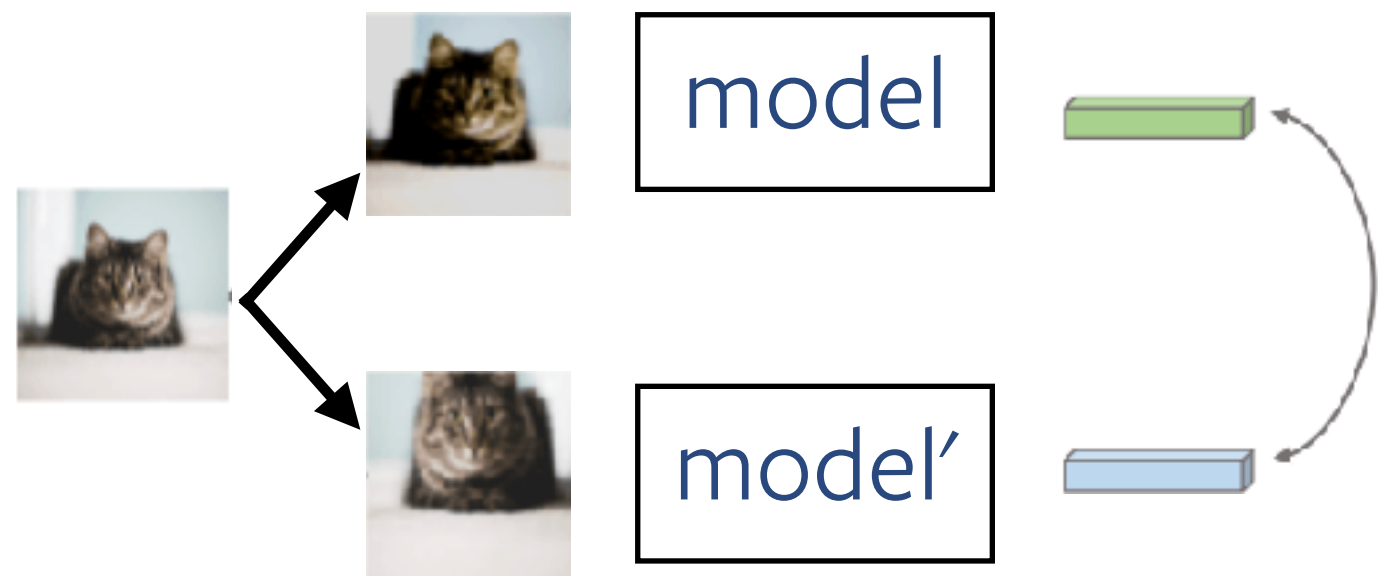
ICLR 2024



UNIVERSITY  
OF AMSTERDAM

# Augmentations are crucial in classic image-SSL, but forcing frames to be invariant is limiting

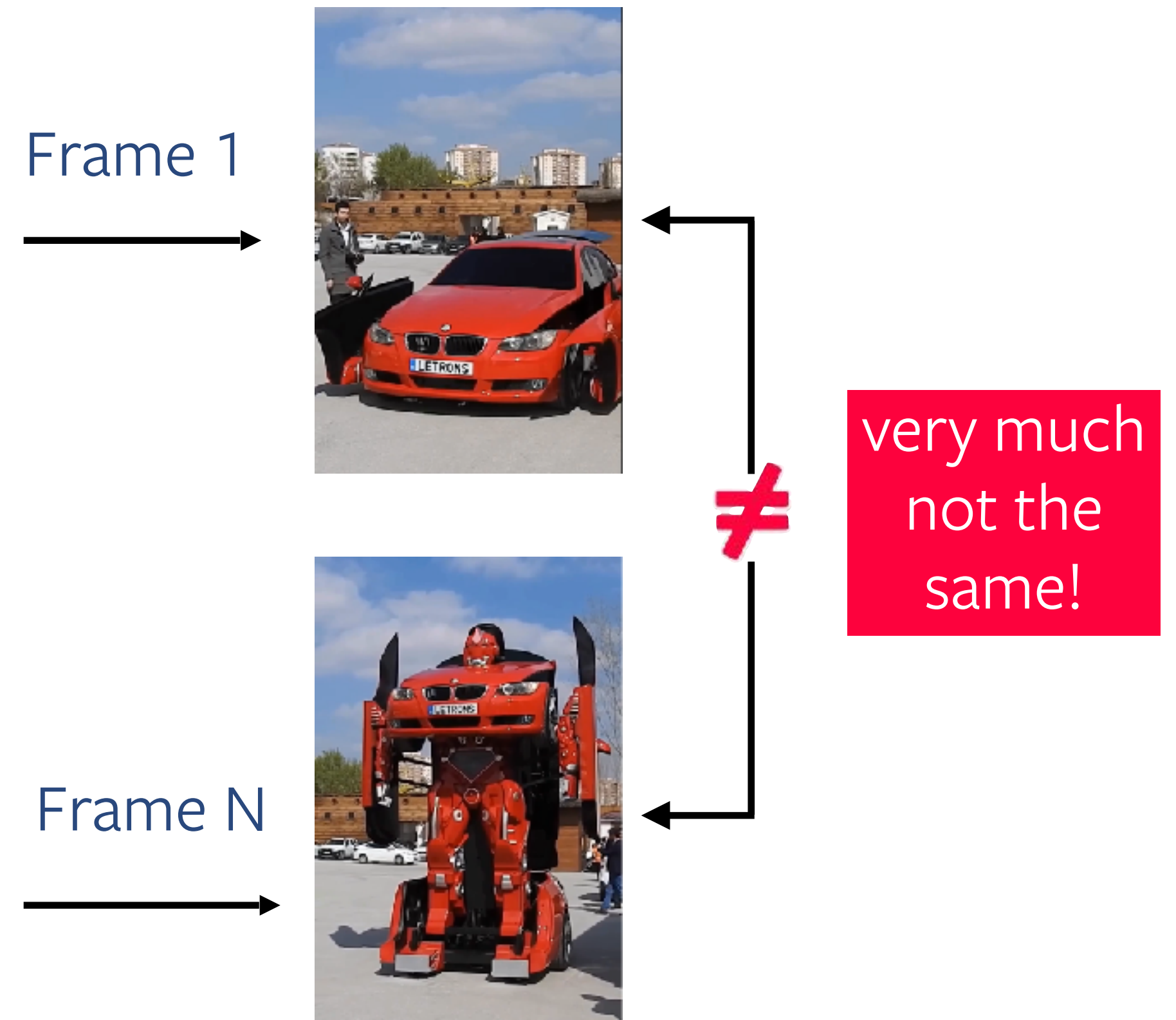
Images: SimCLR, MoCo, SwaAV et al.



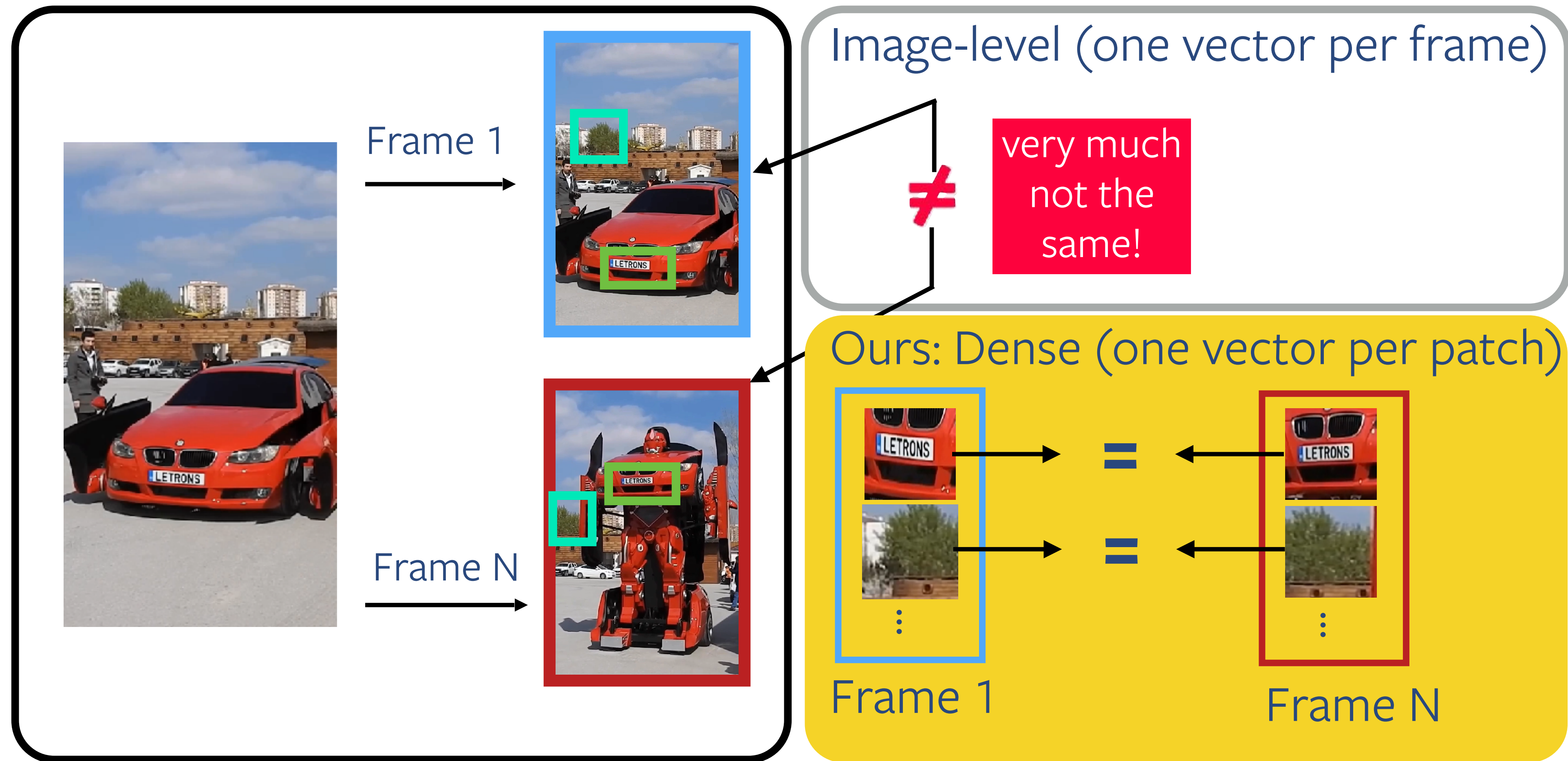
key principle: view-invariance

Might be ok for videos like this:

But does this generally make sense?



# Solution is obvious



## TimeTuning:

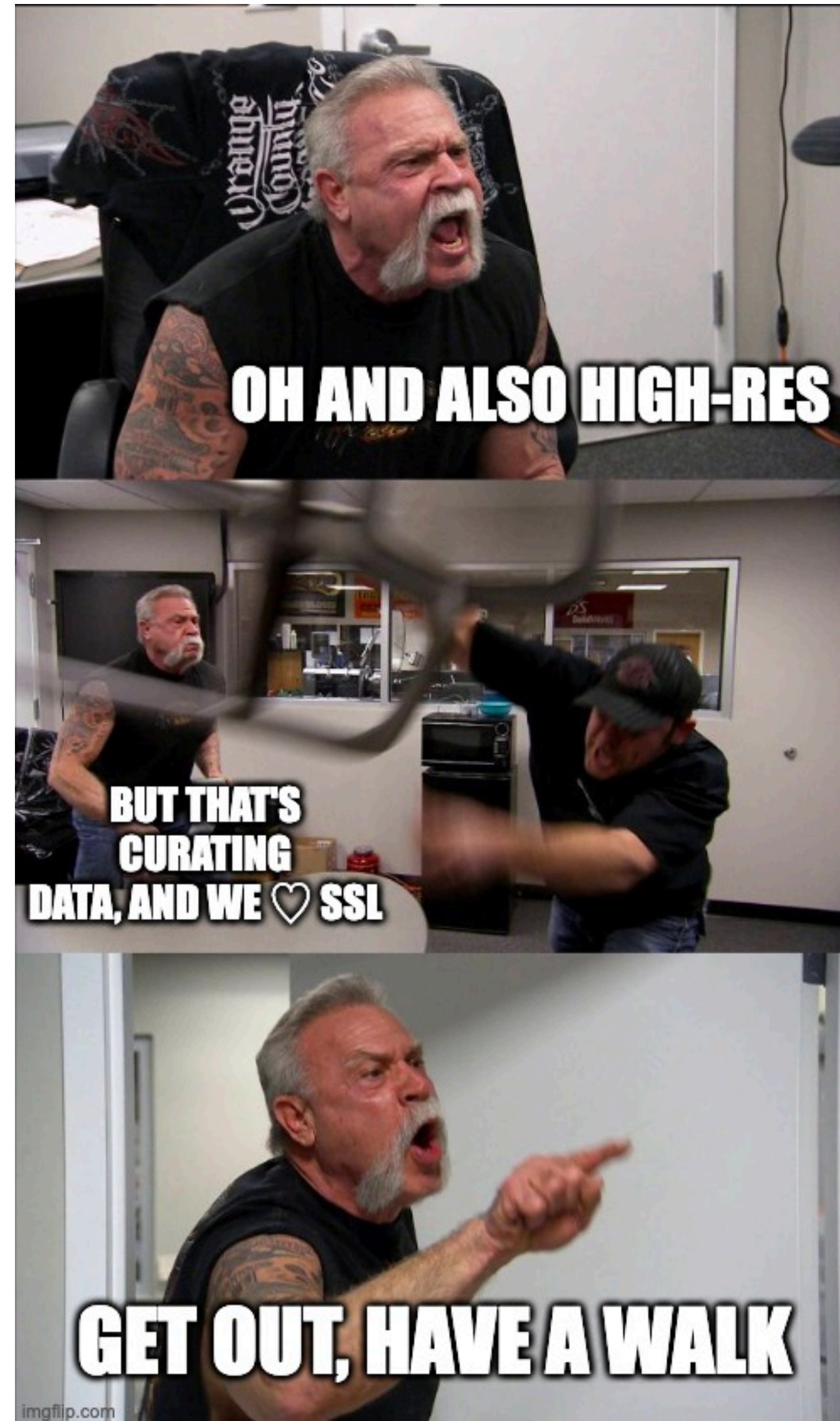
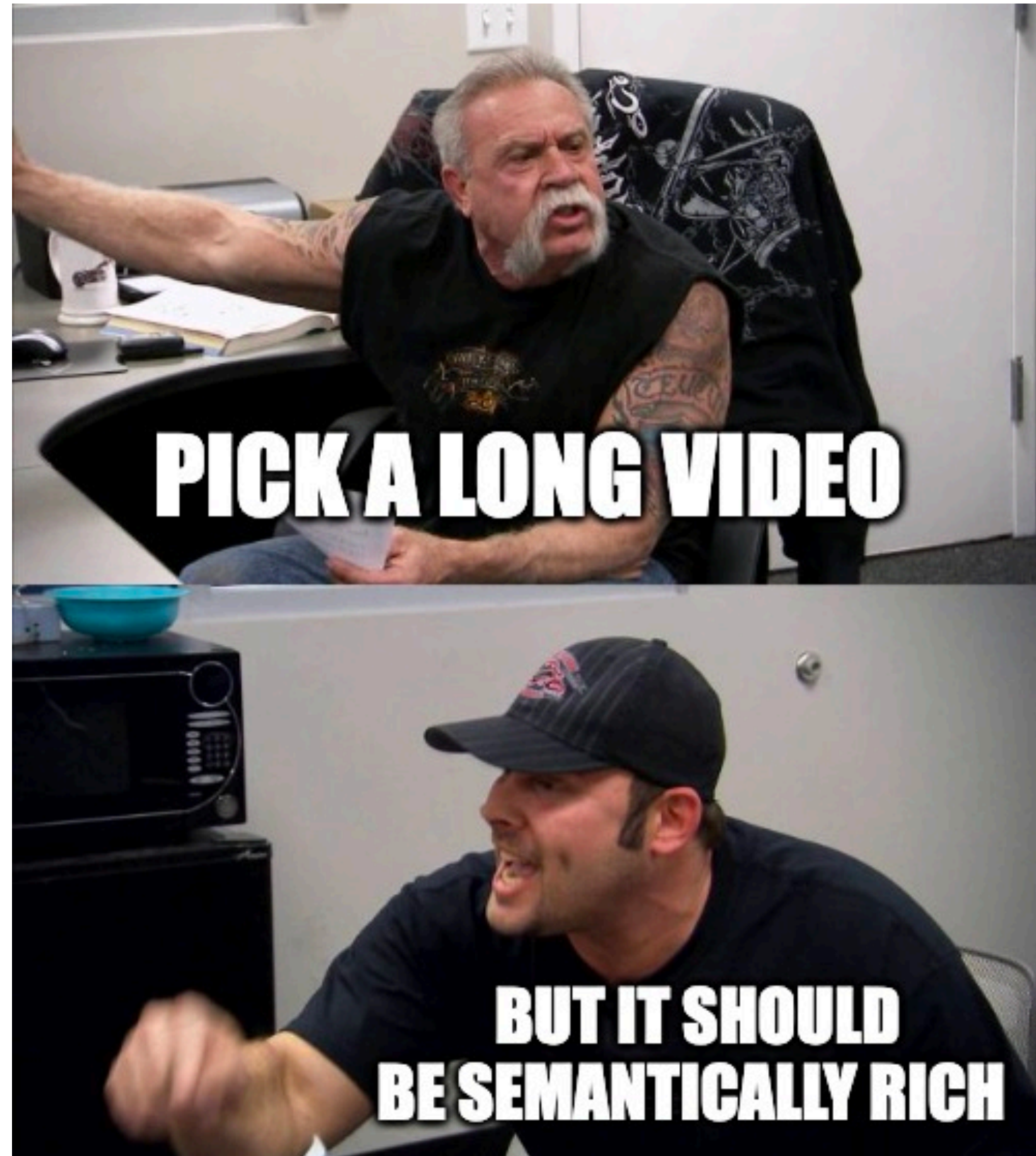
pretrained model &  
use temporal info of  
videos.

How powerful is time  
without image-pretraining?



Study the extreme:  
try to learn from a  
**single video,**  
**from scratch.**

# Us figuring out which video to use



- ✓ Long
- ✓ High-res, smooth
- ✓ Semantically rich
- ✓ Scalable (for SSL)

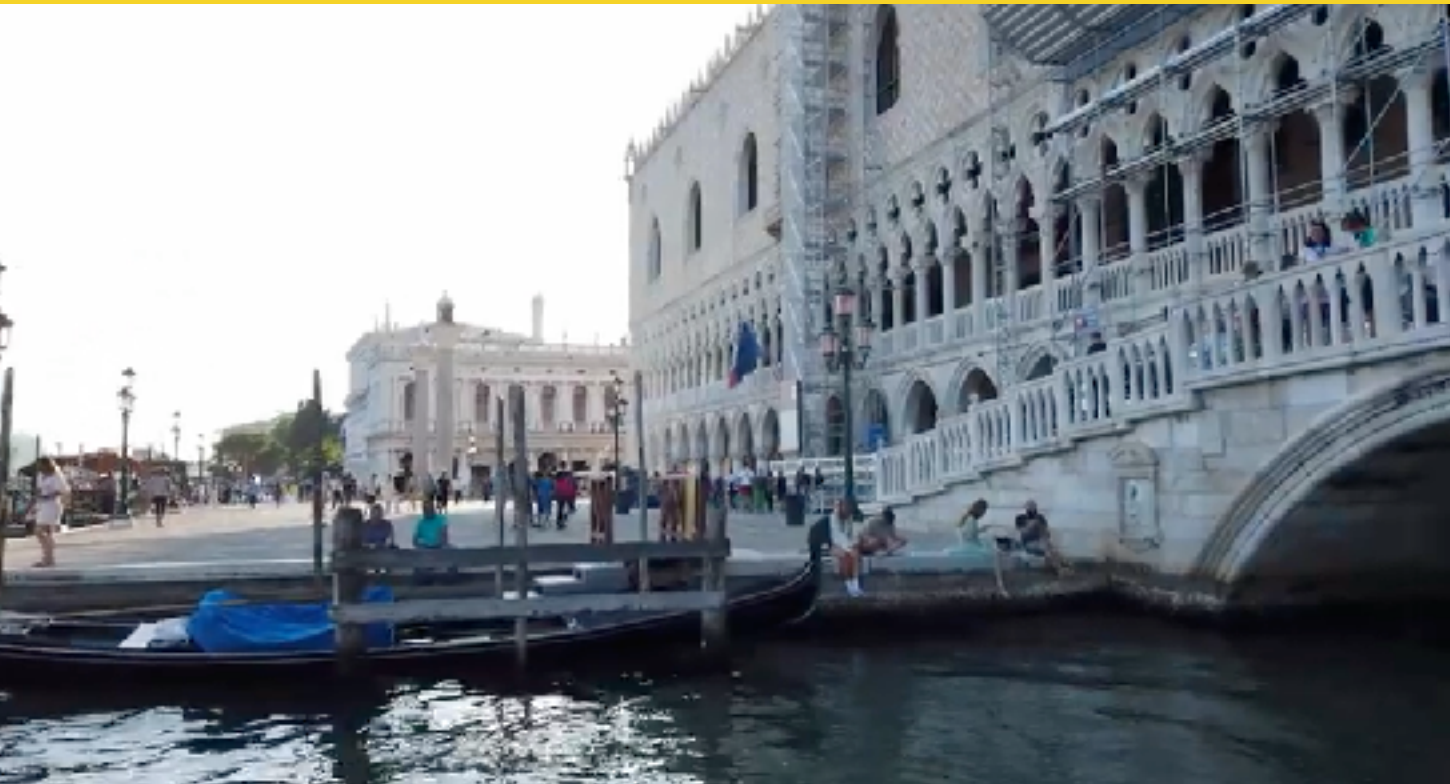


Walking Tours





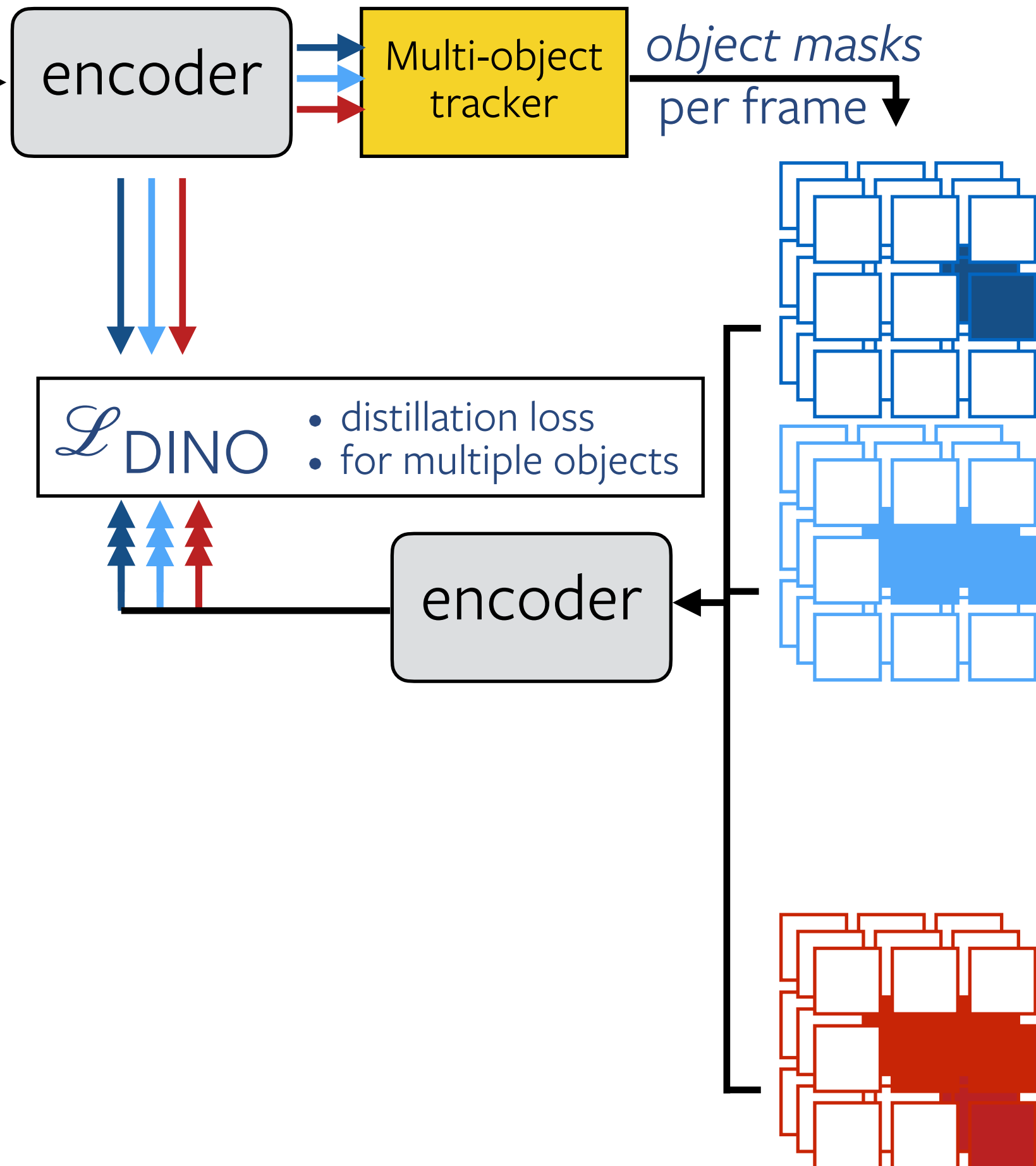
The dataset consists of 10x 4K videos of different cities' Walking Tours.



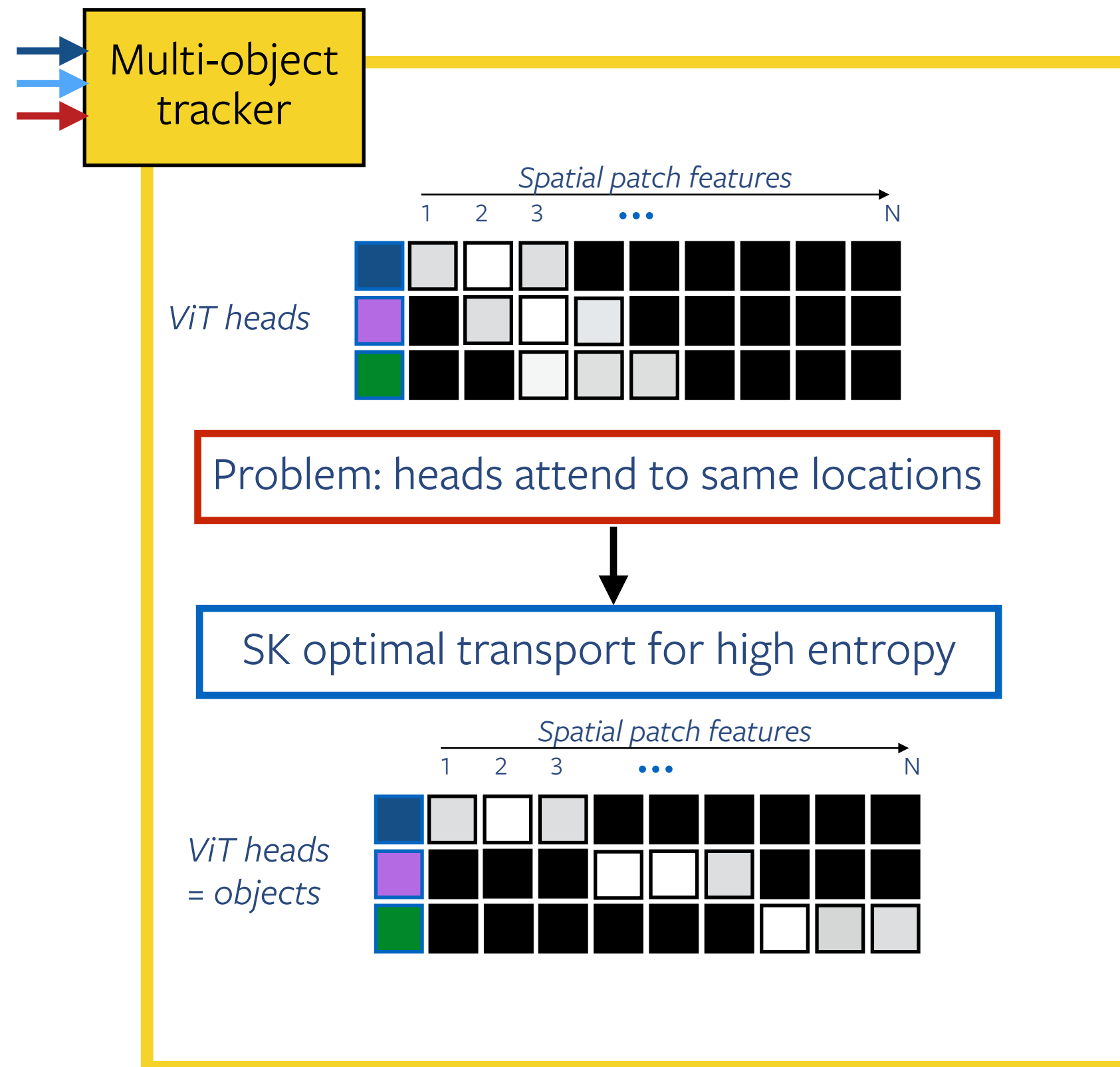
# Dora: **D**iscover and **T**rack



Much like Dora, we walk around and learn from what we see.



# Spreading attention with Sinkhorn-Knopp



Visualise attention of 3 heads with colors R,G,B

$t = 1$

$t = 2$

$t = 3$

$X_t$



without SK

$T_t$

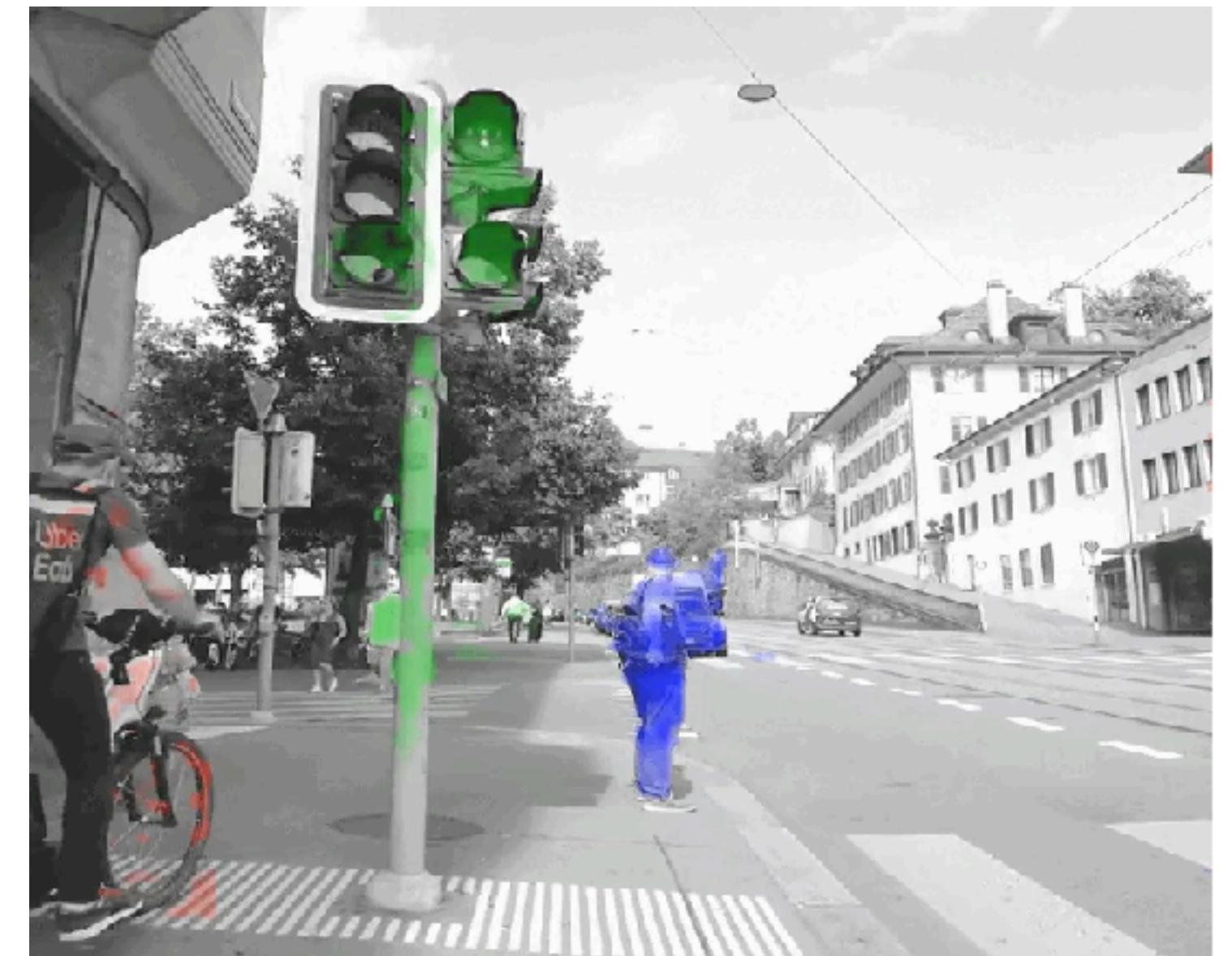
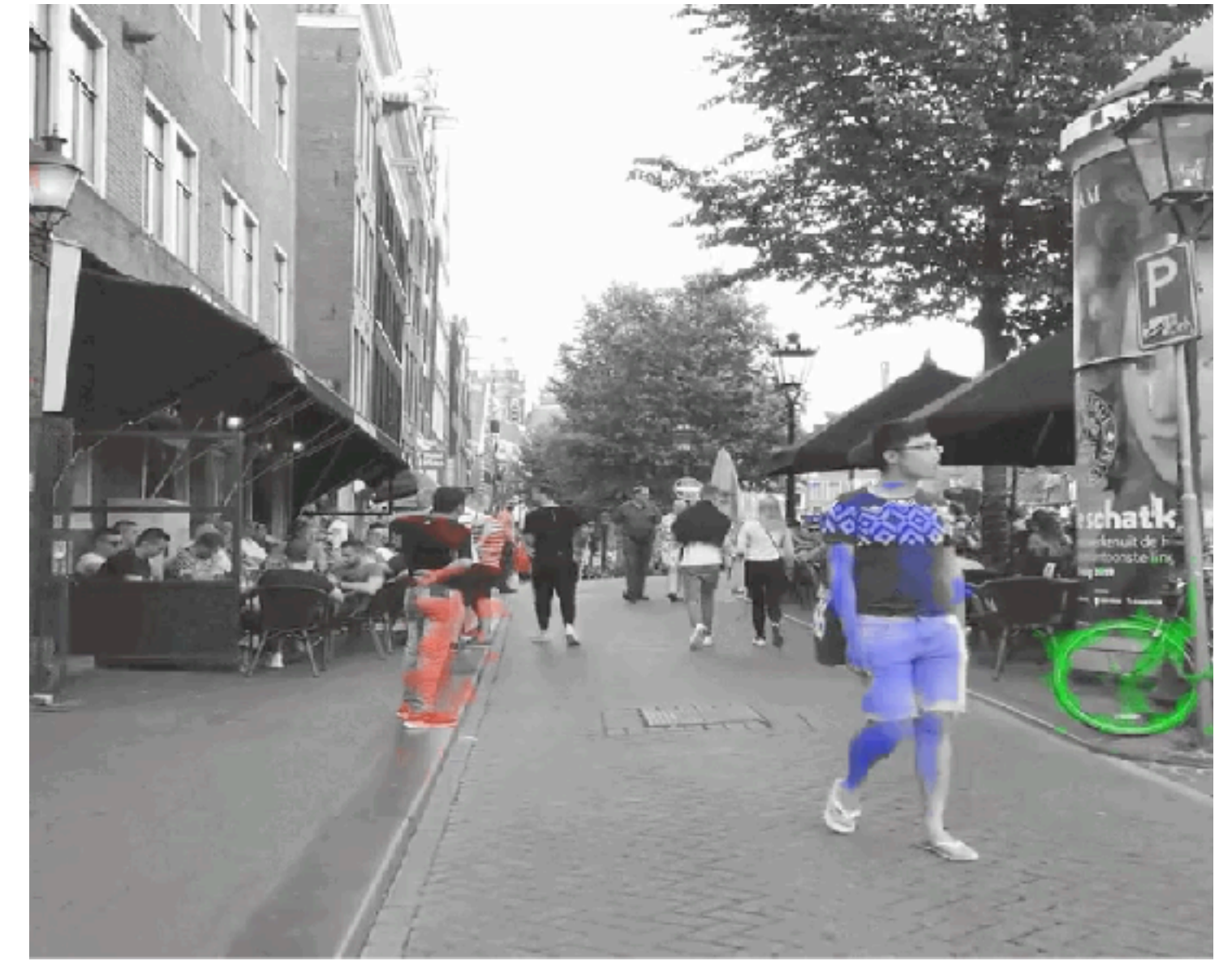


with SK

$T'_t$



# More examples: multi-object tracking in a ViT *emerges*

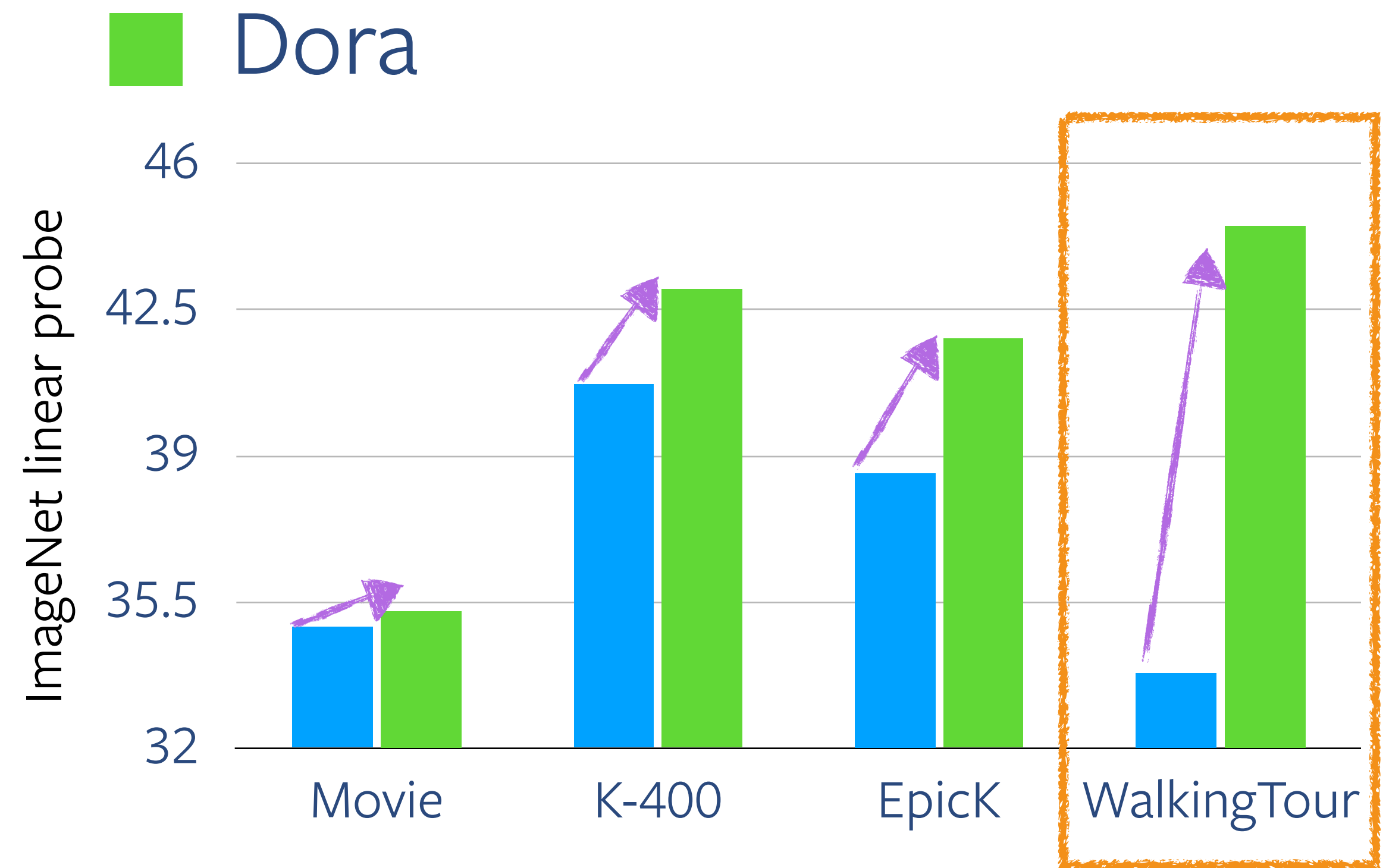
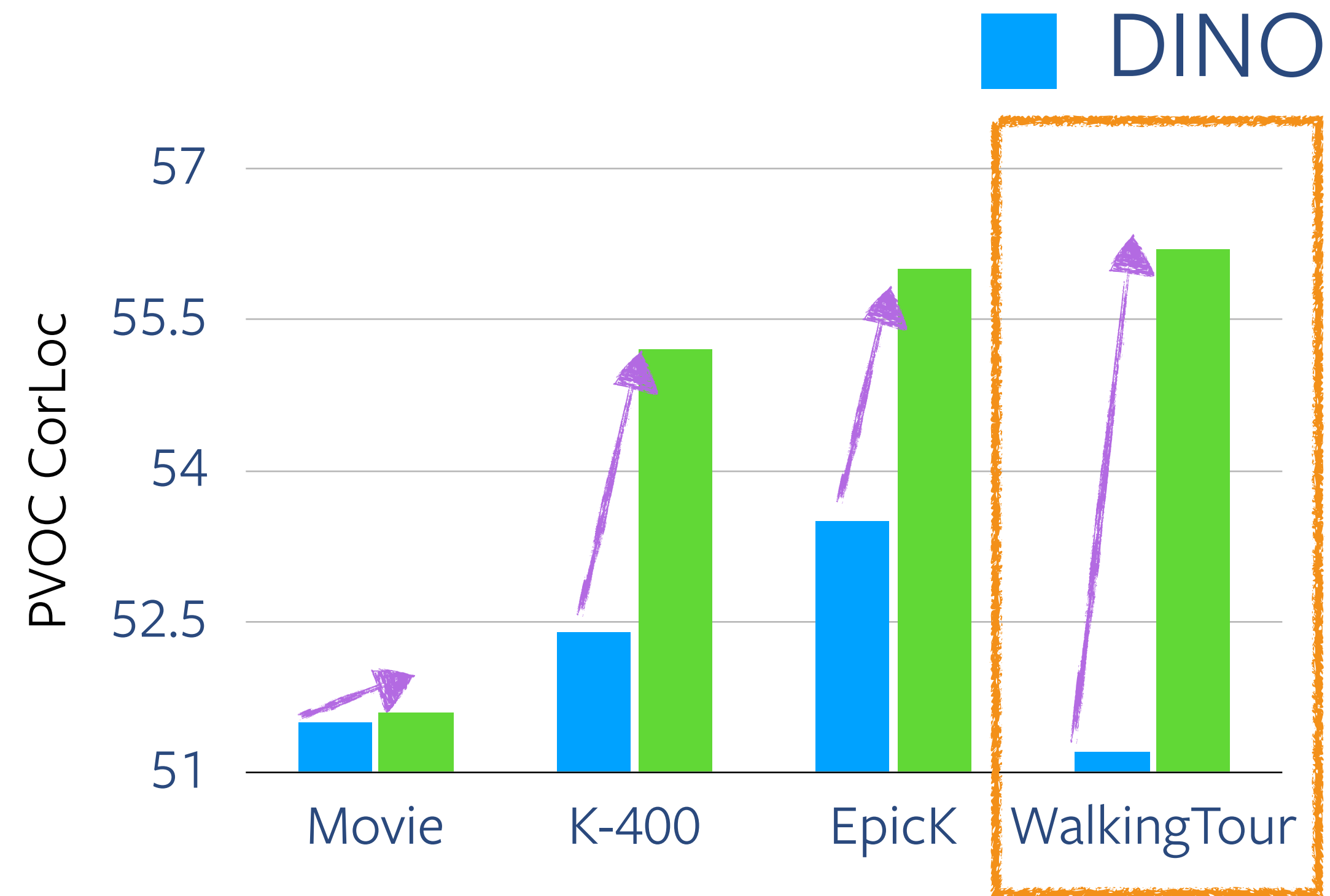


Can we obtain performances better than training on images?



# Dora better than DINO

## WT+ Dora: great match



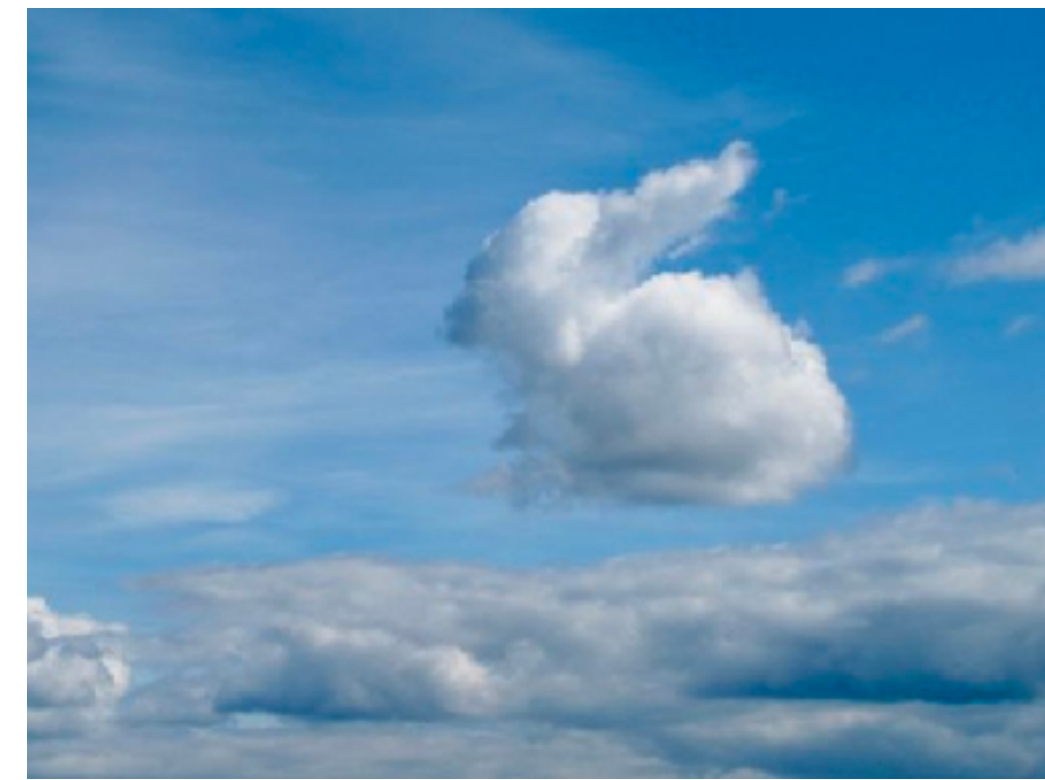
*The Augmented Image Prior: Distilling  
1000 Classes by Extrapolating from a  
Single Image.*

SAEED\*, ASANO\*.

ICLR 2023



UNIVERSITY  
OF AMSTERDAM



≈



n02325366

Wood rabbit



≈

0.4 x

n02087394

0.2 x

n01443537

0.1 x

n07697537

...

n032

!?

n





# How can we test this fairytale?

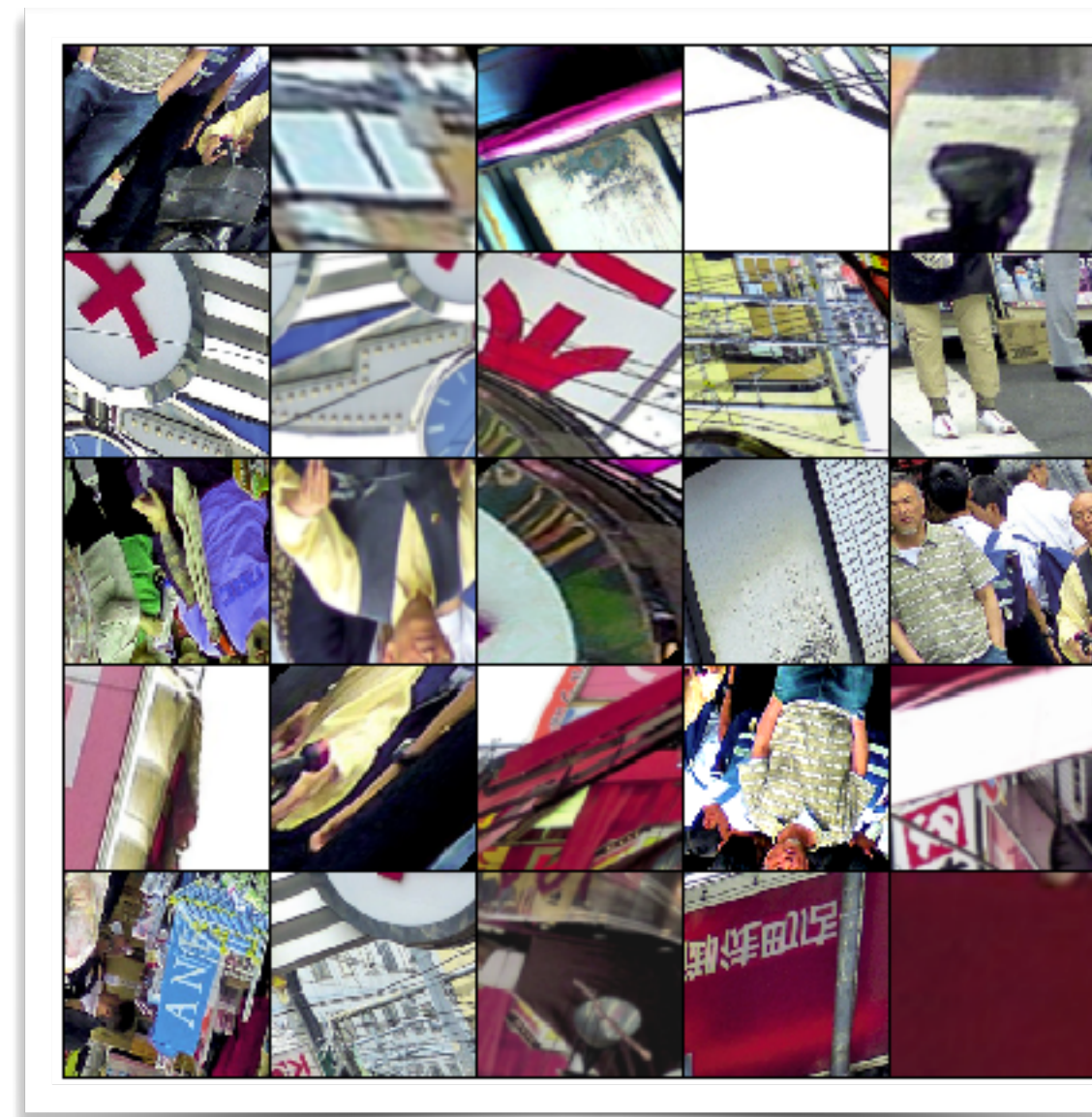
A single image  $I$

Pretrained neural network (teacher)

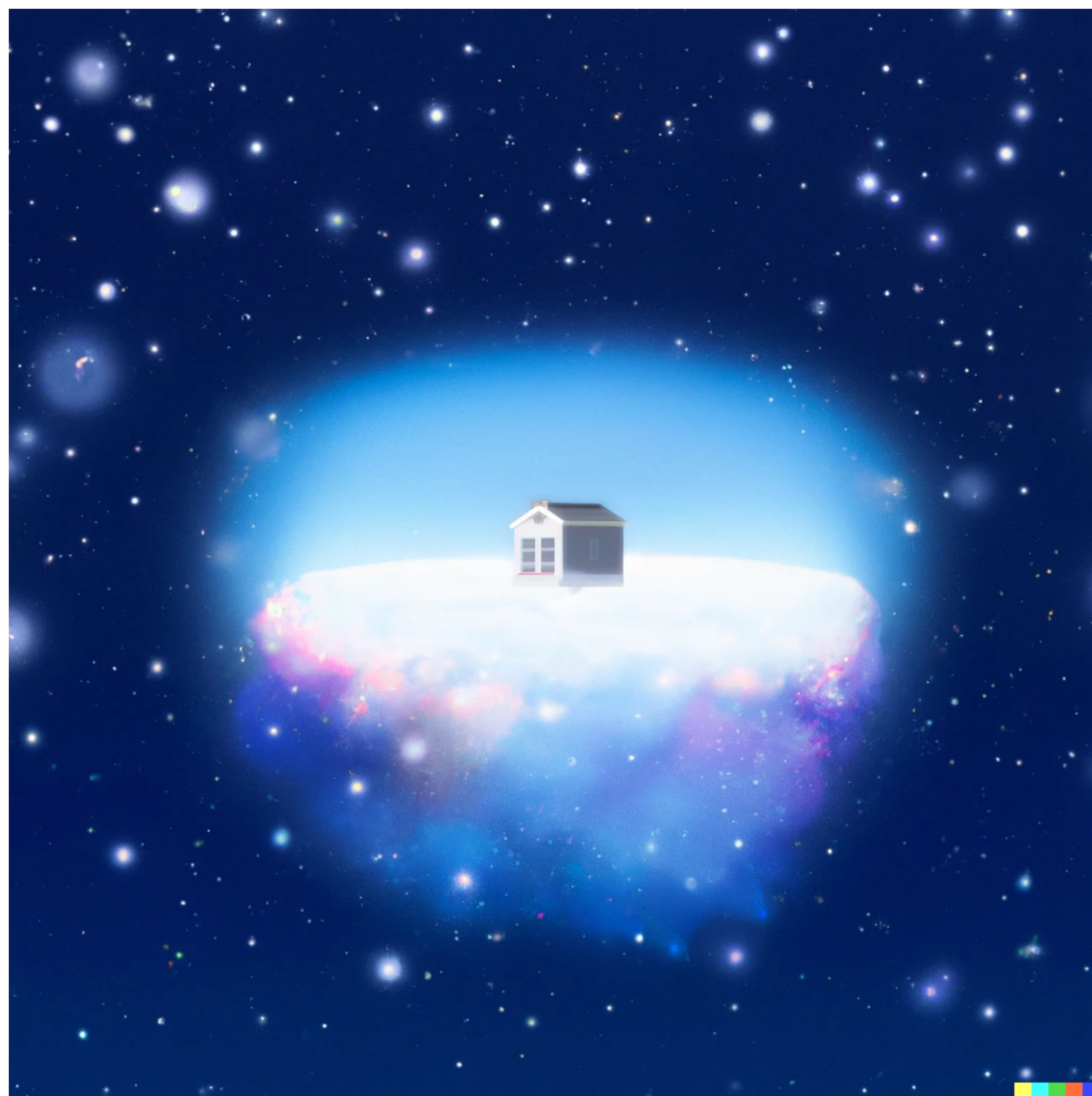
Randomly initialized network (student)

Augmentations  $\mathcal{A}(I)$

$\approx$  ImageNet,  
 $\approx$  Kinetics  
 $\approx$  UCF  
 $\approx$  ..



# The *real* motivation

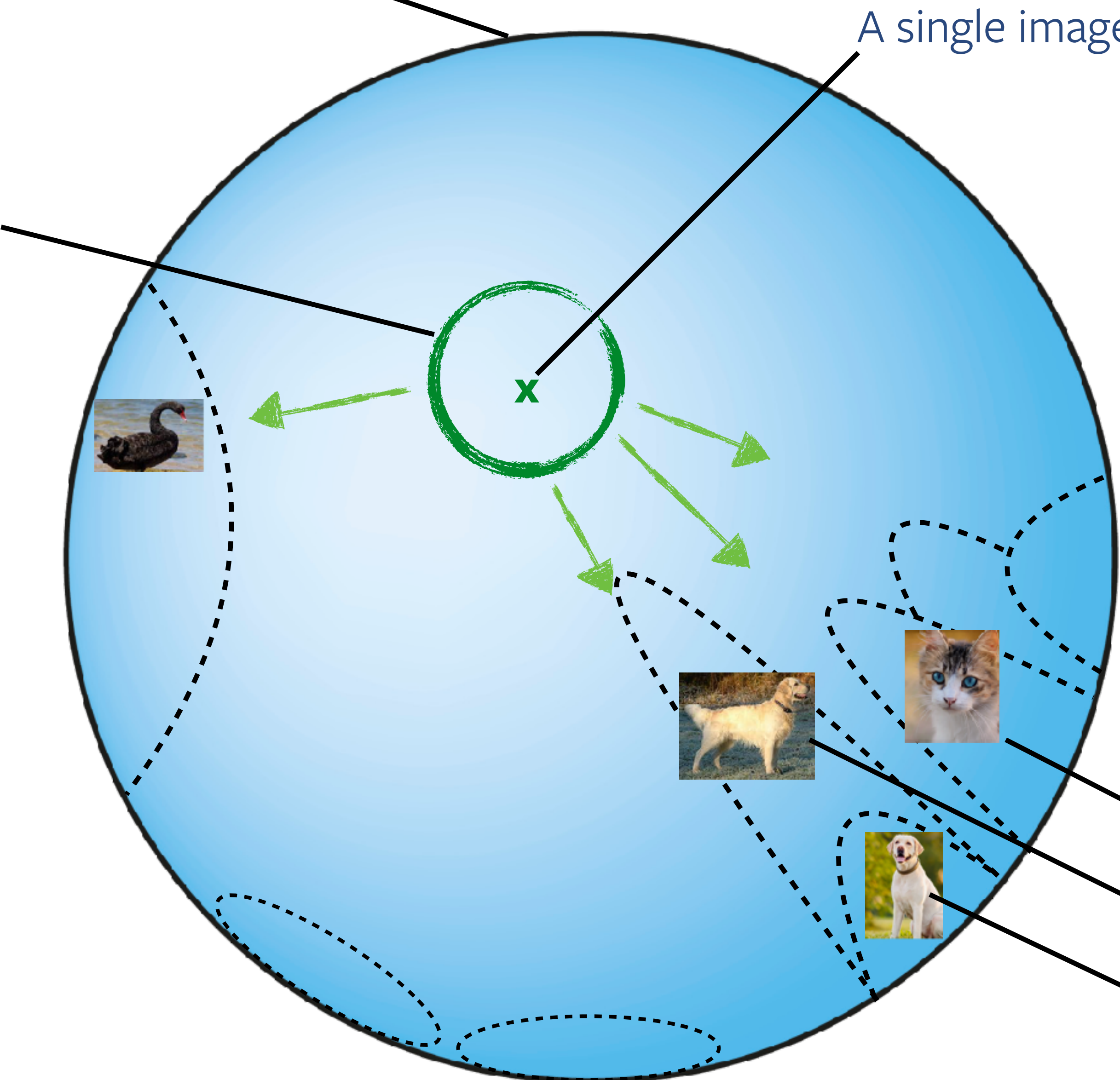
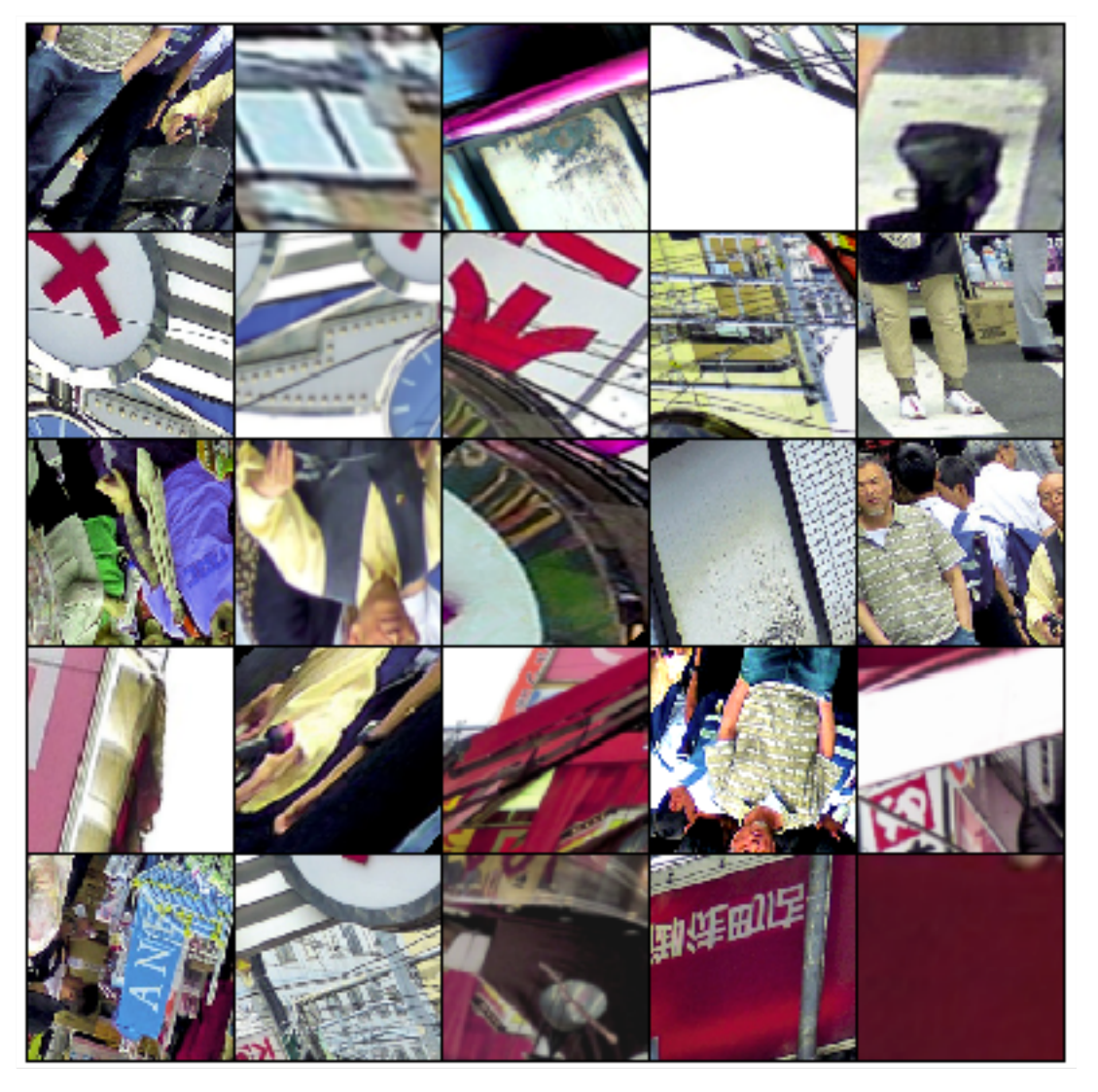


# Can we extrapolate from a single image to semantic categories?

All natural images

A single image  $I$

Augmentations  $\mathcal{A}(I)$



**ImageNet-12:**  
1000 semantic classes

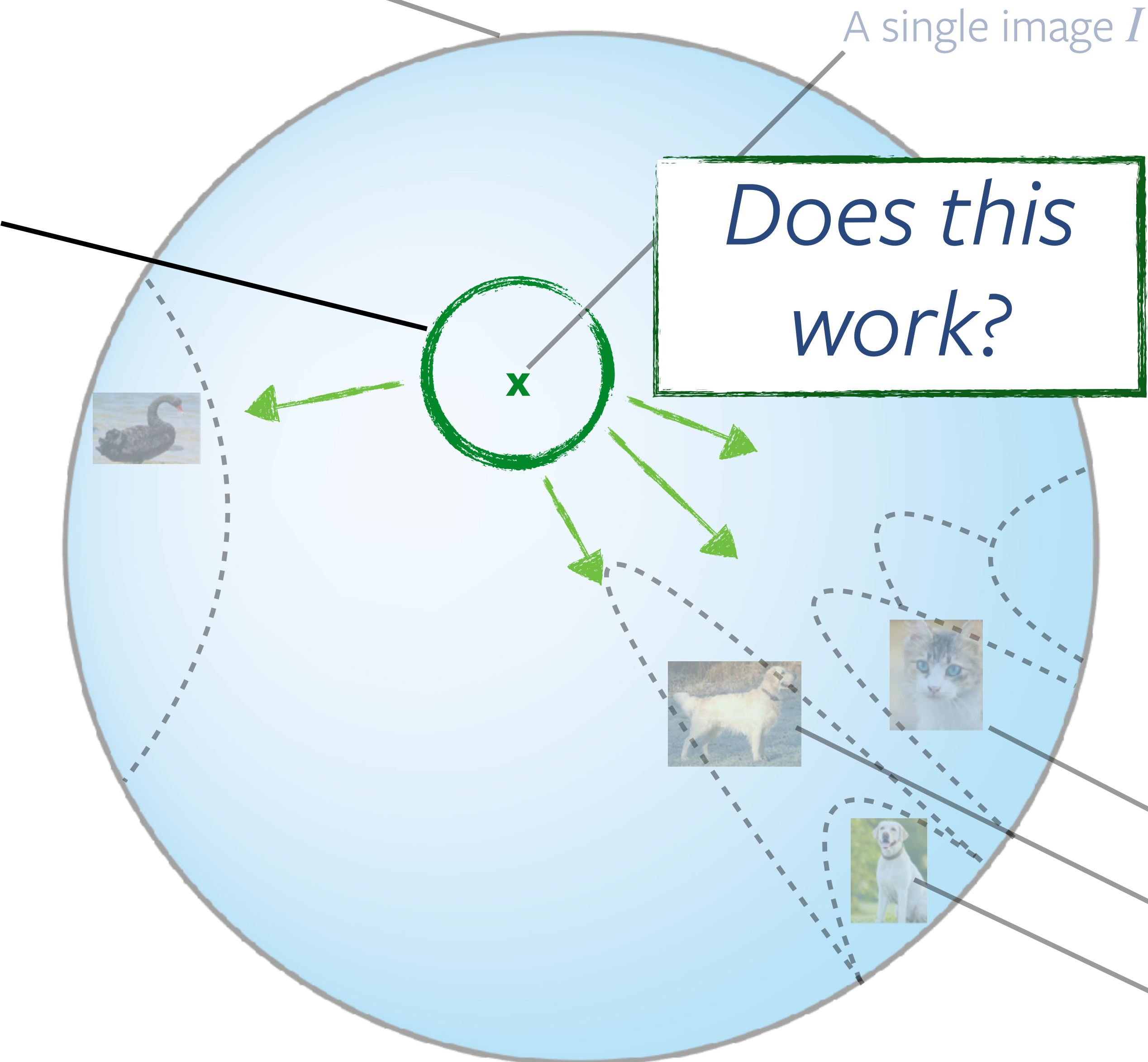
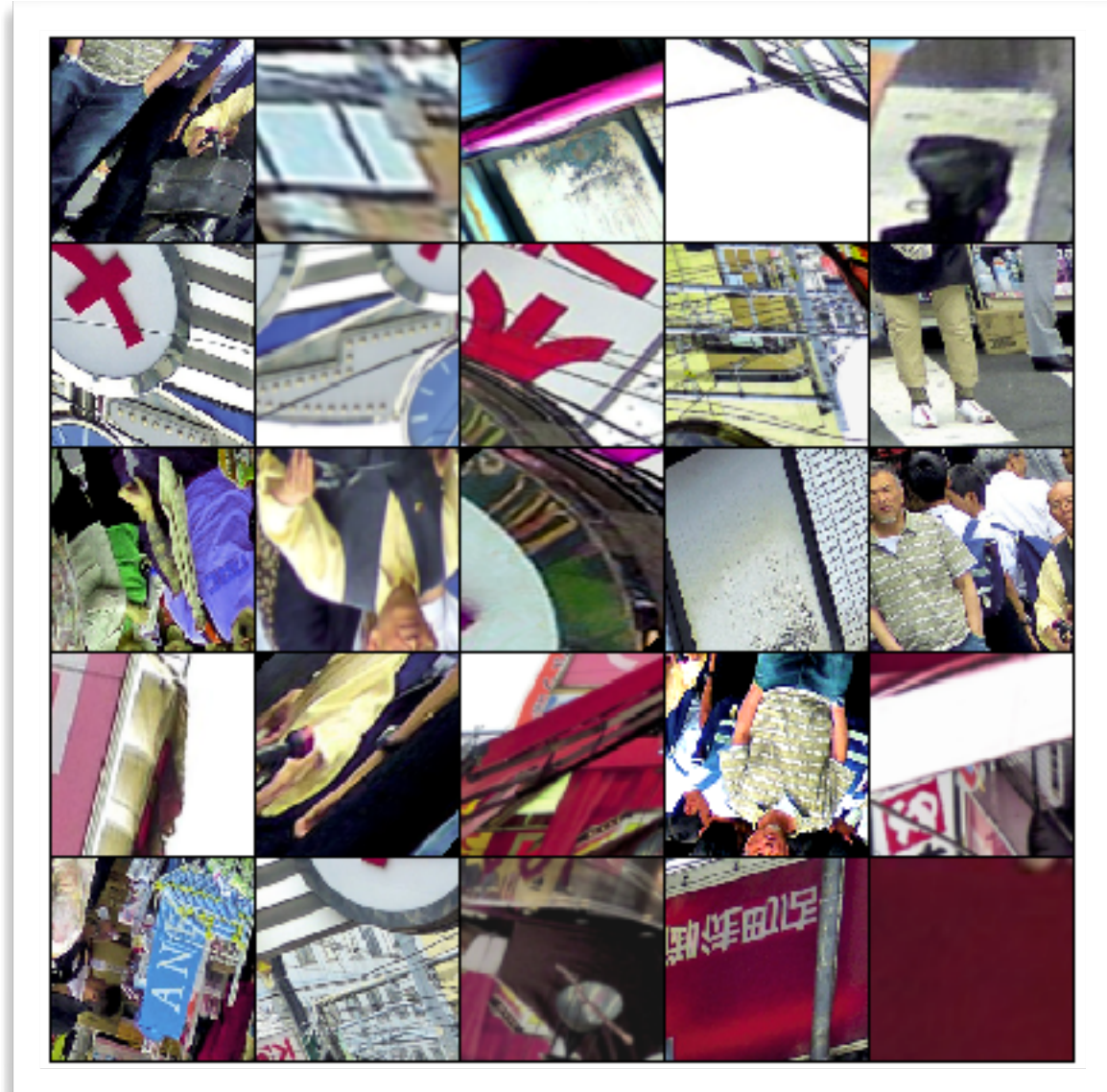
- "cat"
- "Golden retriever"
- "Labrador retriever"

# Can we extrapolate from a single image to semantic categories?

All natural images

A single image  $I$

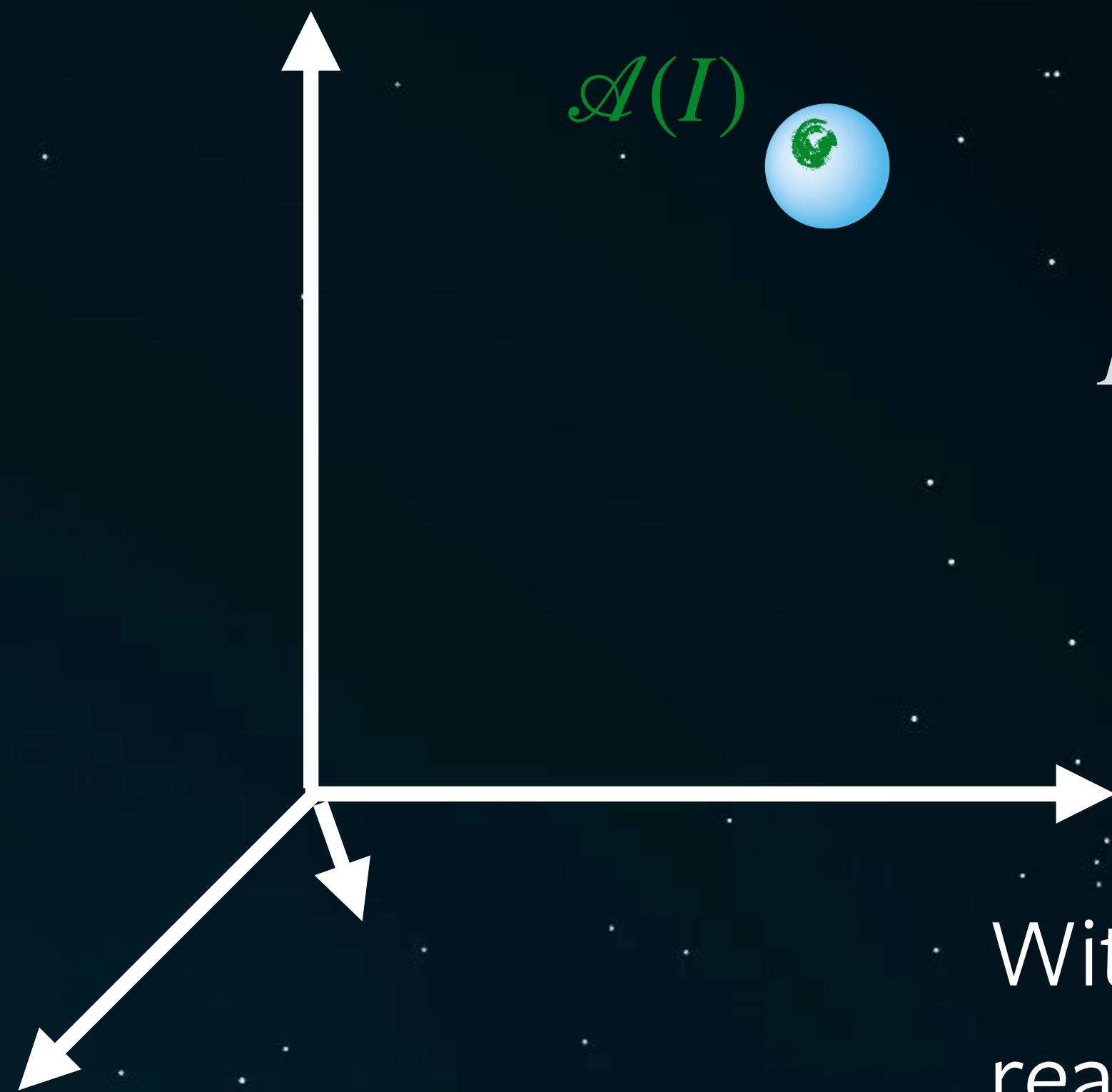
Augmentations  $\mathcal{A}(I)$



**ImageNet-12:**  
1000 semantic classes

- "cat"
- "Golden retriever"
- "Labrador retriever"

# Why it might work: *The augmented image prior*



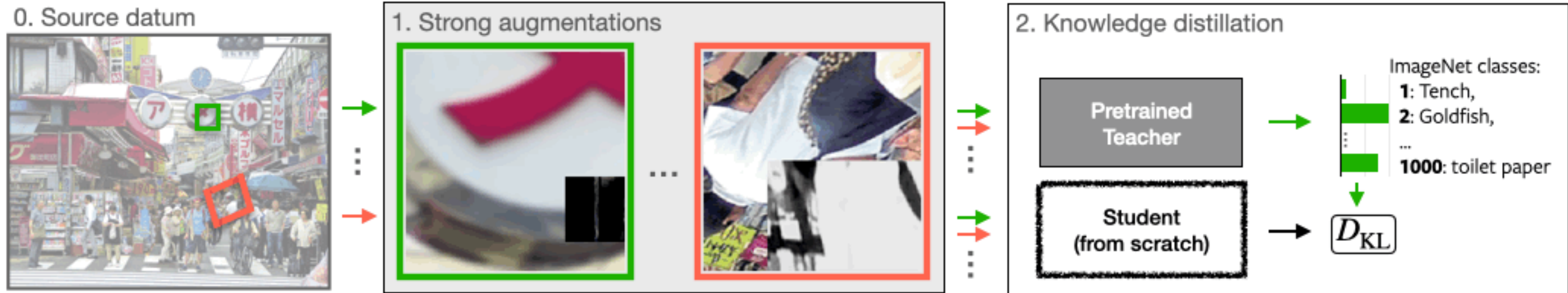
$$I \in \mathcal{I}, \quad \mathcal{I} = \{0, \dots, 255\}^{3 \times 224 \times 224}$$

Within the space of all possible images  $\mathcal{I}$ , a single real image  $I$  and its augmentations  $\mathcal{A}(I)$  provide a very informative prior about all **real images** for extrapolation

# *Method*



Our pipeline is kept as simple as possible.



Training data: very varied, as are the teacher's predictions.



0.2803%: packet (692)  
0.1994%: rubber eraser (767)  
0.1947%: envelope (549)  
0.1910%: Band Aid (419)  
0.1865%: lighter (626)



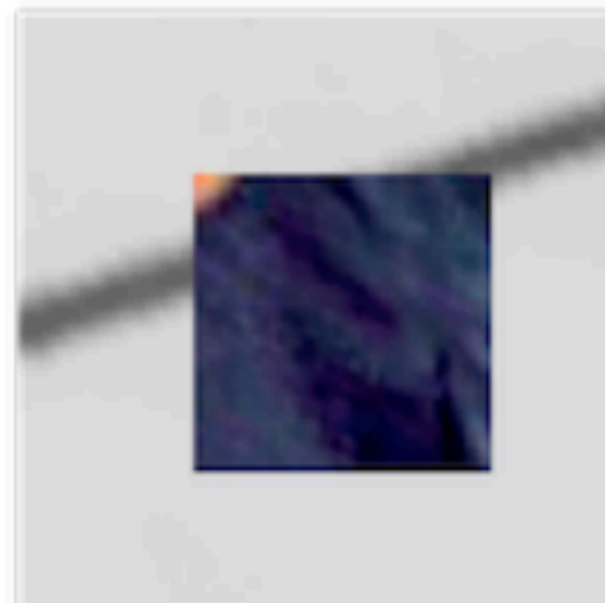
0.1563%: cleaver (499)  
0.1549%: can opener (473)  
0.1521%: whistle (902)  
0.1511%: screw (783)  
0.1507%: spatula (813)



0.2653%: comic book (917)  
0.2453%: sarong (775)  
0.2436%: sombrero (808)  
0.2344%: shopping basket (790)  
0.2336%: toyshop (865)



0.3454%: grocery store (582)  
0.3175%: jinrikisha (612)  
0.2830%: restaurant (762)  
0.2367%: toyshop (865)  
0.2322%: barbershop (424)



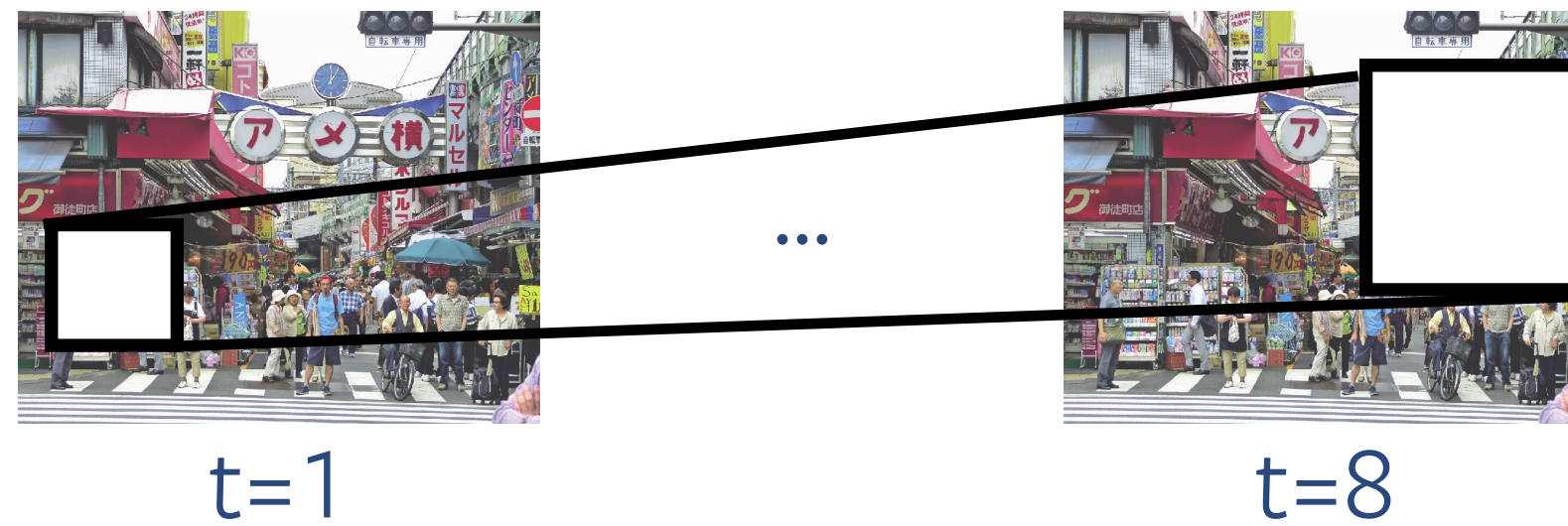
0.2736%: nail (677)  
0.2657%: screw (783)  
0.2121%: hook (600)  
0.2035%: knot (616)  
0.1901%: letter opener (623)



0.2557%: slide rule (798)  
0.2362%: rule (769)  
0.2135%: letter opener (623)  
0.2116%: Windsor tie (906)  
0.2061%: matchstick (644)



We can also make fake-videos out of images.



pick two crops, smoothly transition between them.



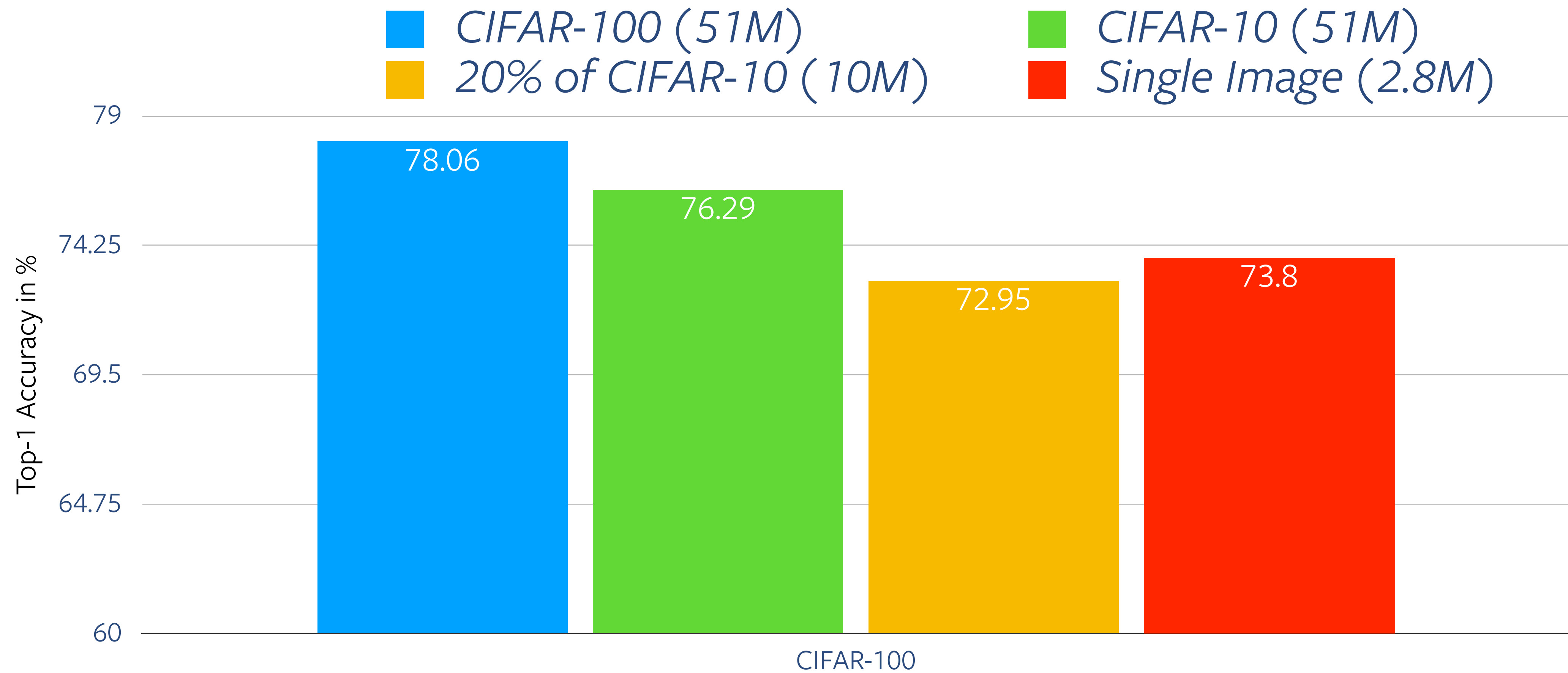
# Guessing game! [the actual training data]



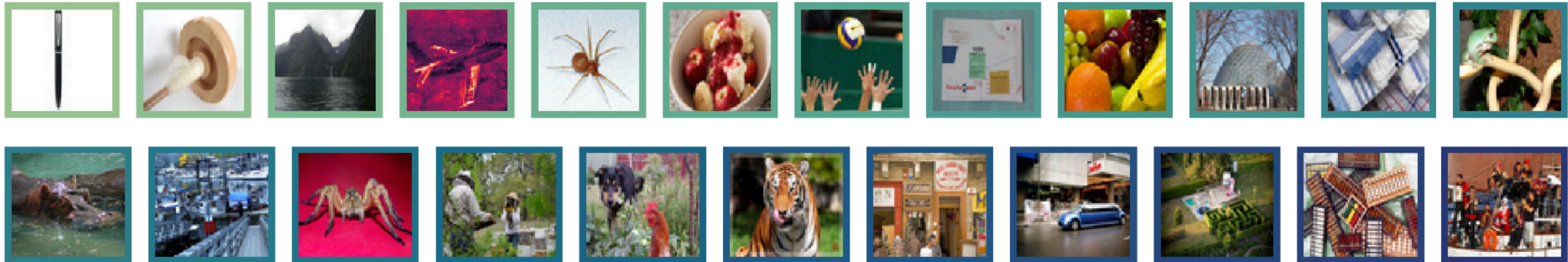
# *Results*



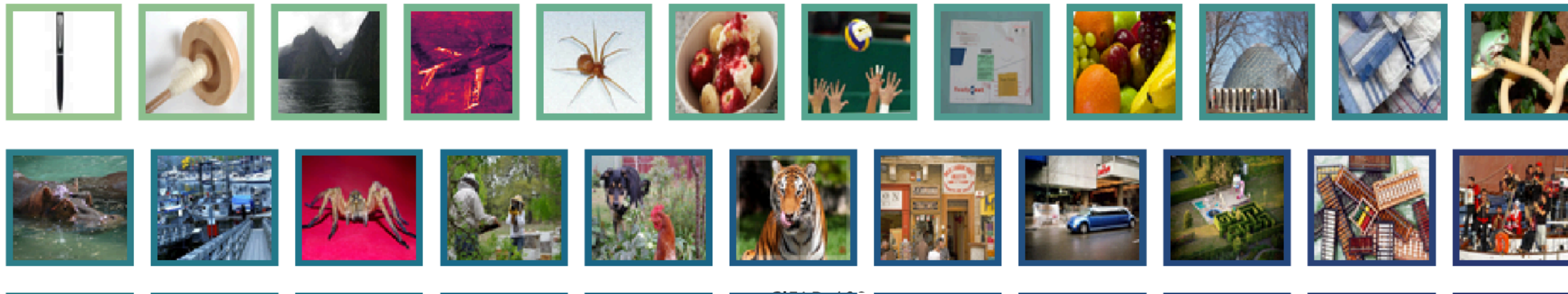
# Comparison of datasets (*number of pixels*)



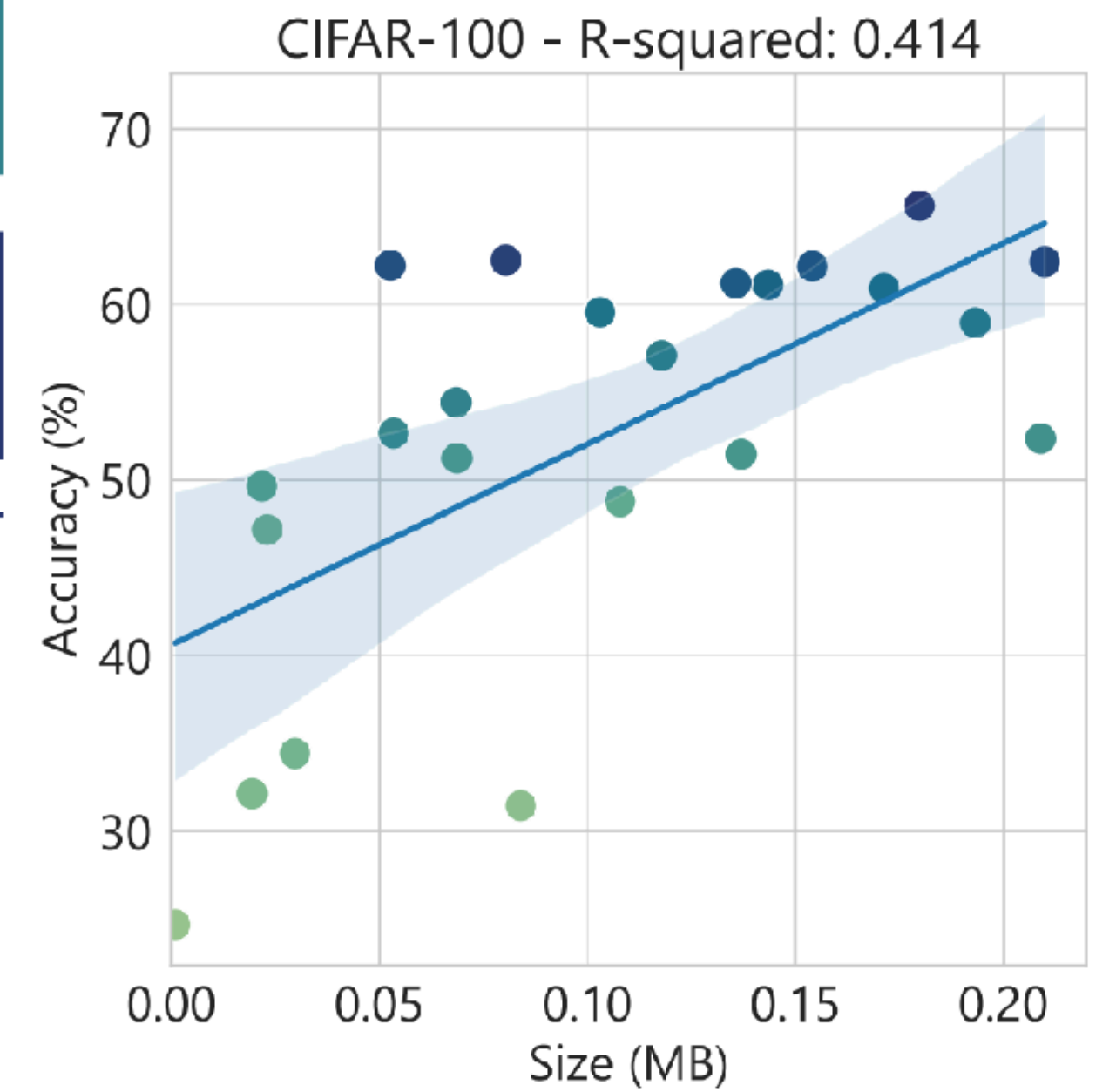
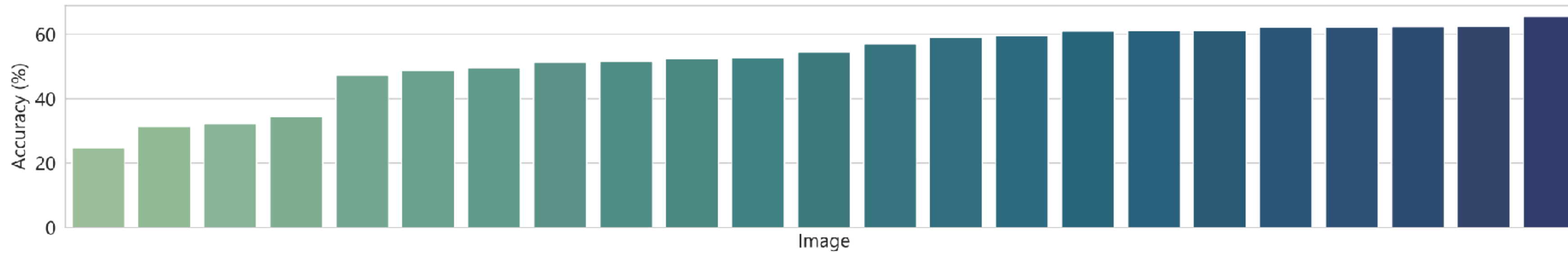
# Image complexity matters



# Image complexity matters



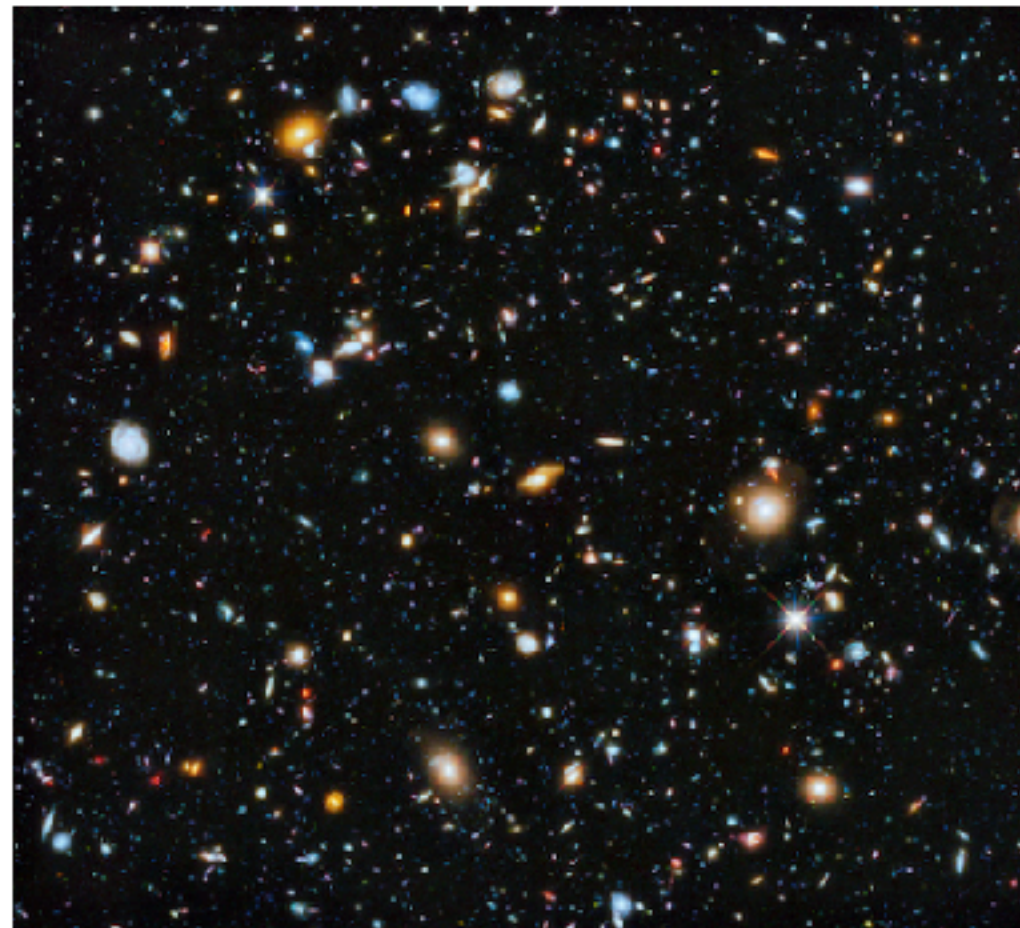
CIFAR-100



# Other images



(a) The “Noise” Image. From uniform noise [0,255]. Size: 2,560x1,920, PNG: 16.3MB.



(b) The “Universe” Image. Size: 2,300x2,100, JPEG: 7.2MB.



(c) The “Bridge” Image. Size: 1,280x853, JPEG: 288KB.



(d) The “City” Image. Size: 2,560x1,920, JPEG: 1.9MB.



(e) The “Animals” Image. Size: 1,300x600, JPEG: 267KB.

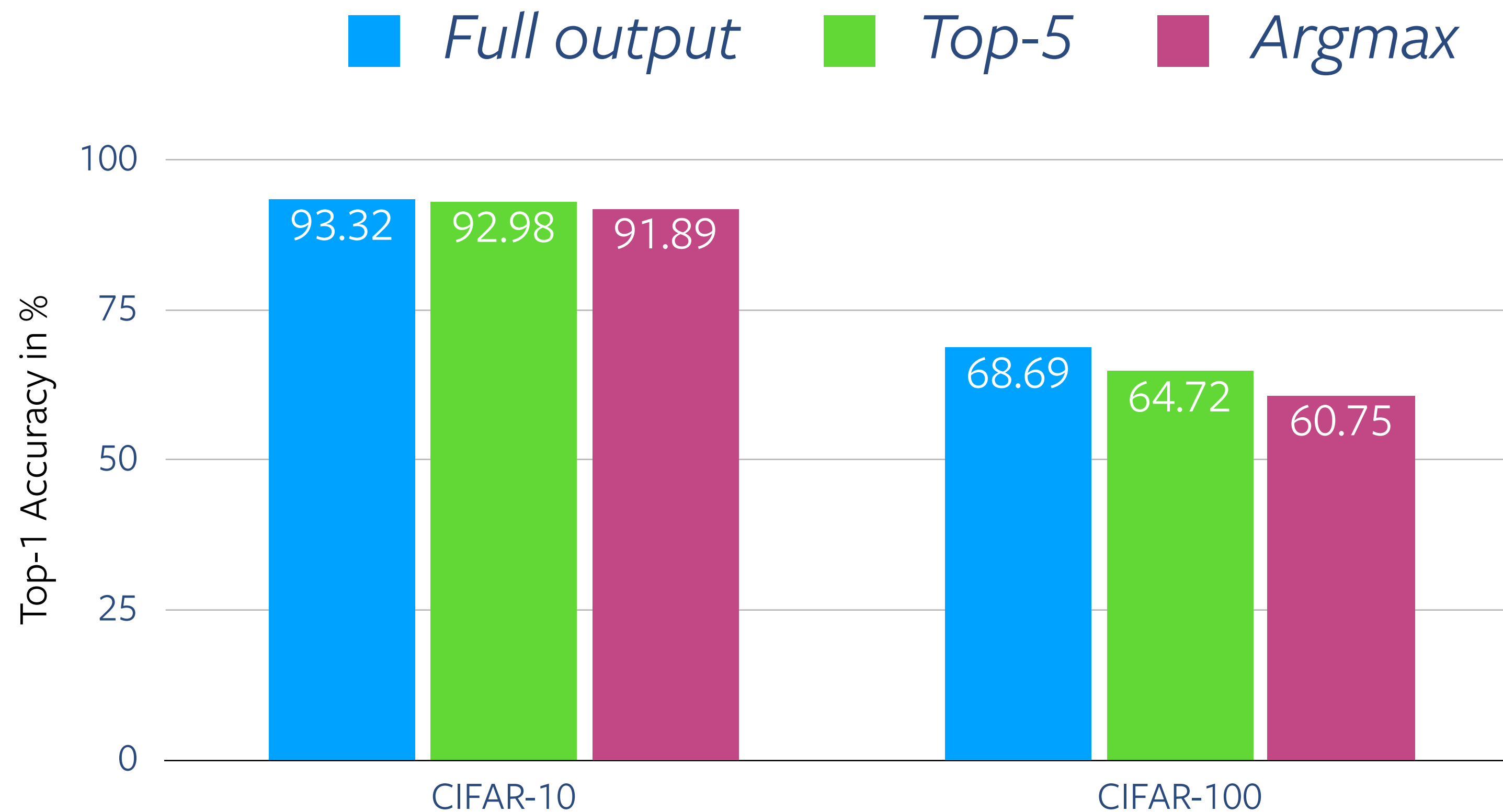
# So what accuracy do we get on ImageNet?



Our model achieves 69.0% top-1 accuracy on ImageNet-1k val.  
Without ever having seen more than that one image.



# Learning signal: even top-5 or argmax works well.

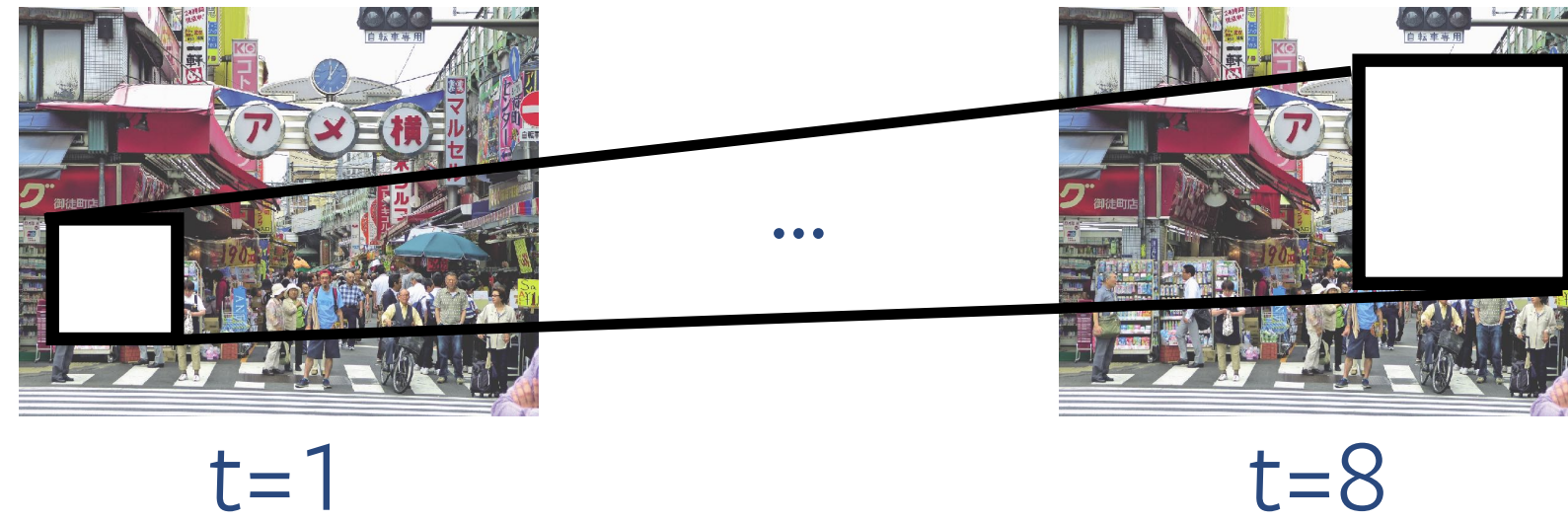


On ImageNet-1k: using argmax:  
**still 44%** top-1 Acc. (compared to 69.%)

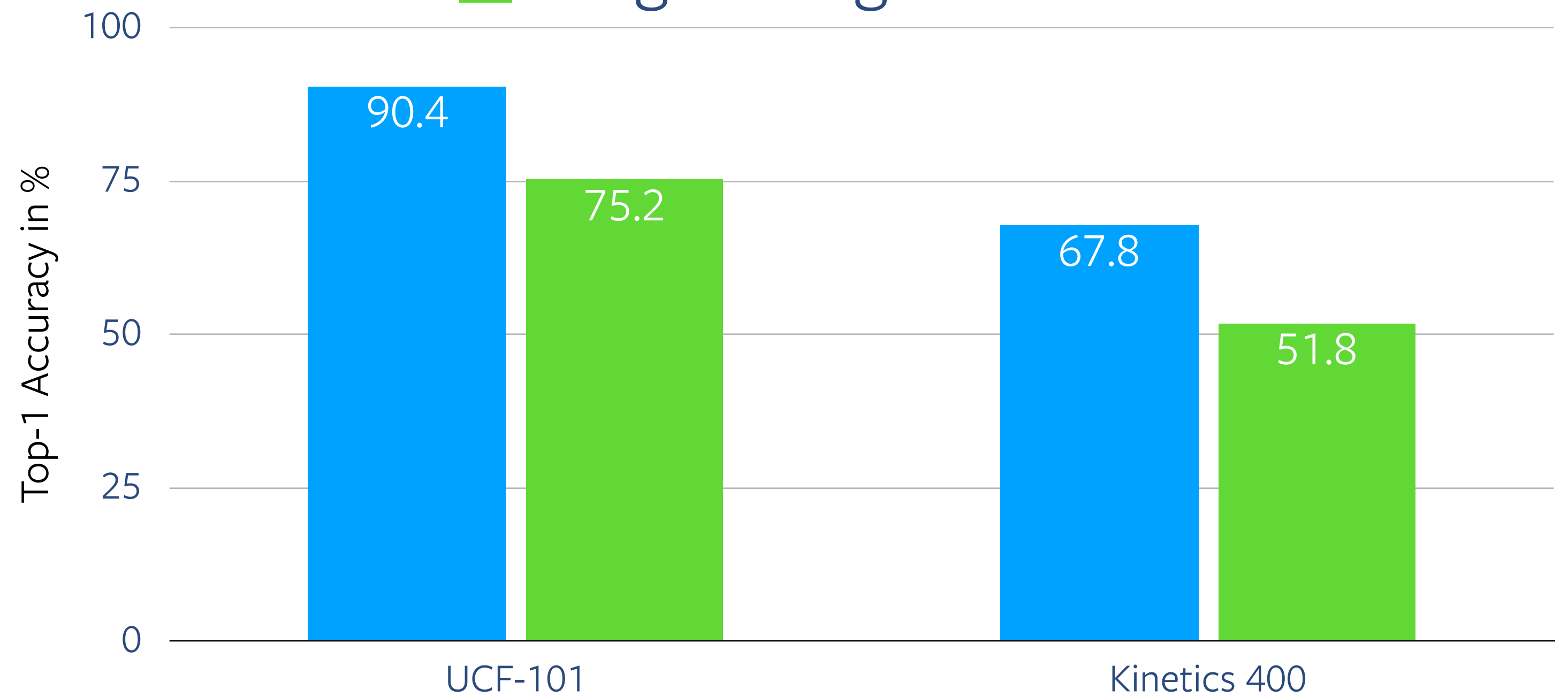
Even with only **top-5 predictions**  
(and confidence) or **hard distillation**,  
performance only slightly degrades.

⚠ API providers! (c.f. Orekondy et al.)

# Performance on video action classification benchmarks.

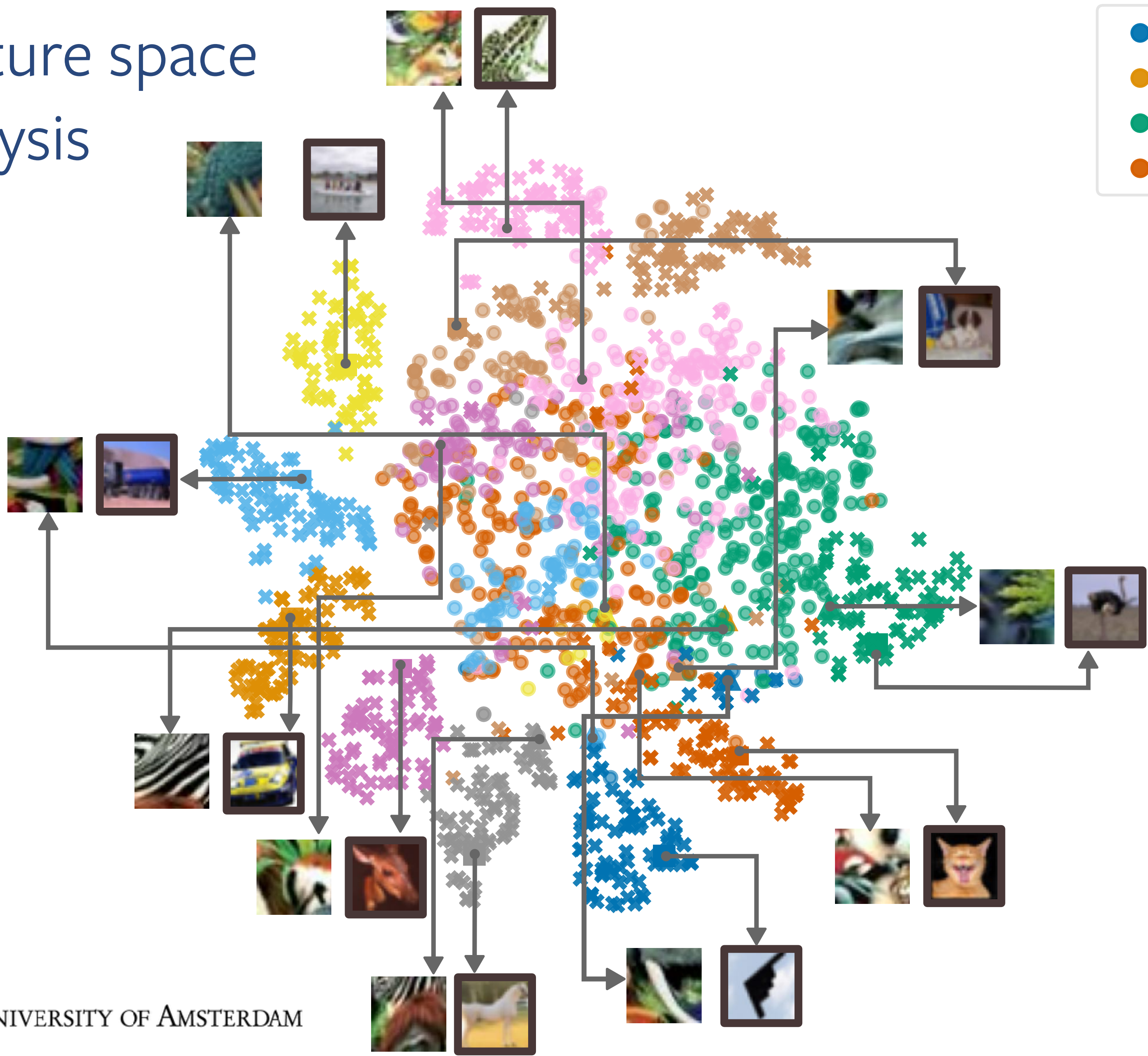


■ *Teacher (X3D-S)*  
■ *Single "Image" trained student*



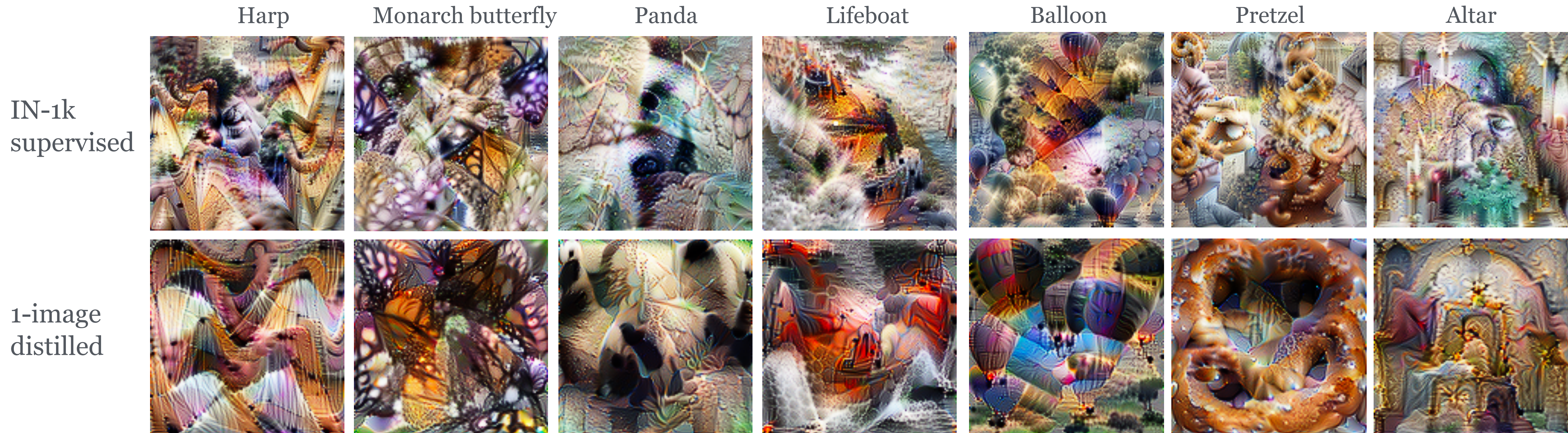
# Feature space analysis

● Airplane	● Deer	● Horse
● Automobile	● Dog	● Ship
● Bird	● Frog	● Truck
● Cat		



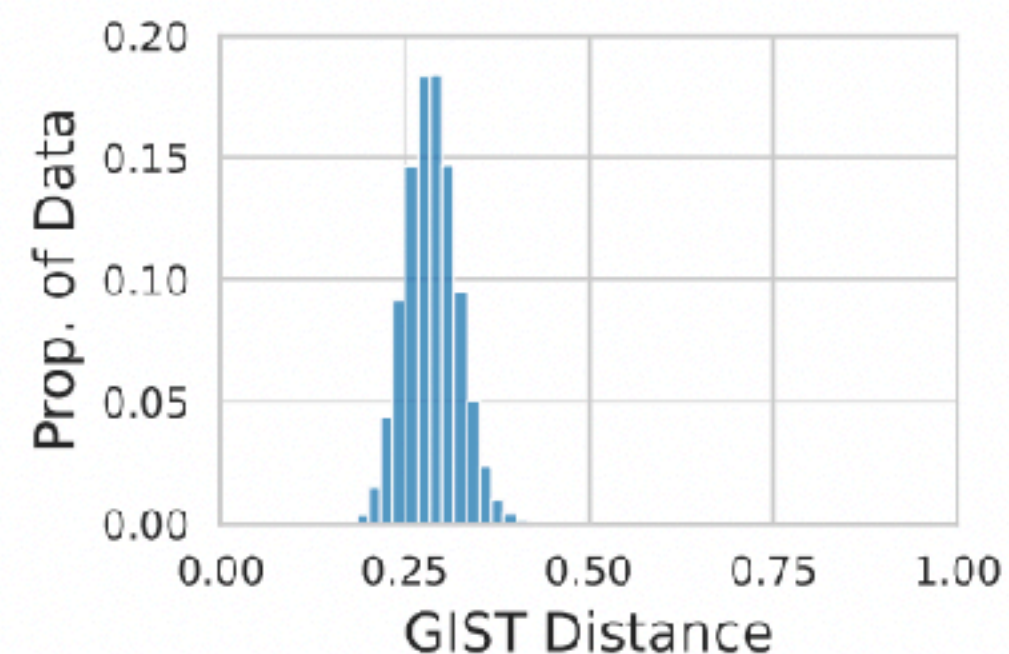
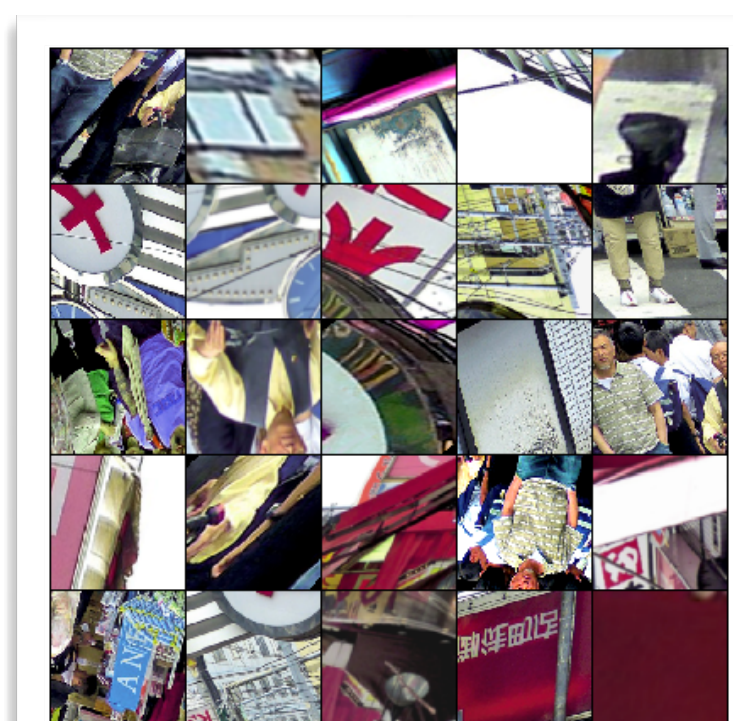
- Semantic extrapolation works
- patches (·) are "inside", real images (×) "outside"

# Neuron activation maximisation: learning pandas from the city-jungle.

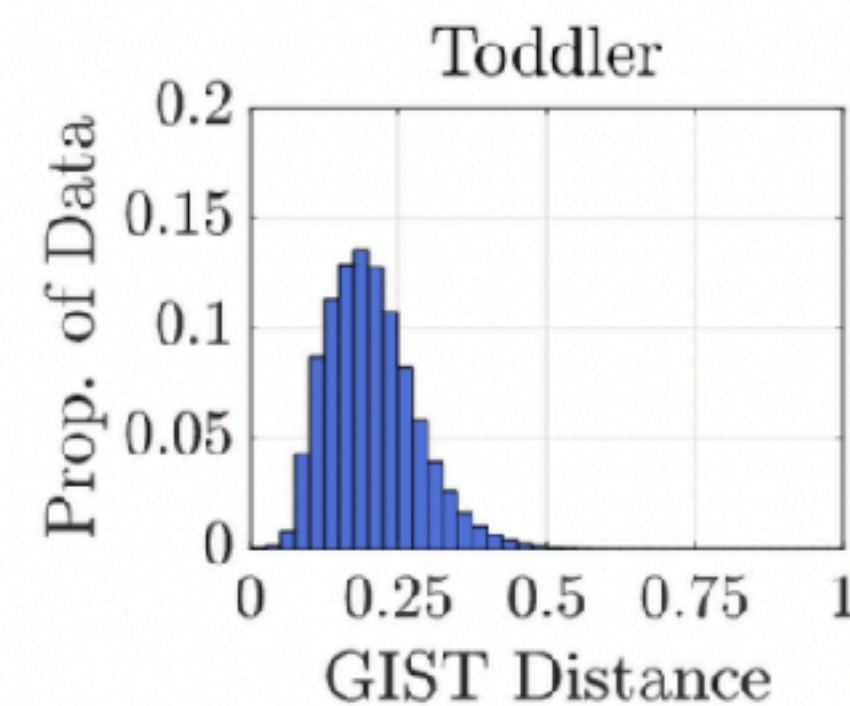
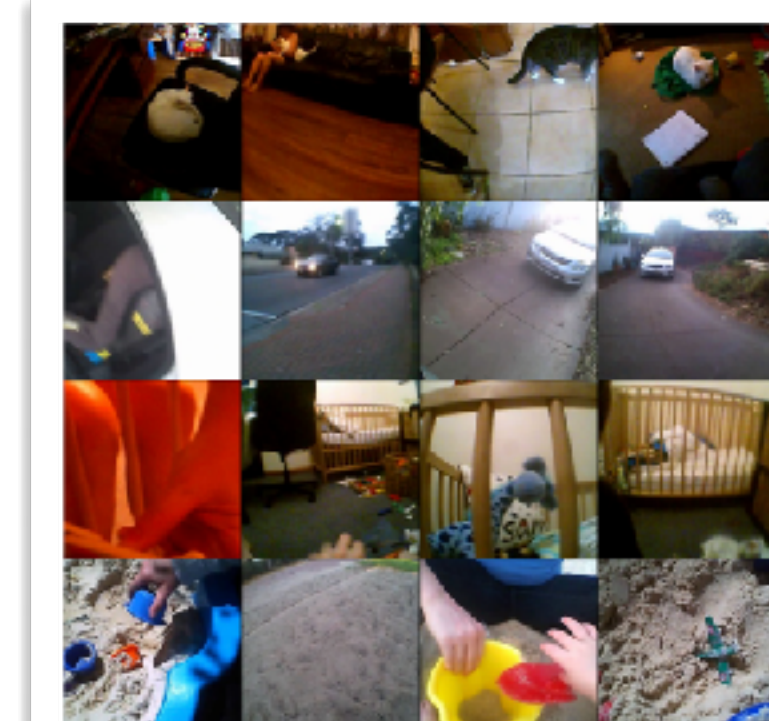




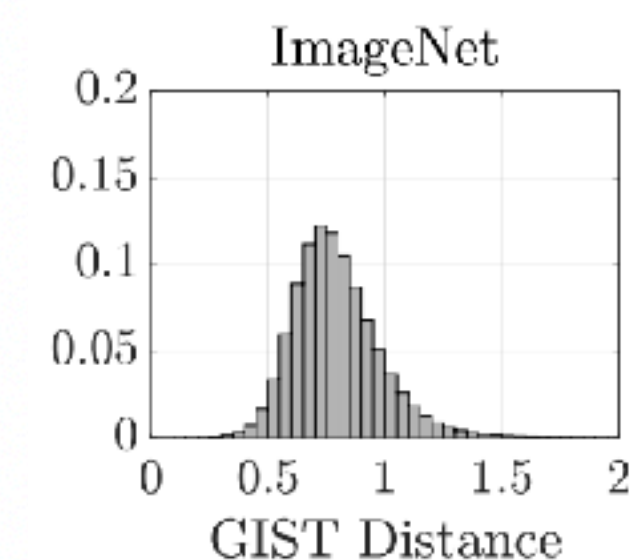
## One hypothesis



(a) **Our 1-image dataset.**



(b) **Toddler data.** Figure from [3].



Toddlers data looks much more like 1-augmented image, compared to ImageNet. Maybe this high correlation helps with learning?

# Team for the works presented

## *Time Tuning*



Mohammadreza Salehi



Efstratios Gavves



Cees G. M. Snoek



Yuki M. Asano

## *WTour Dora*



Shashanka Venkataramanan



Mamshad N Rizve



Joao Carreira



Yuki M. Asano\*

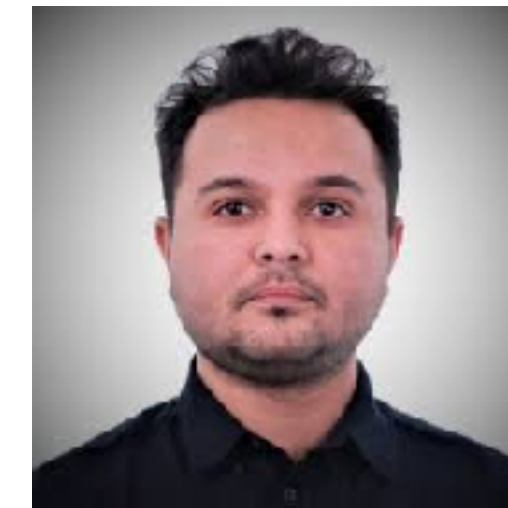


Yannis Avrithis\*

## *1-image distill*



Yuki M. Asano<sup>1</sup>



Aaqib Saeed<sup>1</sup>

Salehi, Gavves, Snoek, Asano. *Time does tell: self-supervised time-tuning of dense image representations*. ICCV 2023

Venkataramanan, Rizve, Carreira, Avrithis\*, Asano\*. *Is ImageNet worth 1 video? Learning strong image encoders from 1 long unlabelled video*. ICLR 2024 [oral]

Asano & Saeed. *The Augmented Image Prior: Distilling 1000 Classes by Extrapolating from a Single Image*. ICLR 2023.

<sup>1</sup>: co-first authors; \*: co-last authors