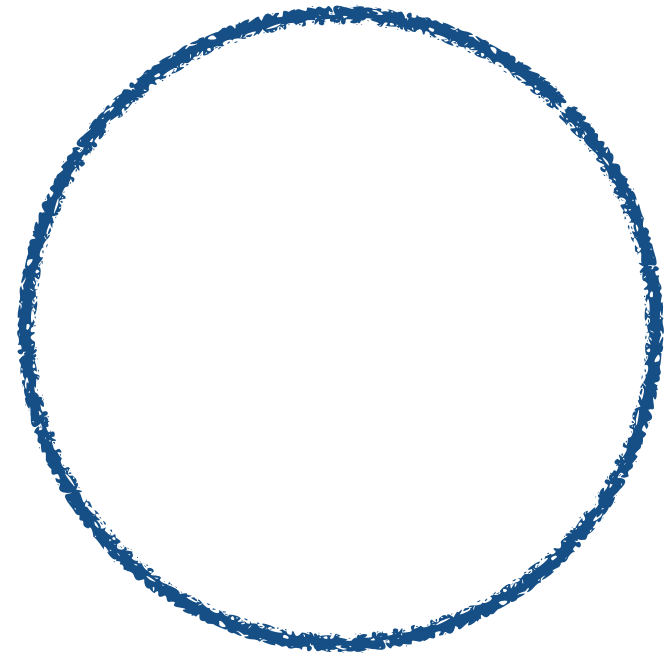


Modern Self-supervised Image Representation Learning from Videos

YUKI ASANO
NEURIPS 2023

Why do we want Self-supervised Learning in the age of CLIP *et al*?

Massive scale



sup. \ll weak sup. \ll raw

Cost of (re)labelling



Problems of labels



Fundamentals



Especially videos open exiting new directions



Visual development for AI



"Get" physics

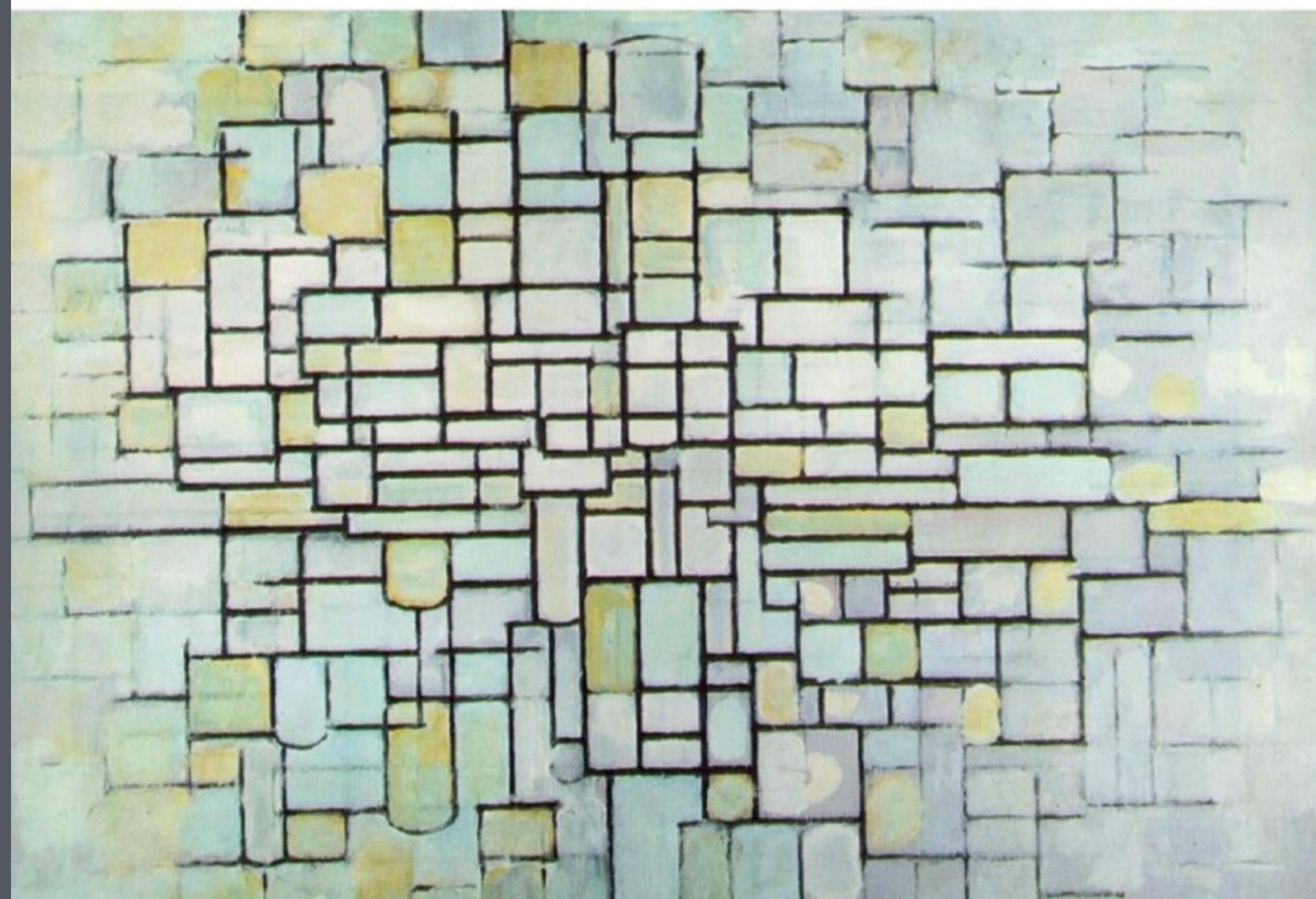


Embodied AI

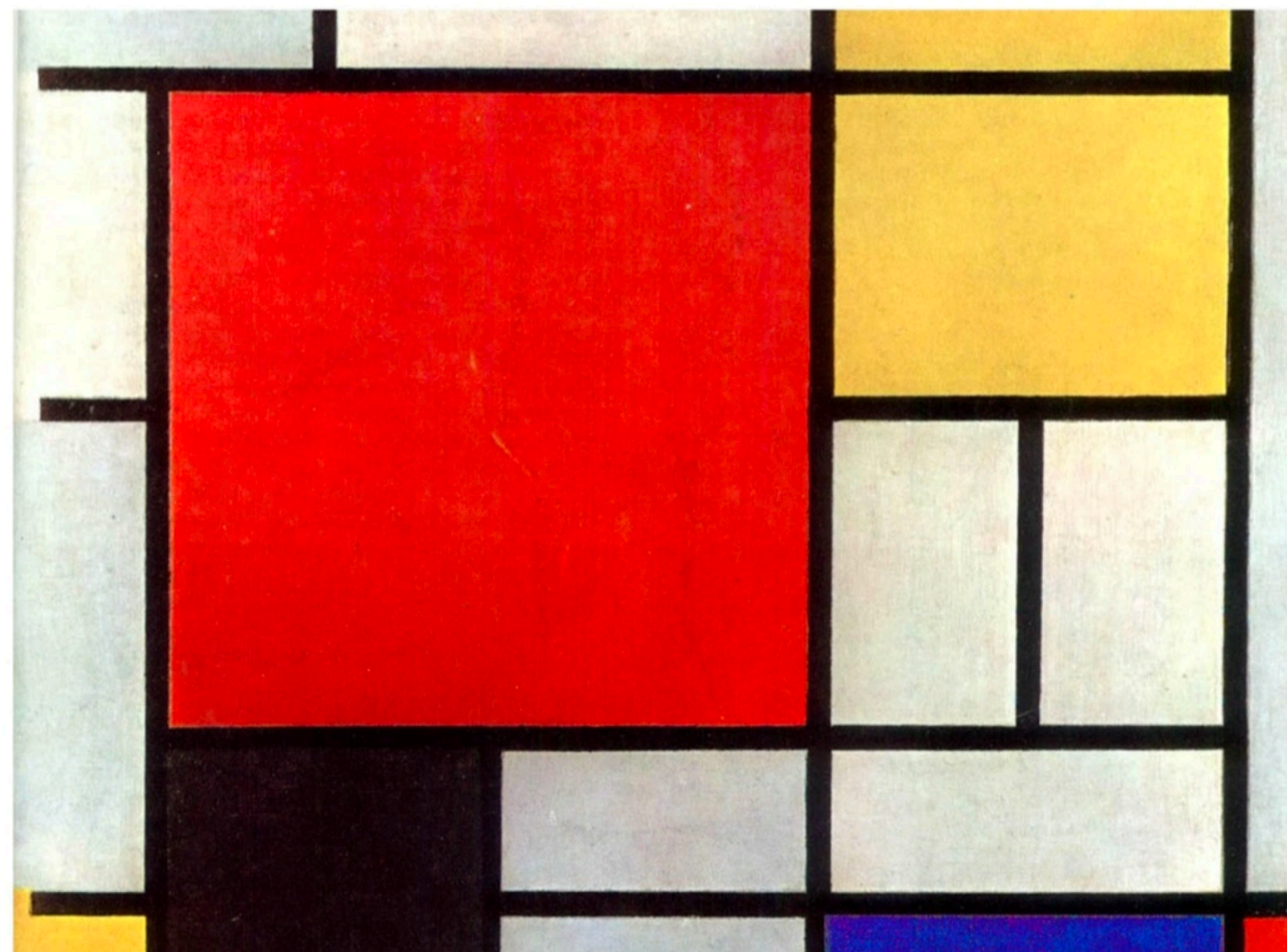
Bonus: *insane* scale:



LEARNING
IMAGE
ENCODERS
FROM TIME

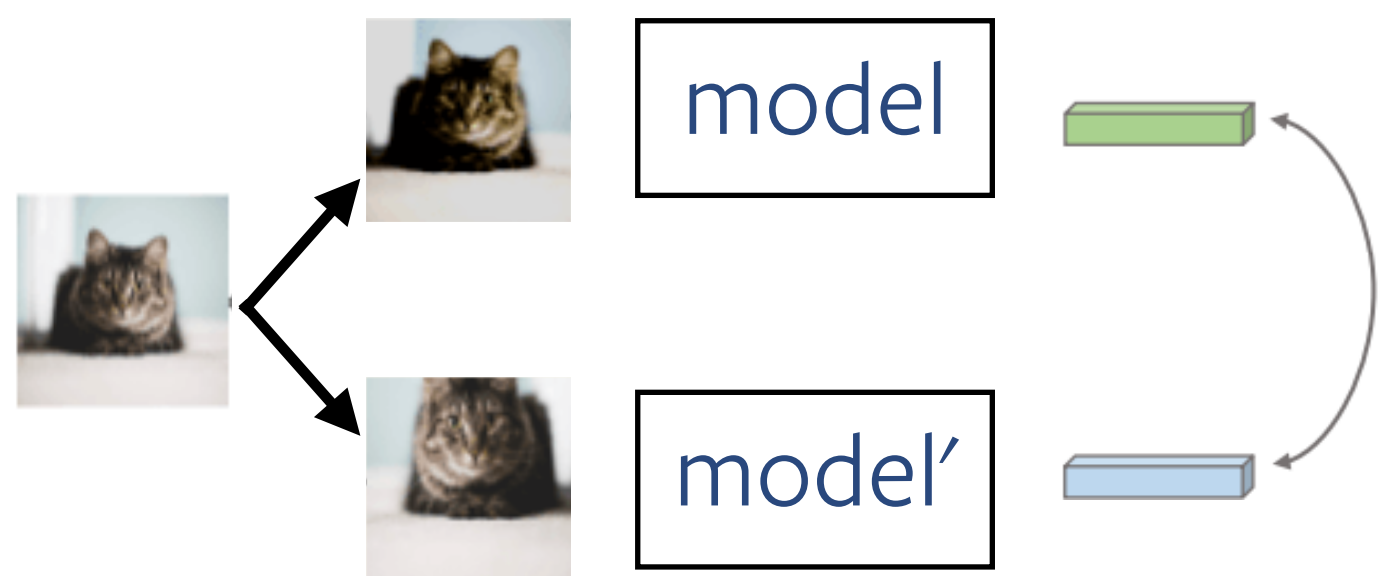


Piet Mondrian, paintings 1910, 1911, 1913, 1930



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

Images: SimCLR, MoCo, SwaAV et al.

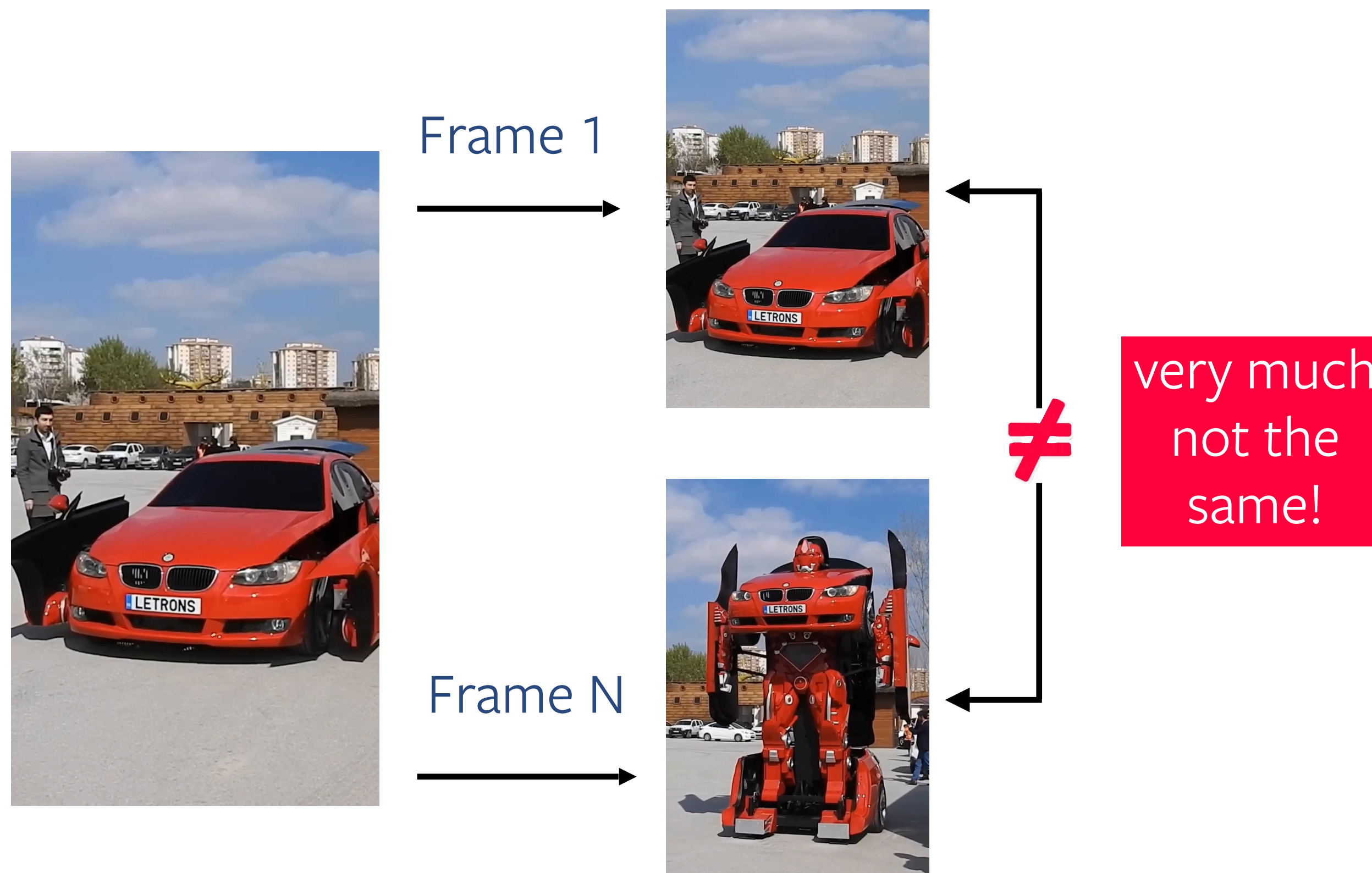


key principle: view-invariance

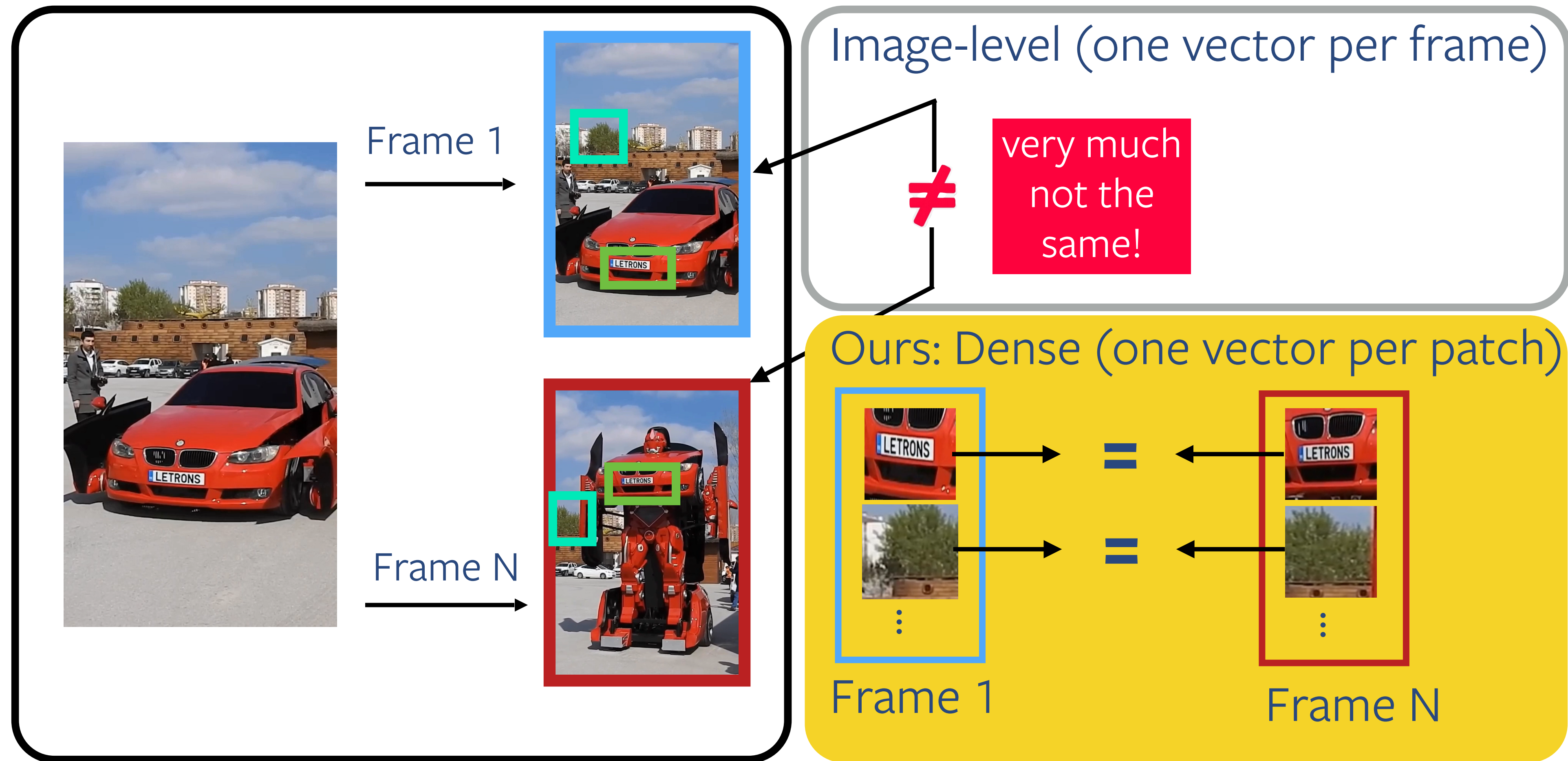
Might be ok for videos like this:



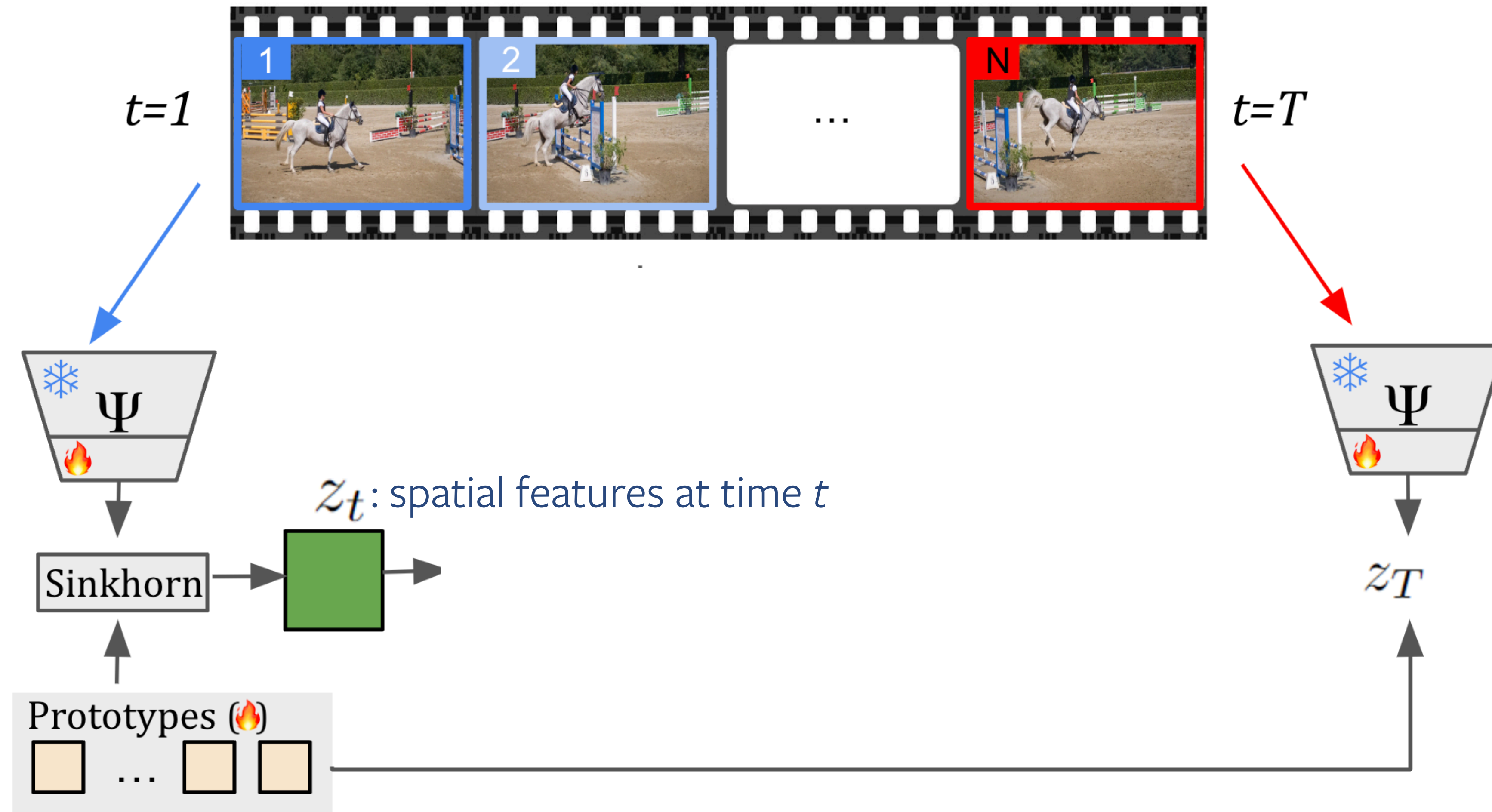
But does this generally make sense?



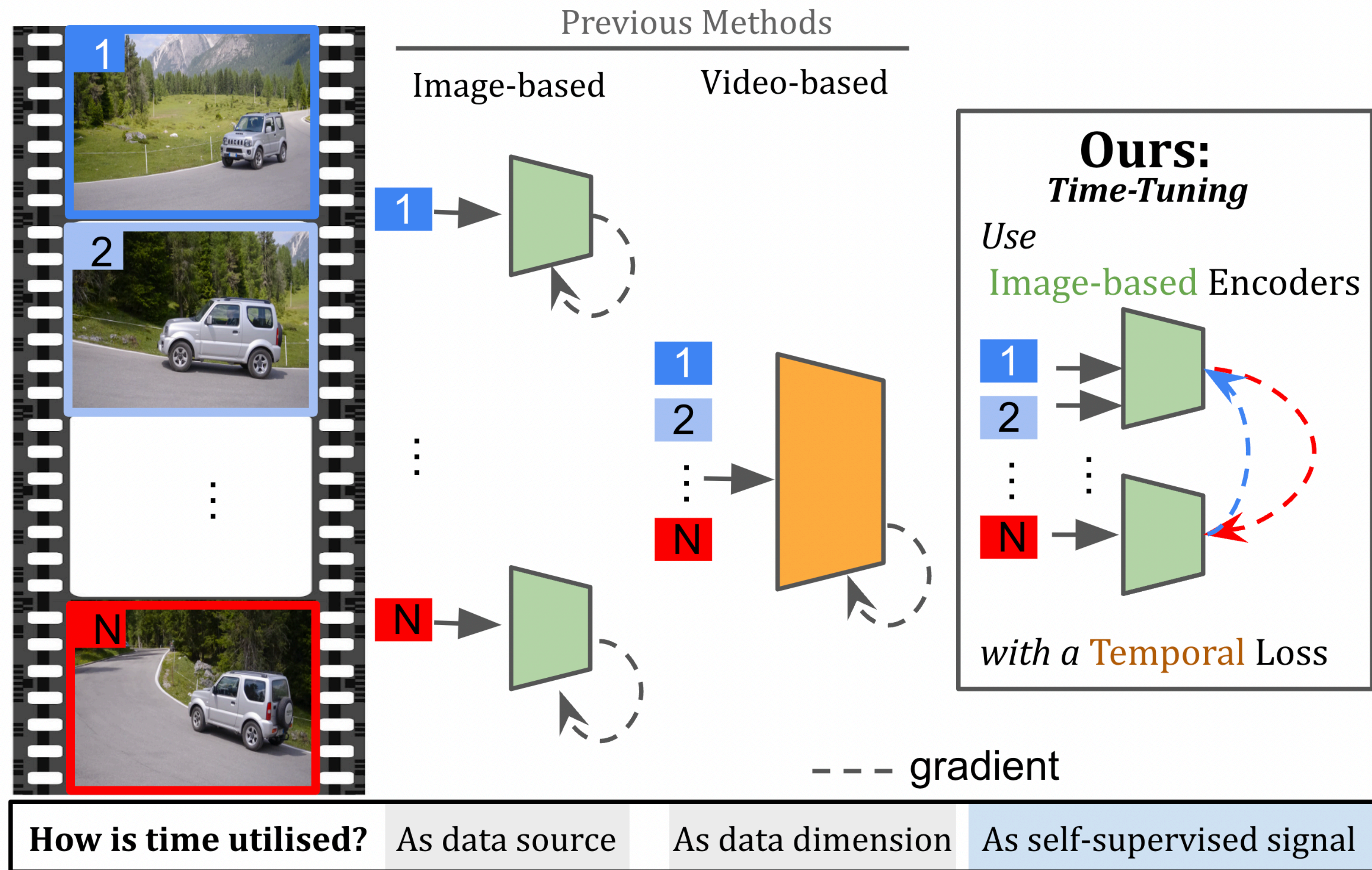
Solution is obvious



We *model* a video by tracking image patches, and aligning their clustered features

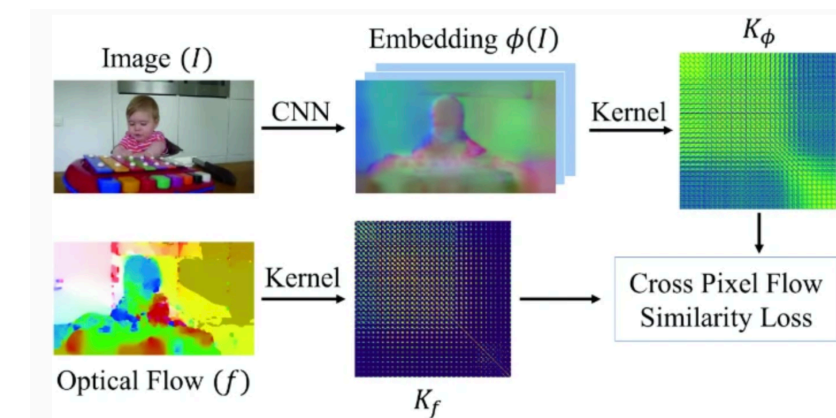
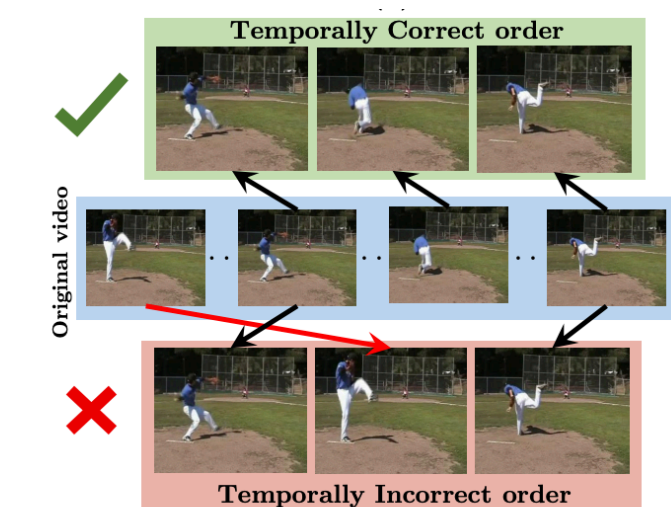
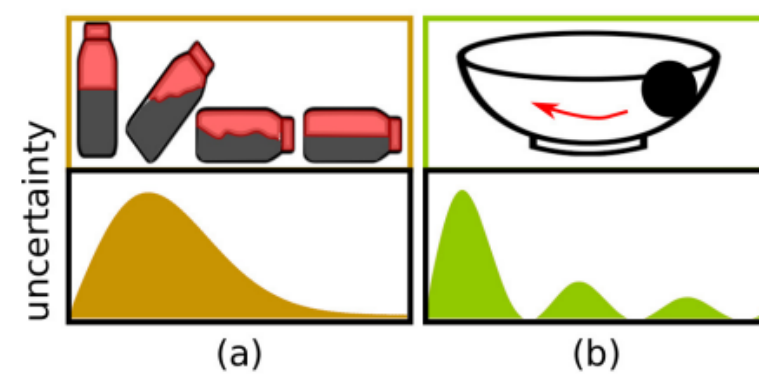
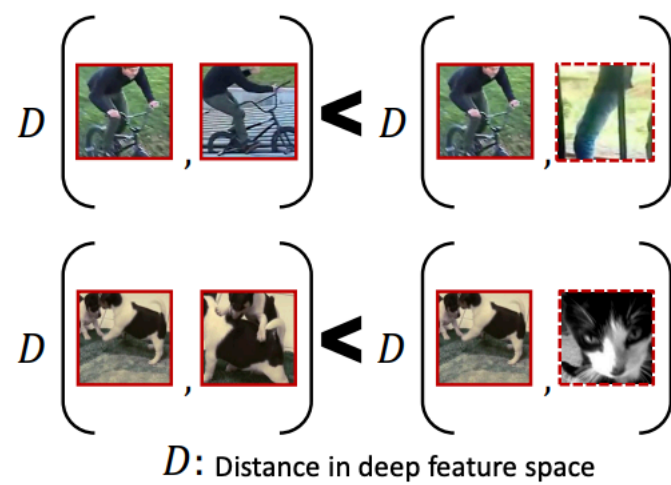
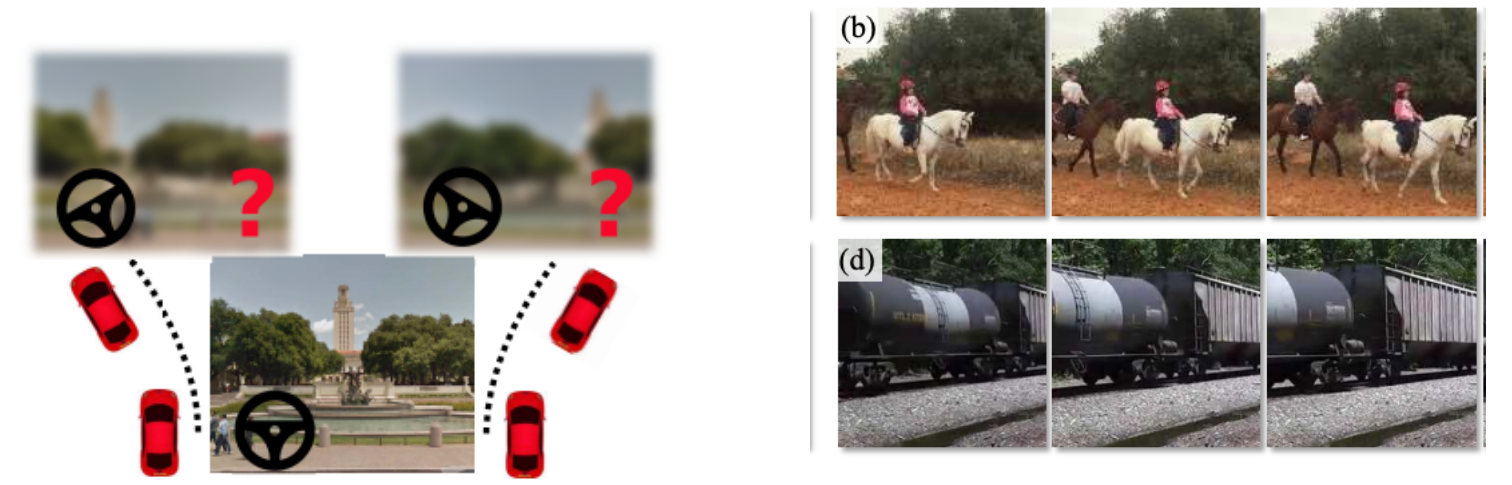


Using videos to learn self-supervised image encoders

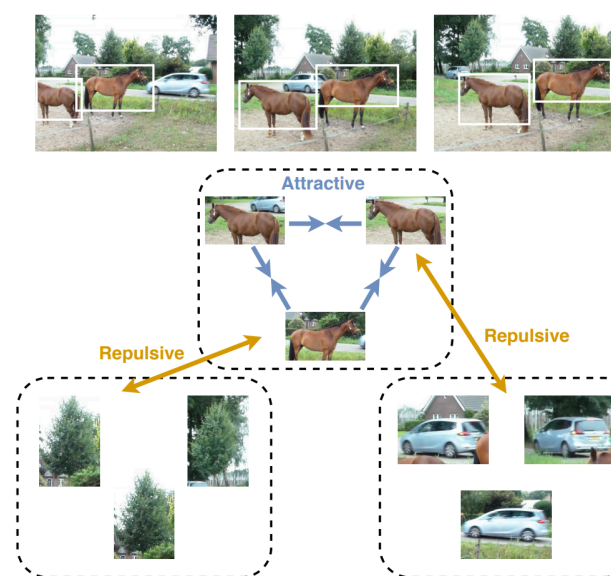
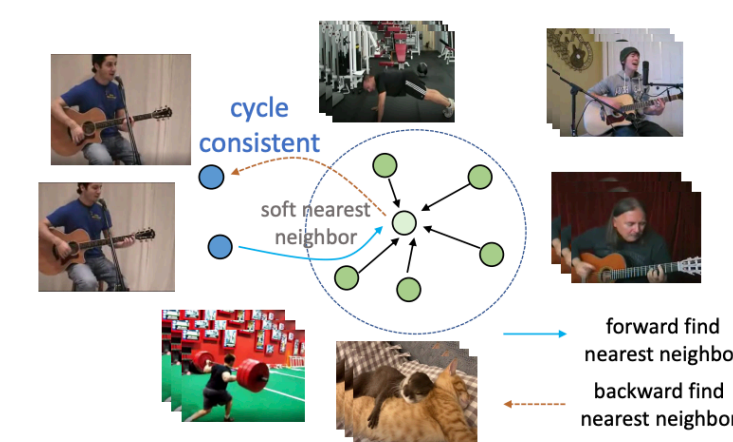


This field has a rich history. And now it's time to get back to it

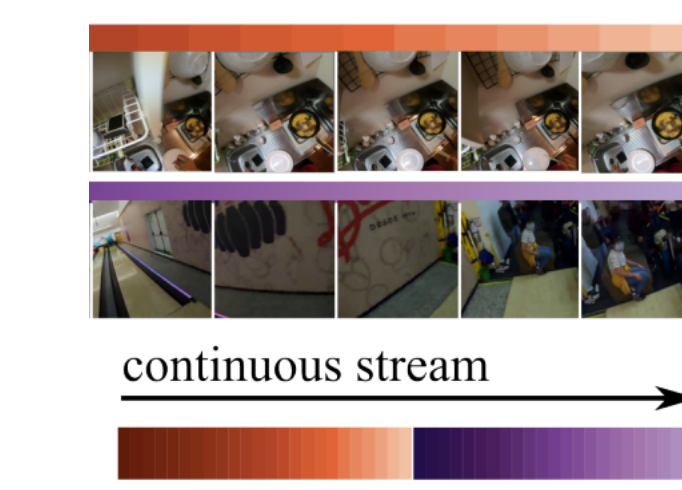
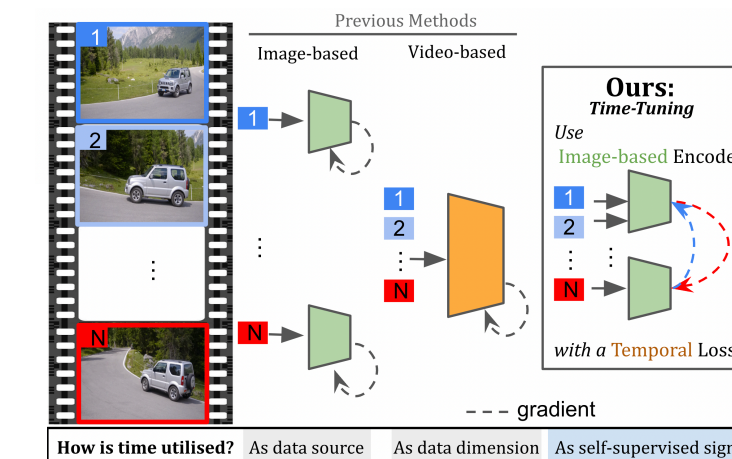
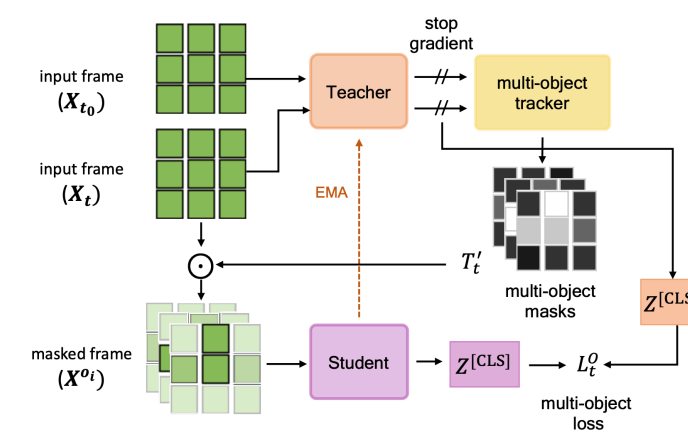
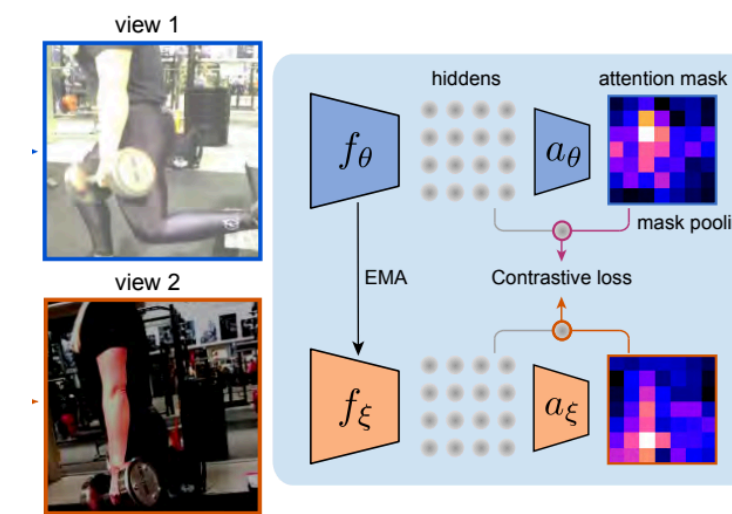
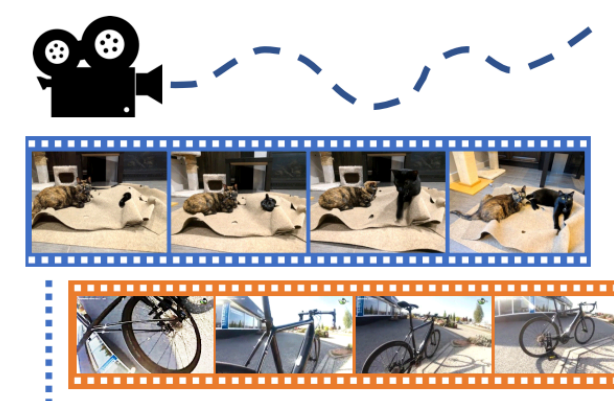
2015-2019



2020 onwards getting competitive to ImageNet



Video Noise Contrastive Estimation

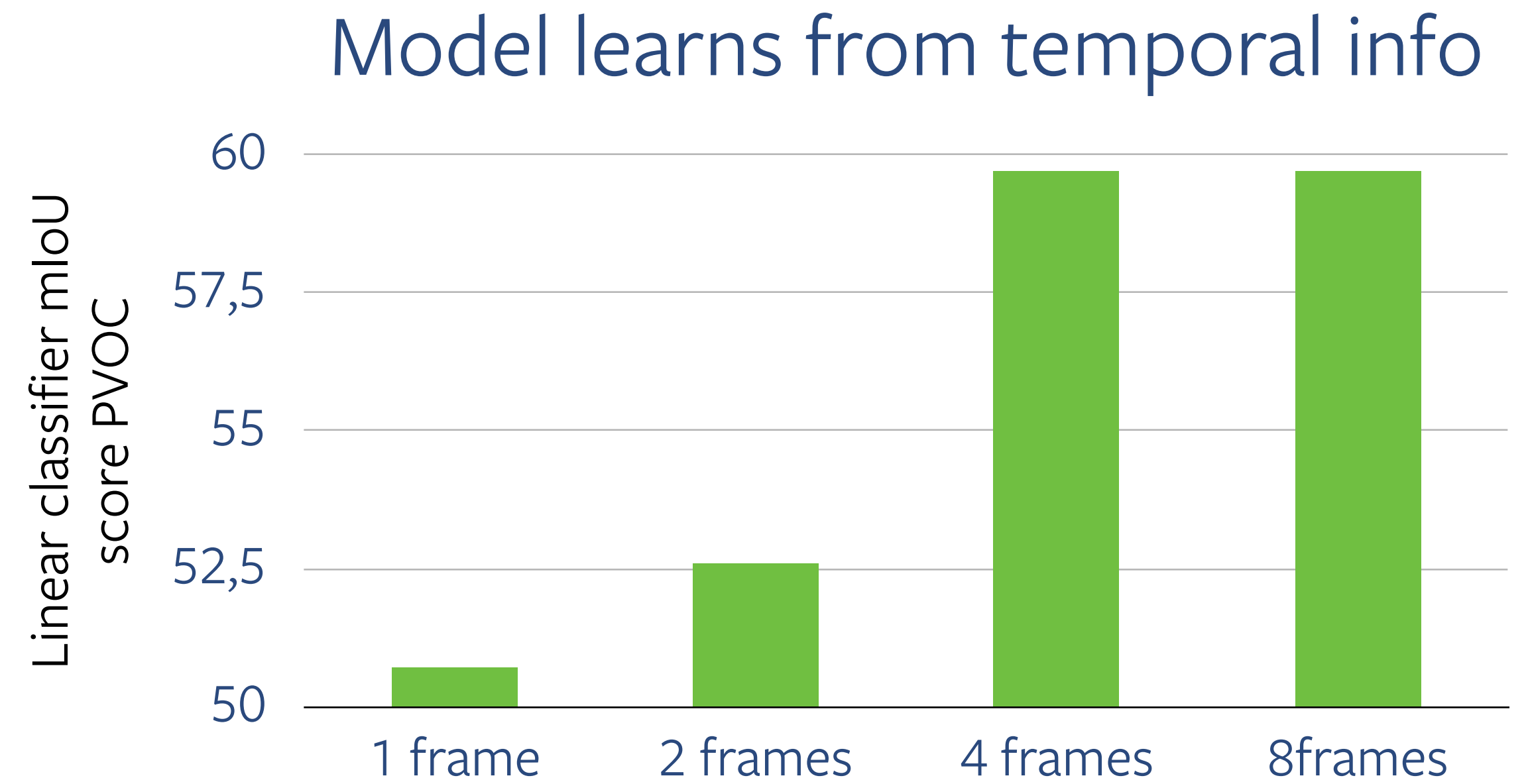
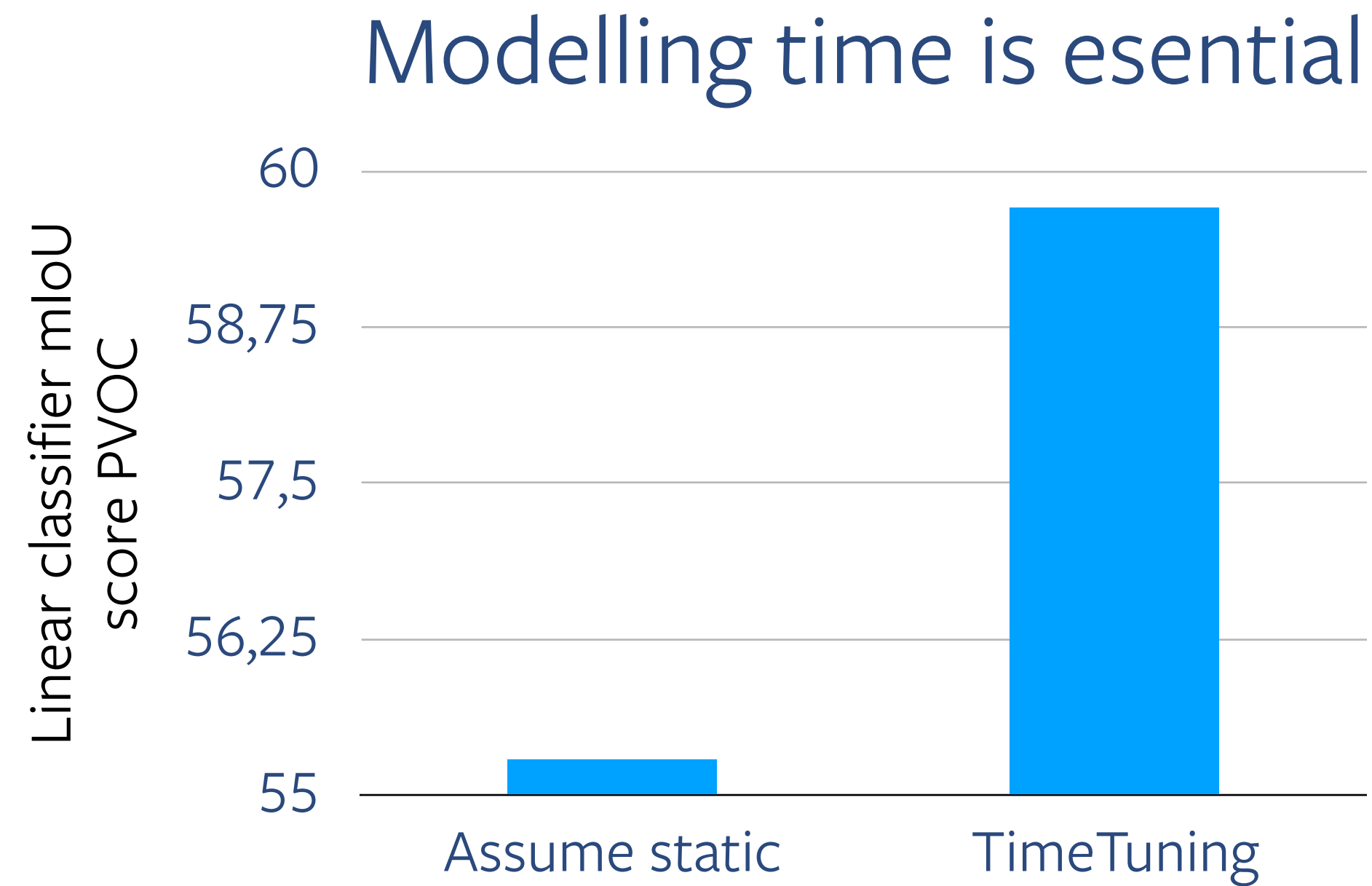


Some references:

2002 Wiskott, Sejnowski. Slow Feature Analysis: Unsupervised Learning of Invariances
 2015 Agrawal, Carreira, Malik. Learning to See by Moving: predict egomotion from frames
 2015 Wang, Gupta. Unsupervised Learning of Visual Representations using Videos
 2015 Goroshin, Bruna, Eigen, LeCun. Unsupervised feature learning from temporal data
 2015 Ramanathan, Tang, Mori, Fei-Fei. Learning Temporal Embeddings for Complex Video Analysis
 2016 Gao, Jayaraman, Graumann. Object-Centric Representation Learning from Unlabeled Videos
 2017 Wang, Kaiming, Gupta. Transitive Invariance for Self-supervised Visual Representation Learning
 2018 Wei, Lim, Zisserman, Freeman. Learning and Using the Arrow of Time
 2019 Jayaraman, Ebert, Efron, Levine. Time-Agnostic Prediction: Predicting Predictable Video Frames
 2019 Mahendran, Thewlis, Vedaldi. focus on motion: cross-pixel flow

2020 Tschannen, Djolonga, Ritter, Mahendran, Zhai, Houlsby, Gelly, Lucic. Self-Supervised Learning of Video-Induced Visual Invariances.
 2021 Wu, Wang. Contrastive Learning of Image Representations with Cross-Video Cycle-Consistency
 2020 Gordon, Ehsani, Fox, Farhadi. Watching the World Go By: Representation Learning from Unlabeled Videos
 2023 Parthasarathy, Eslami, Carreira, Henaff. Self-supervised video pretraining yieldshuman-aligned visual representations
 2023 Salehi, Gavves, Snoek, Asano. Time does tell: self-supervised time-tuning of dense image representations
 2023 Venkataramanan, Rizve, Carreira, Avrithis*, Asano*. Is ImageNet worth 1 video? Learning strong image encoders from 1 long unlabelled video.
 2023 Carreira et al. Learning from One Continuous Video Stream

Ablations demonstrate using time helps learn better features



Unsupervised Semantic Segmentation on videos

[simply running k-means on a couple of videos' spatial features, $k=10$]

DINO



✗ part-centric maps

STEGO



✗ noisy maps

Ours



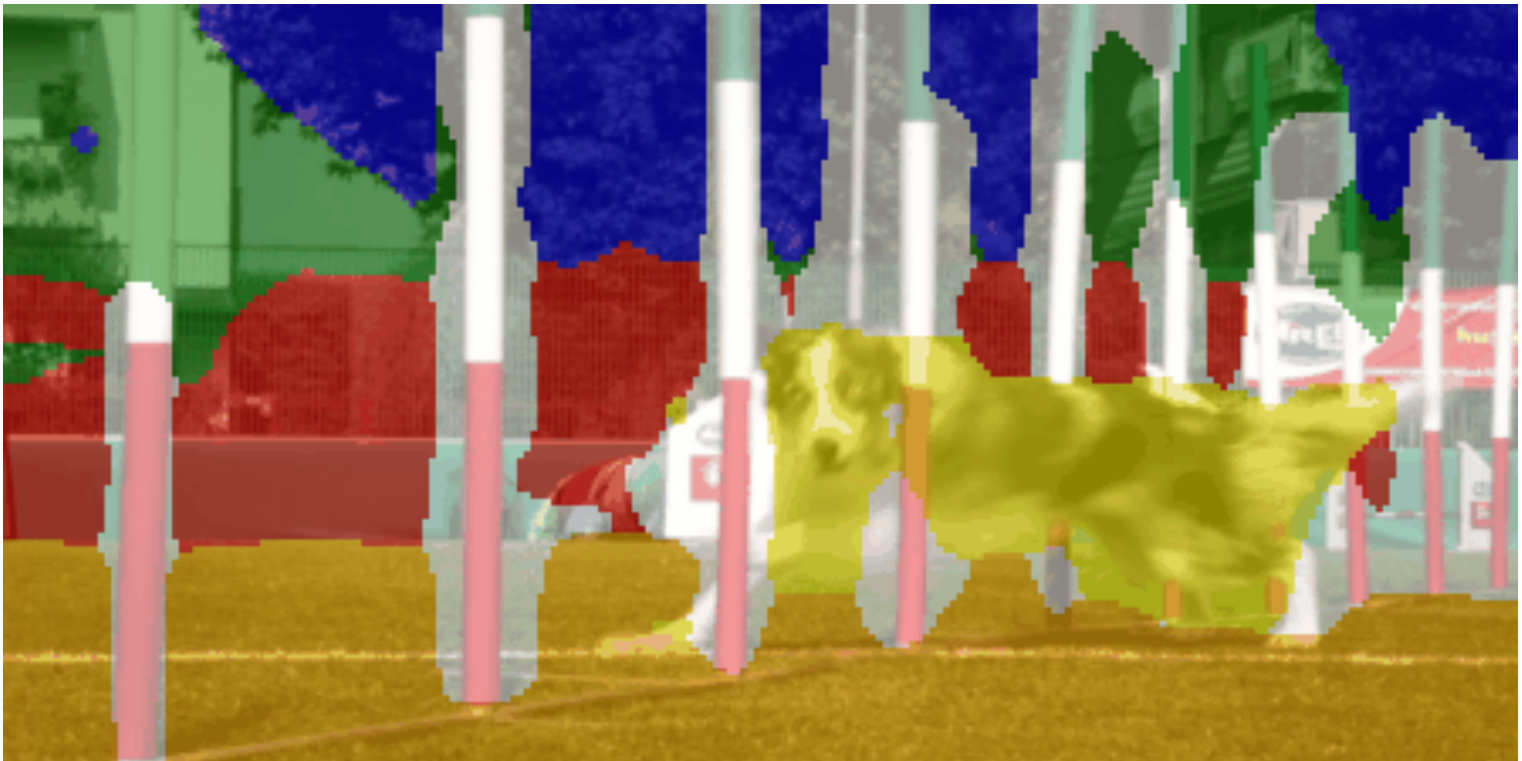
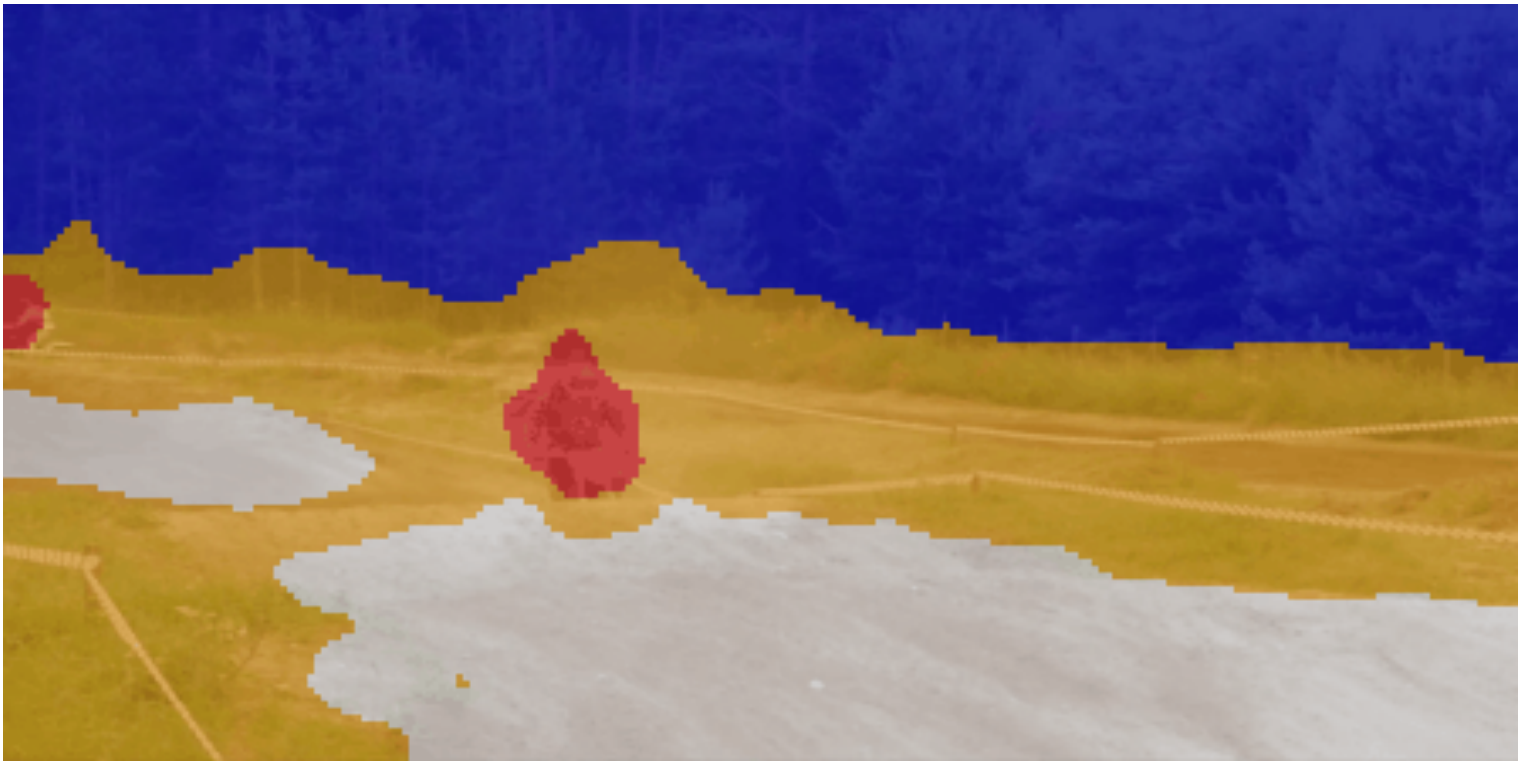
✓ crisp semantic maps

Unsupervised Semantic Segmentation on videos

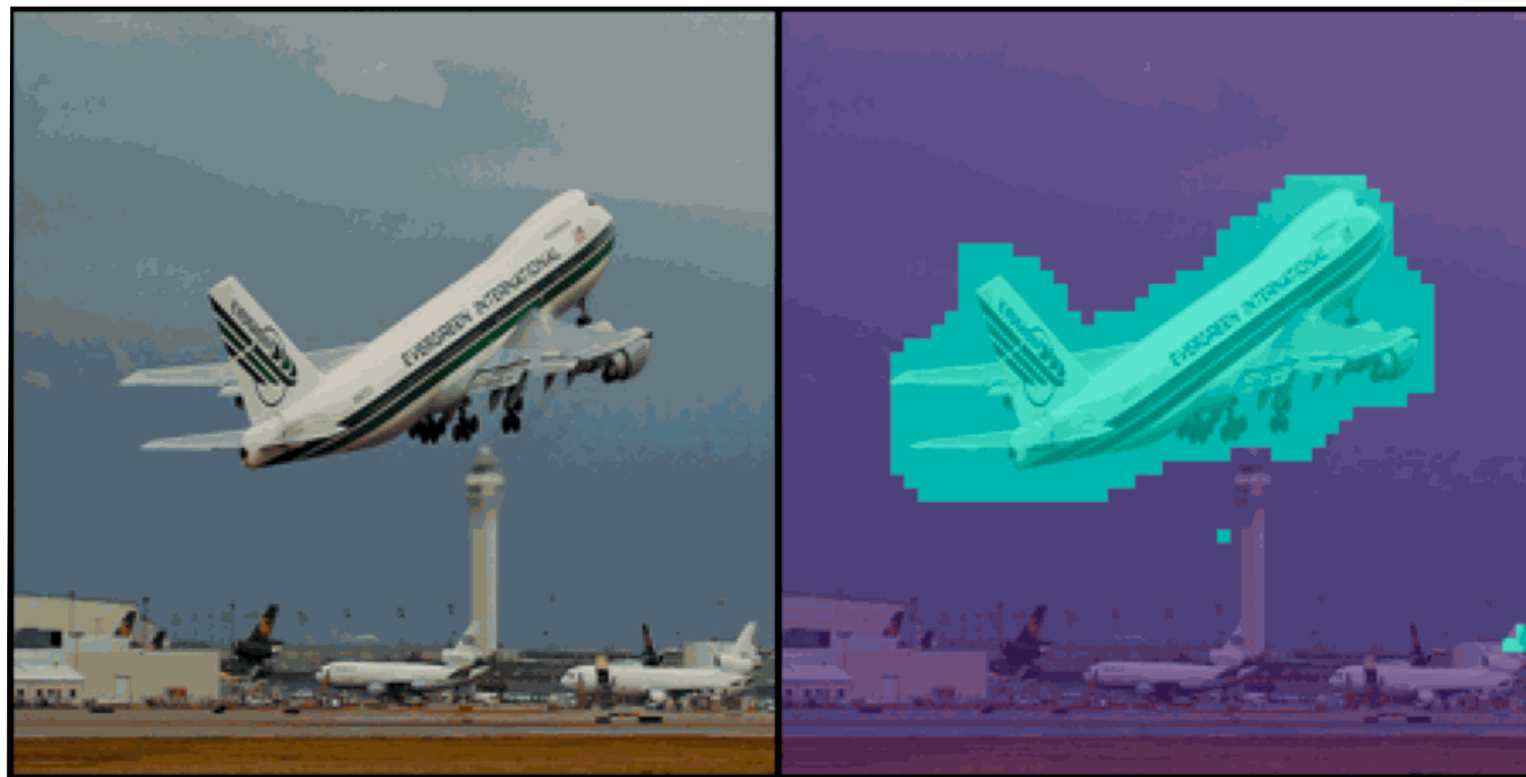
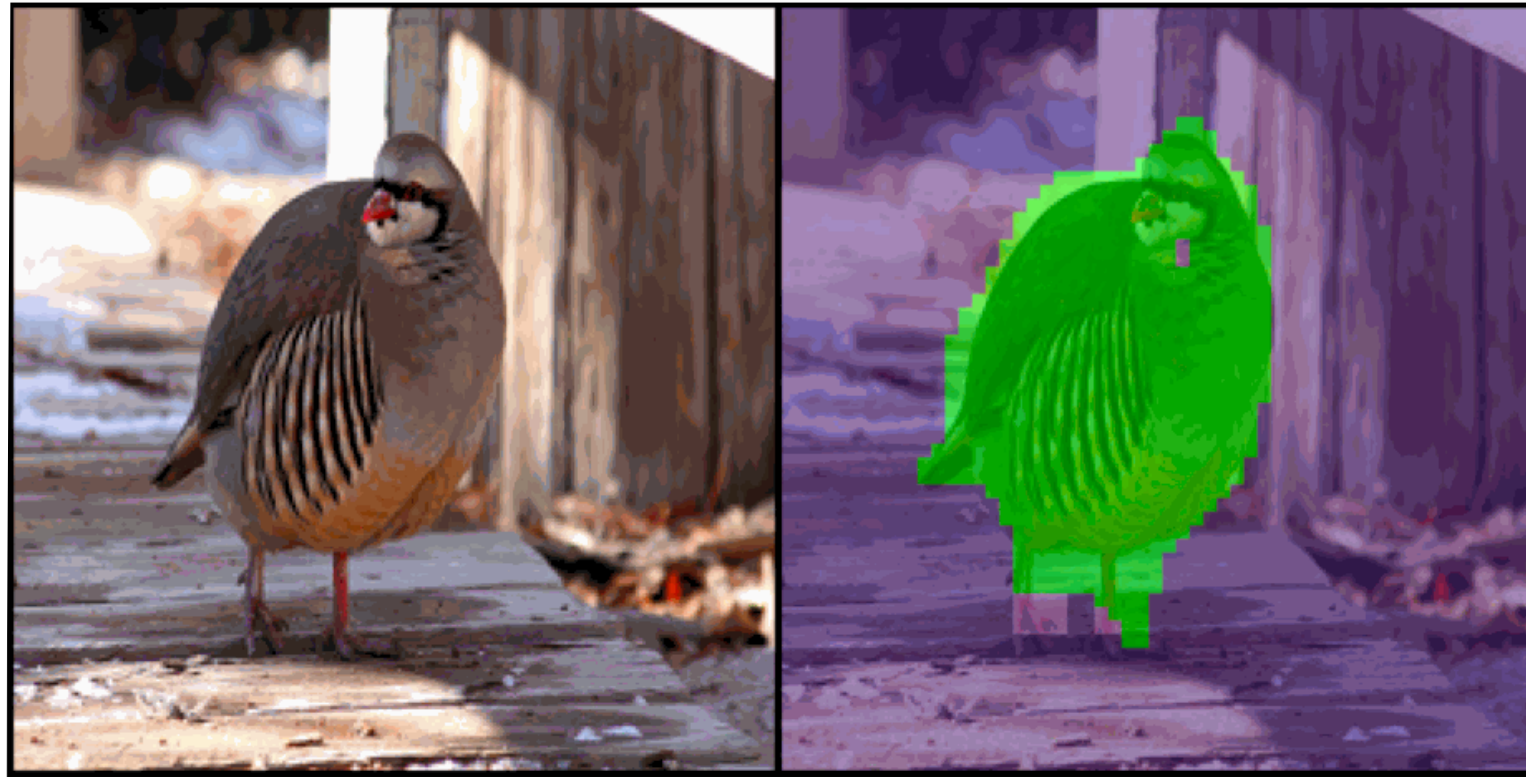
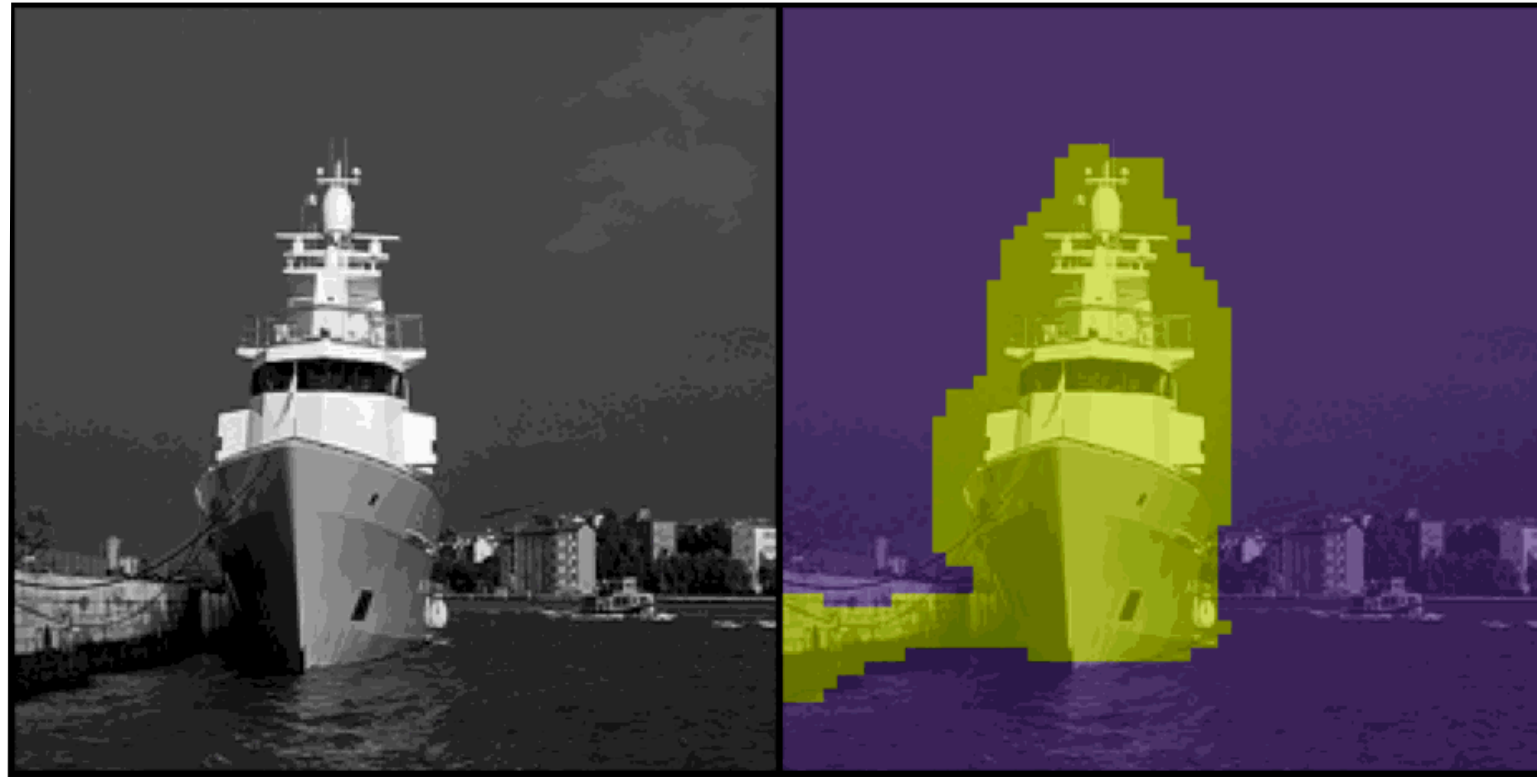
[here: running k-means on the whole video's spatial features, k=5]



More examples



Also good performance on images, despite having been tuned on videos.



TimeTuning:

DINO as init & use temporal info of videos.

How powerful is time
without image-pretraining?



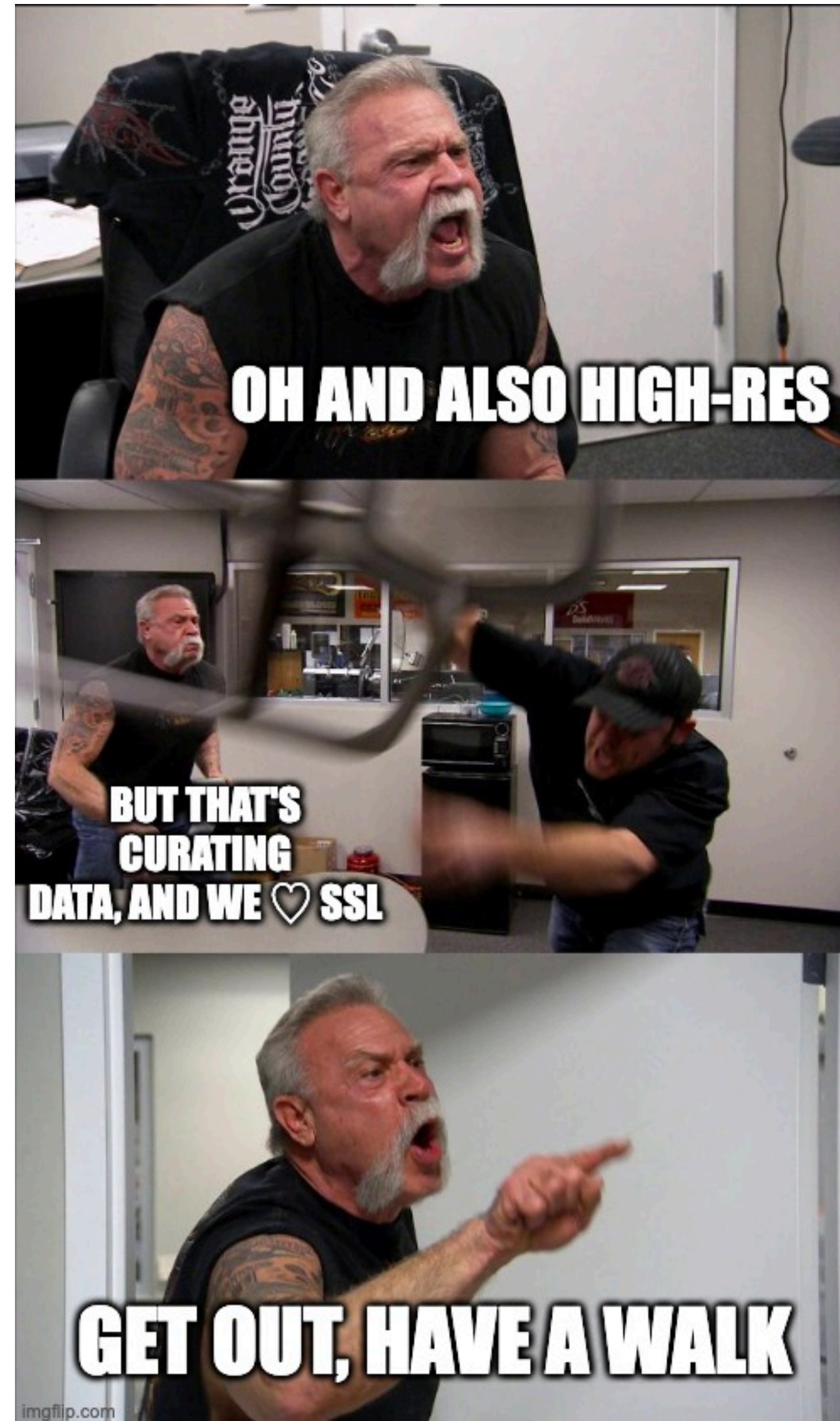
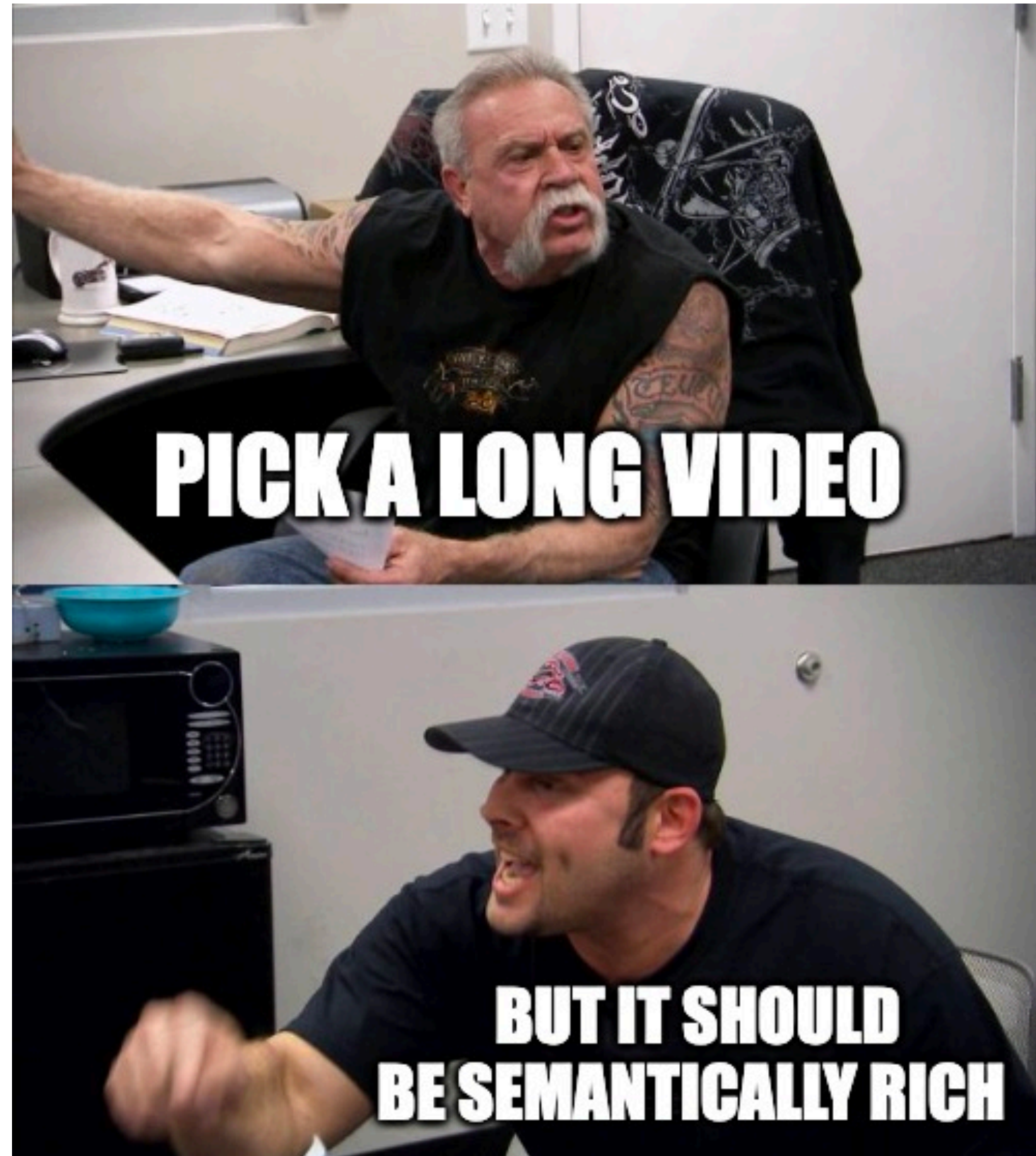
Study the extreme:
try to learn from a
single video,
from scratch.



IS ALL WE NEED ONE
LONG VIDEO WITH
MANY DETAILS?

Vermeer, The Milkmaid 1660

Us figuring out which video to use



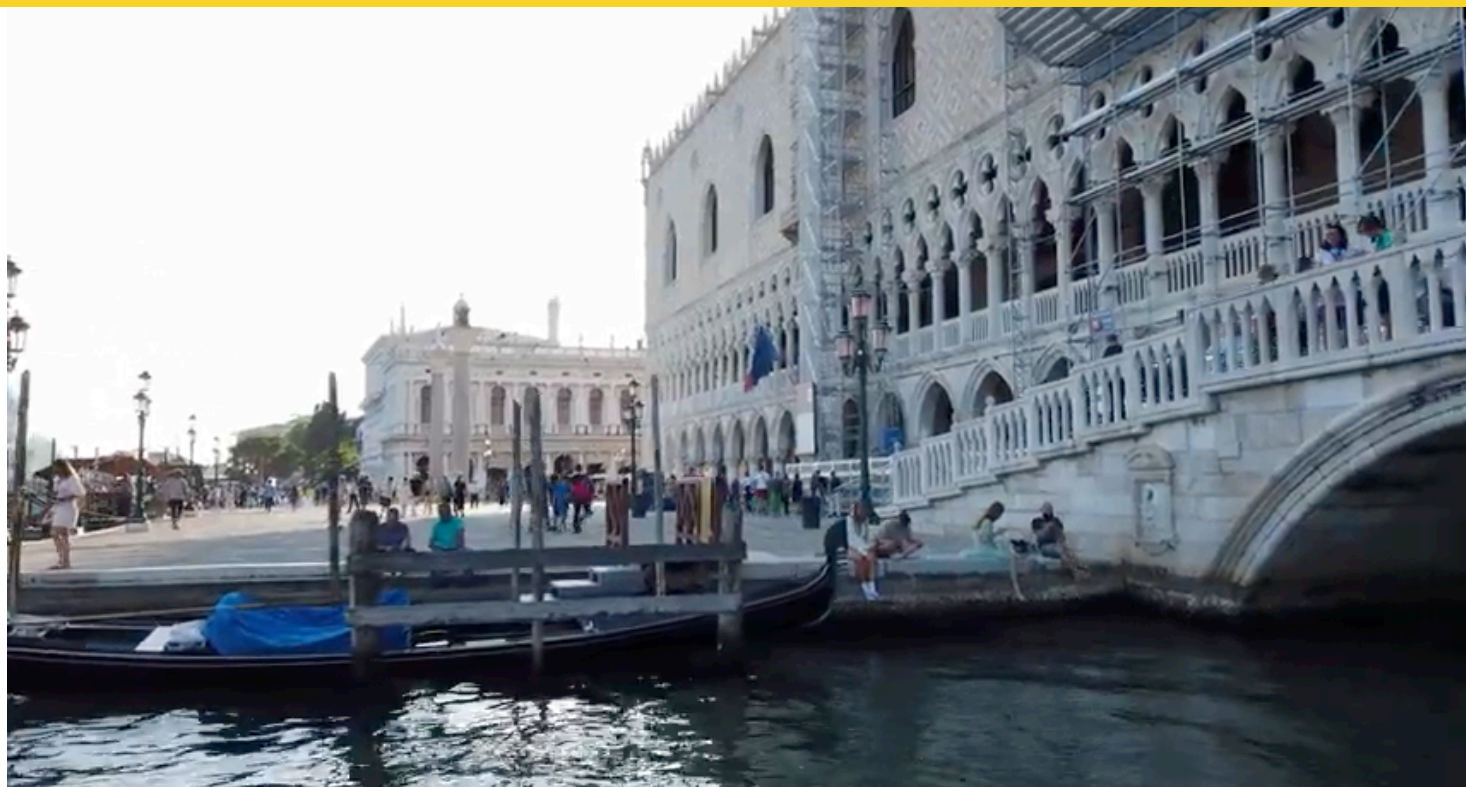
- ✓ Long
- ✓ High-res, smooth
- ✓ Semantically rich
- ✓ Scalable (we ♥ SSL)



Walking Tours



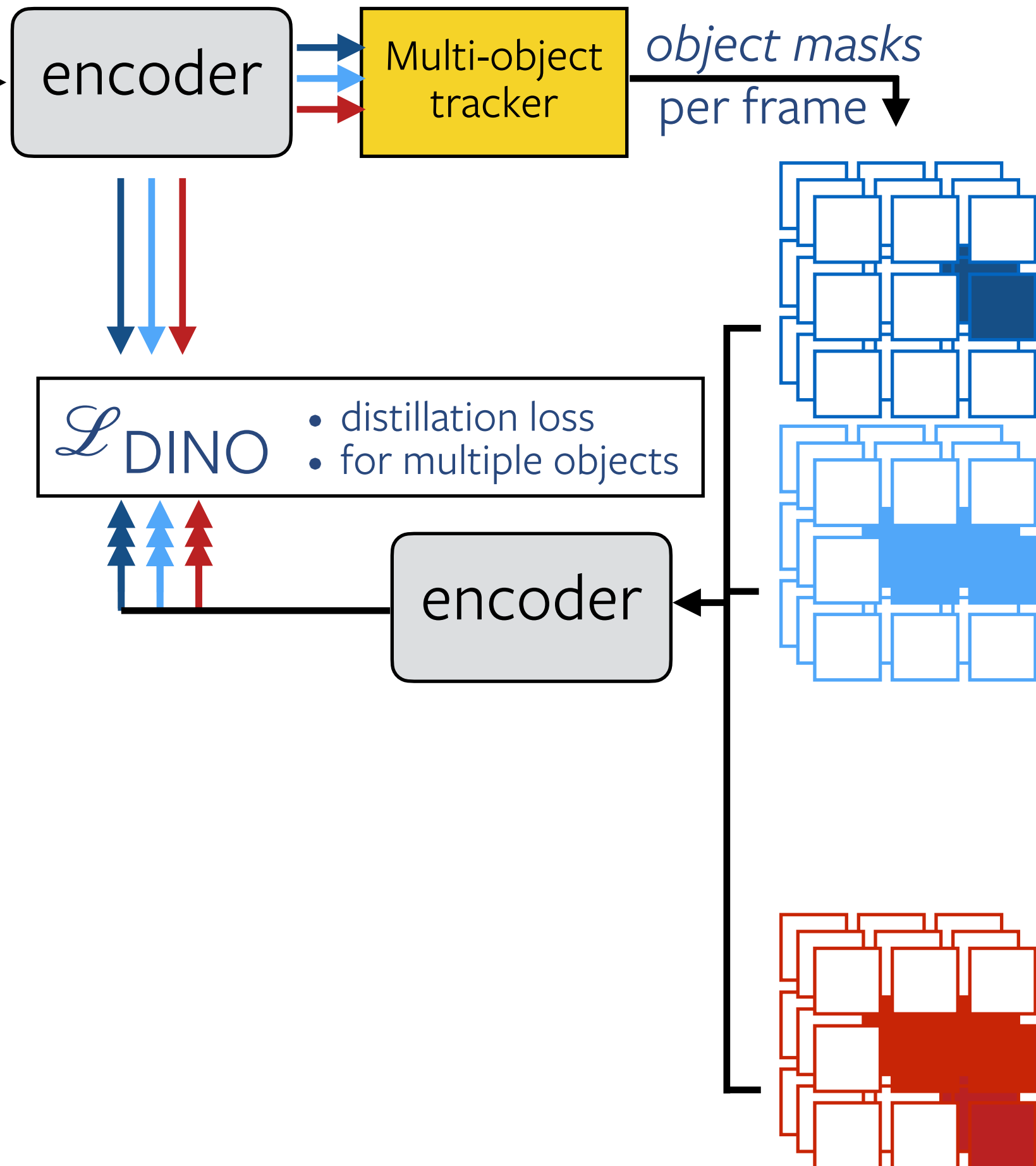
The dataset consists of 10x 4K videos of different cities' Walking Tours.



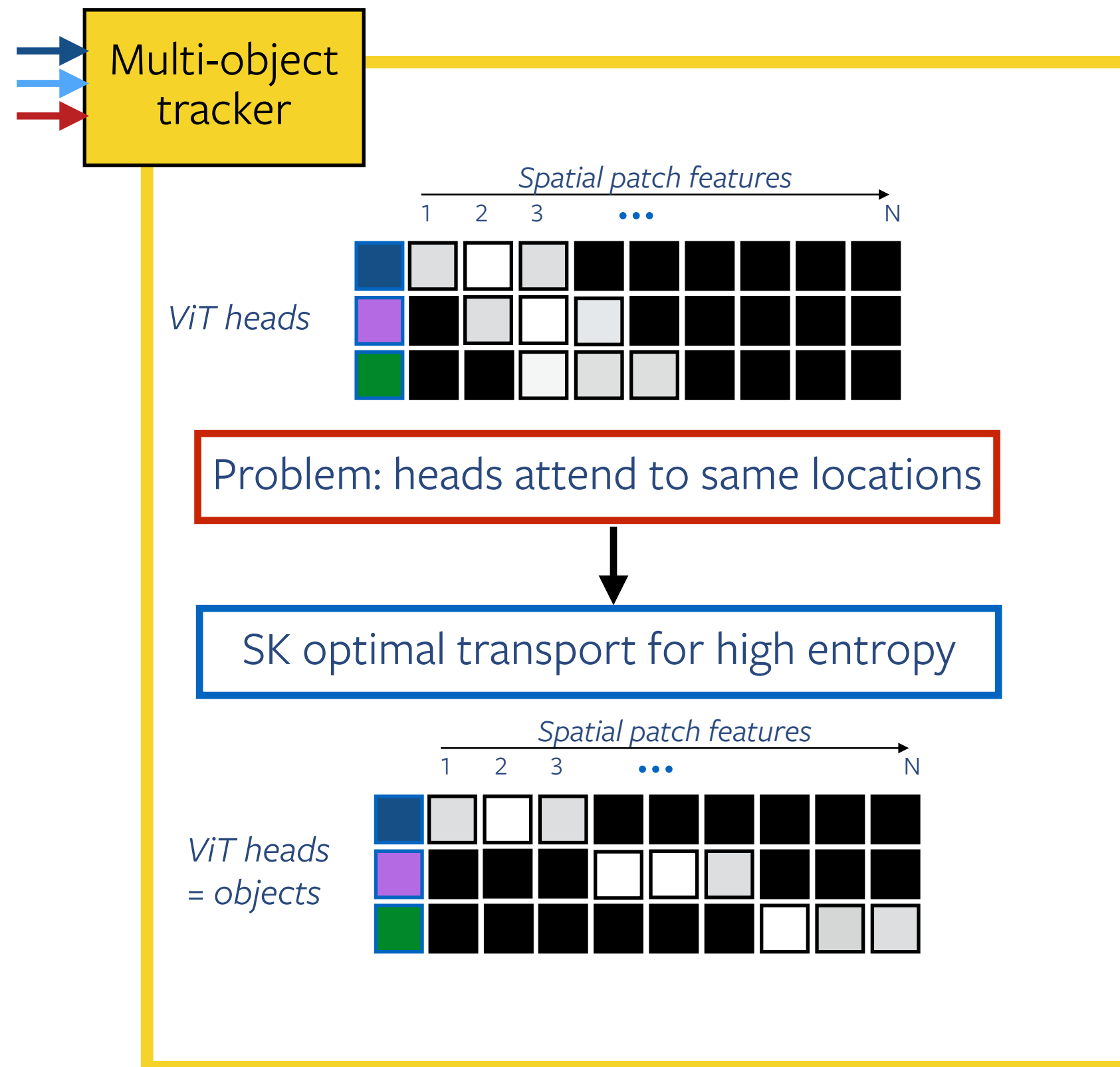
Dora: Discover and Track



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

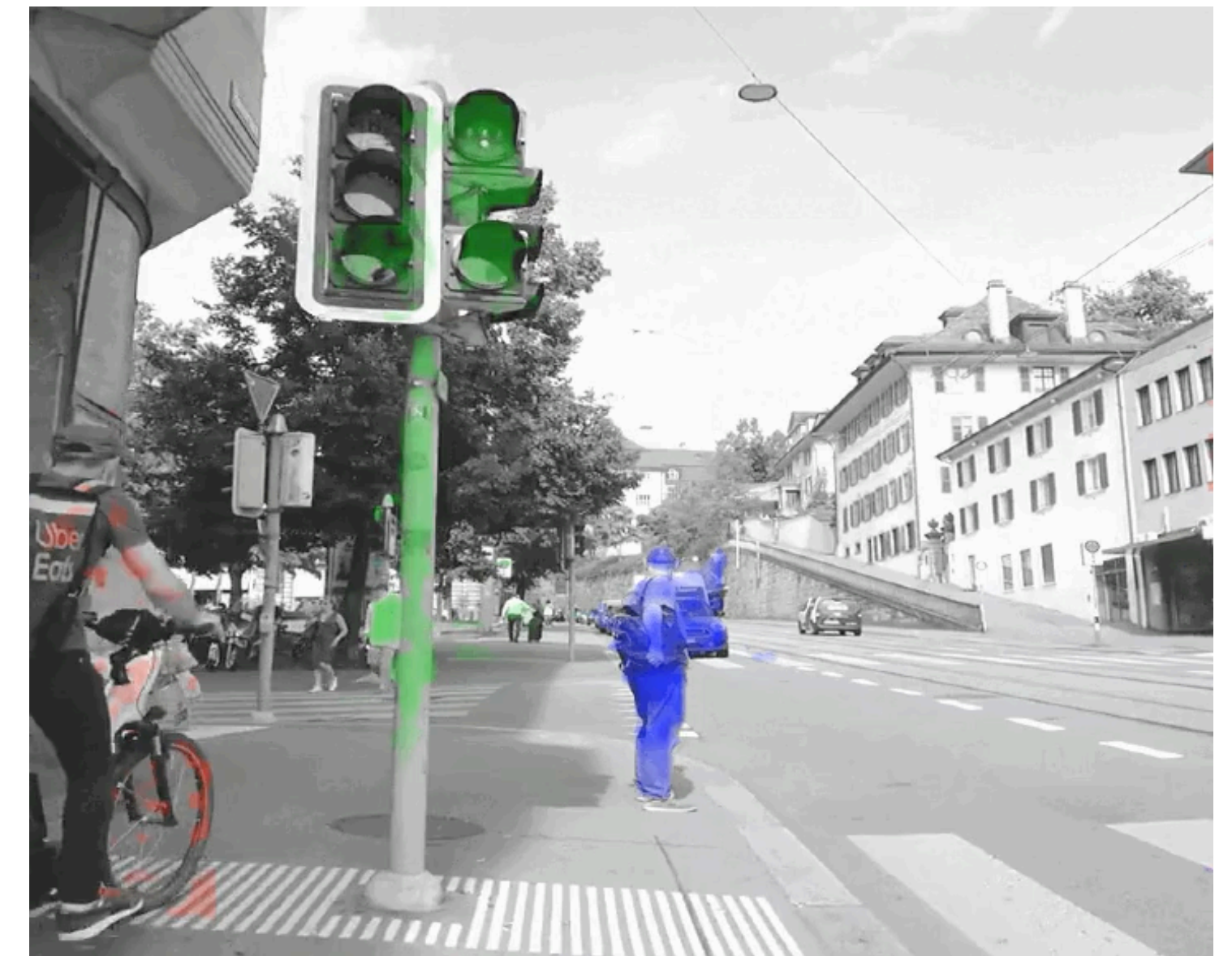
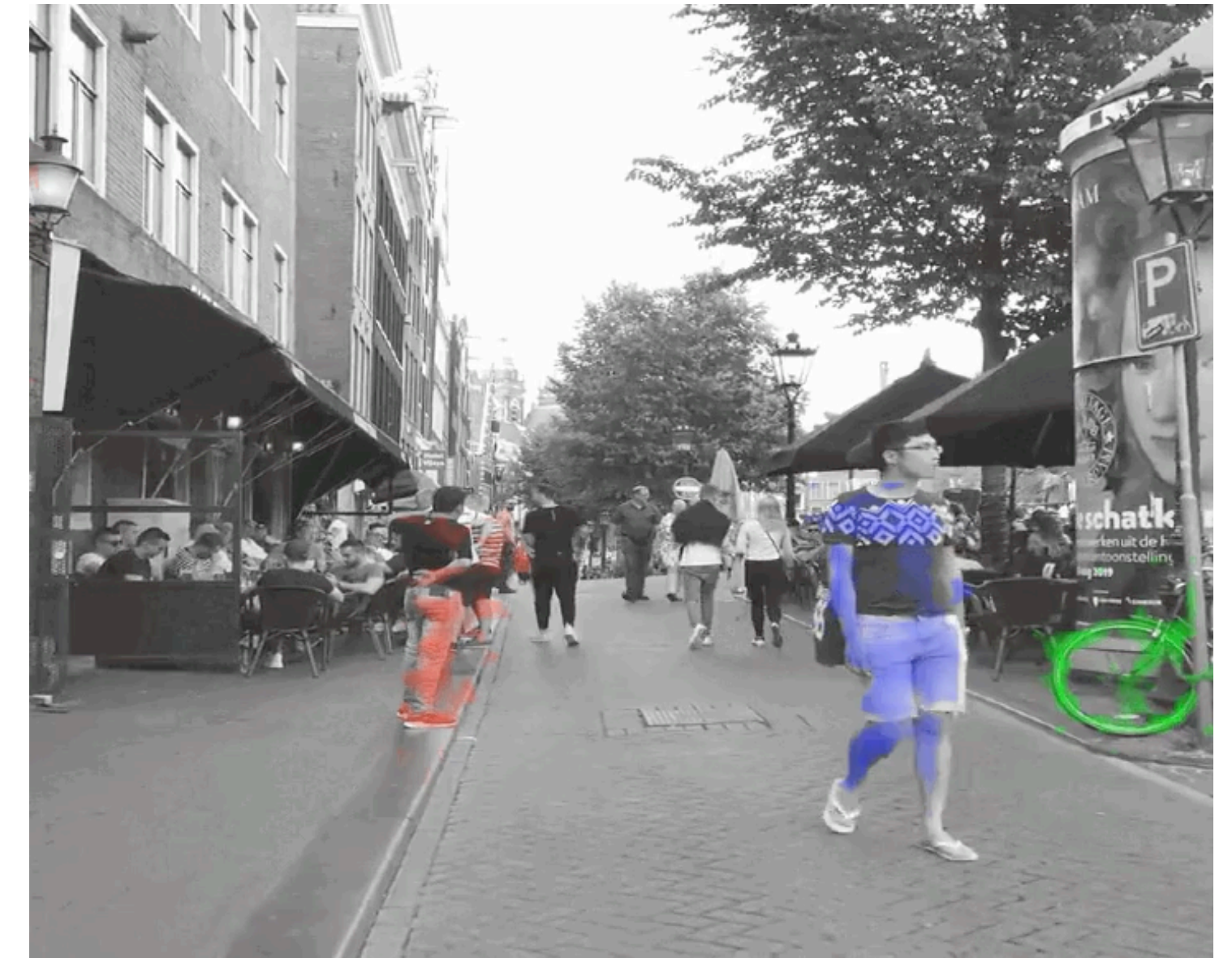


with SK

T'_t

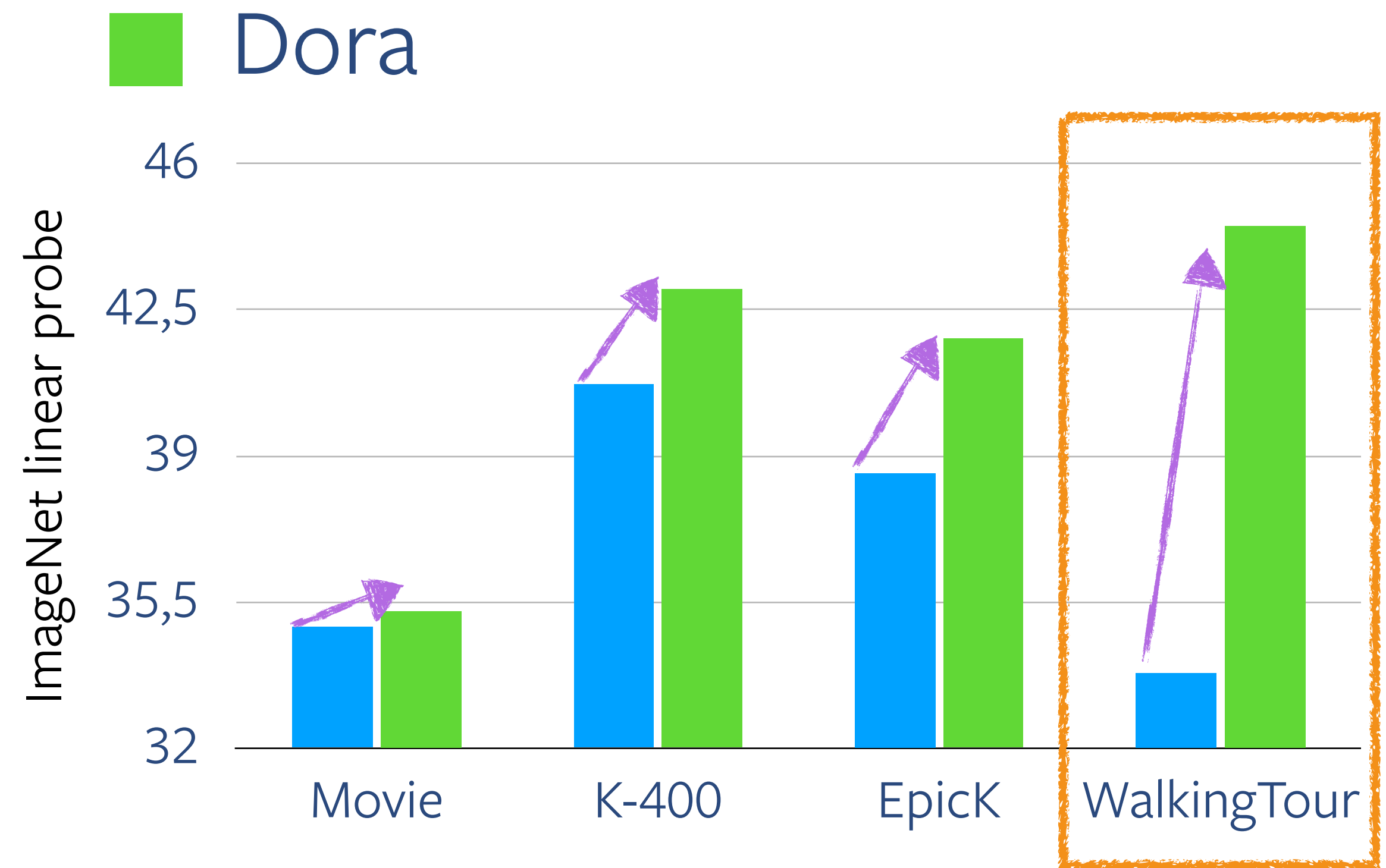
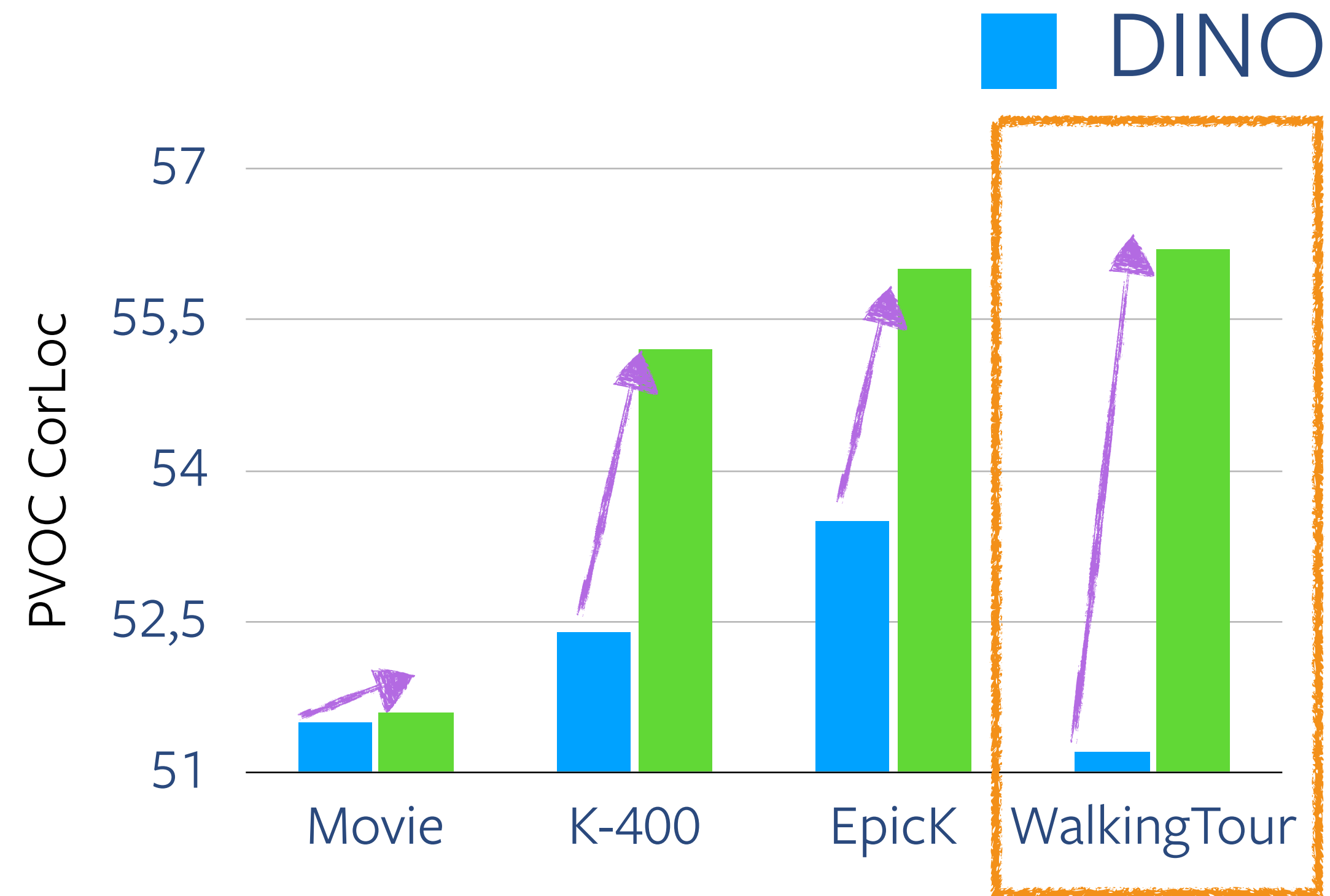


More examples: multi-object tracking in a ViT *emerges*



Dora better than DINO

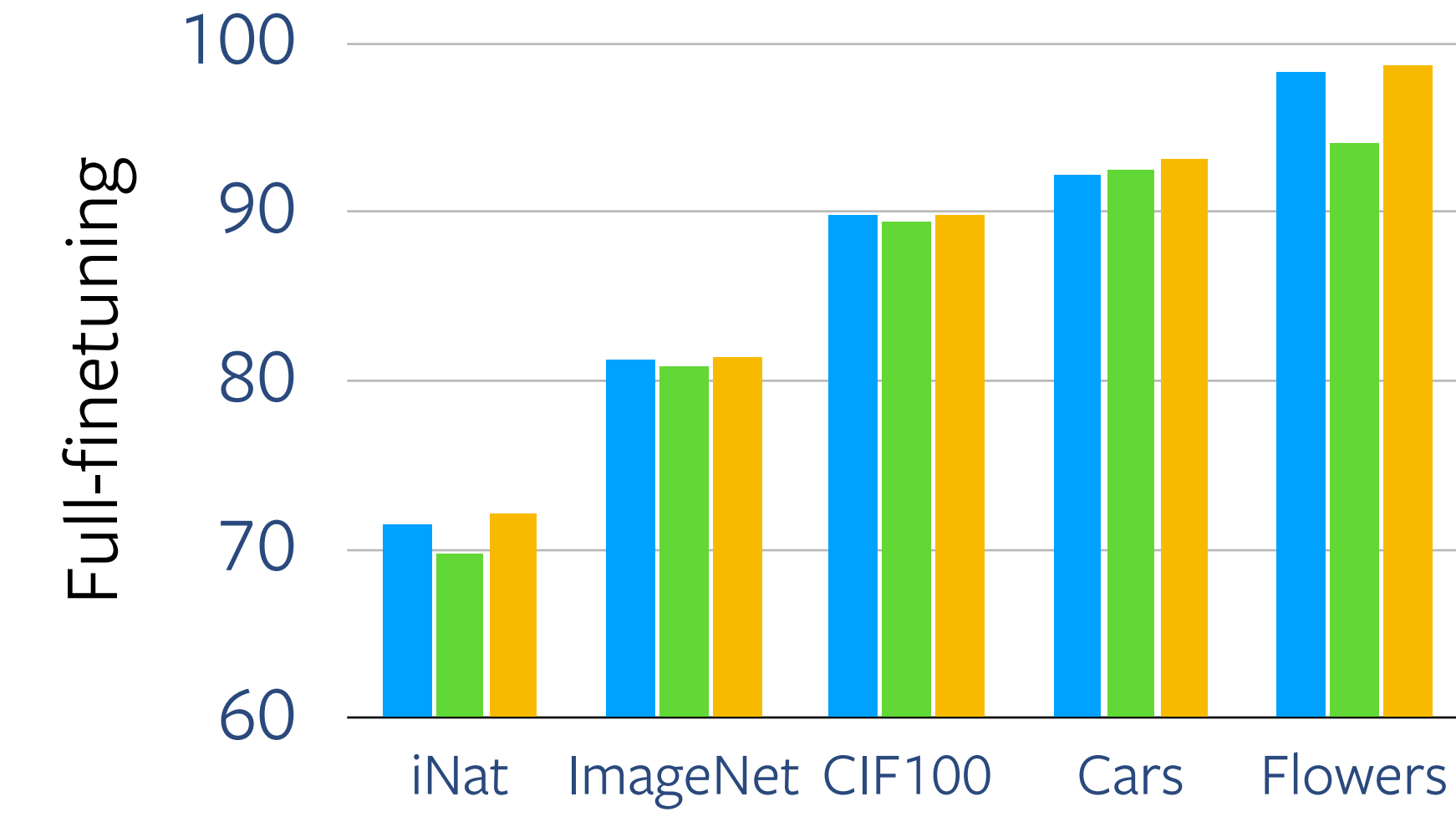
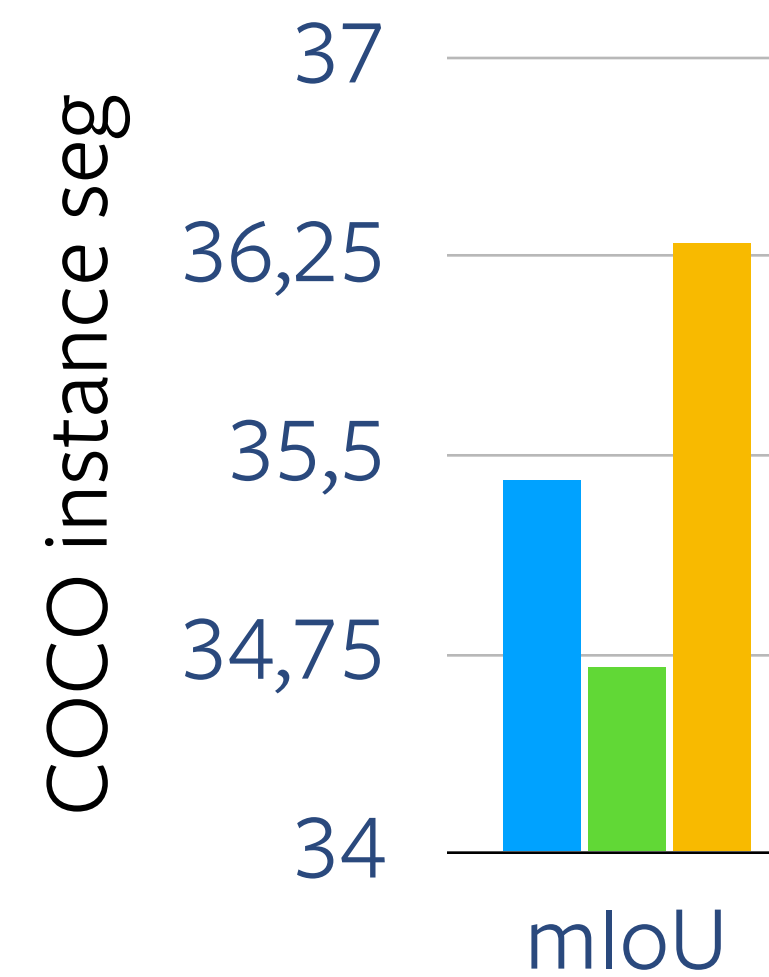
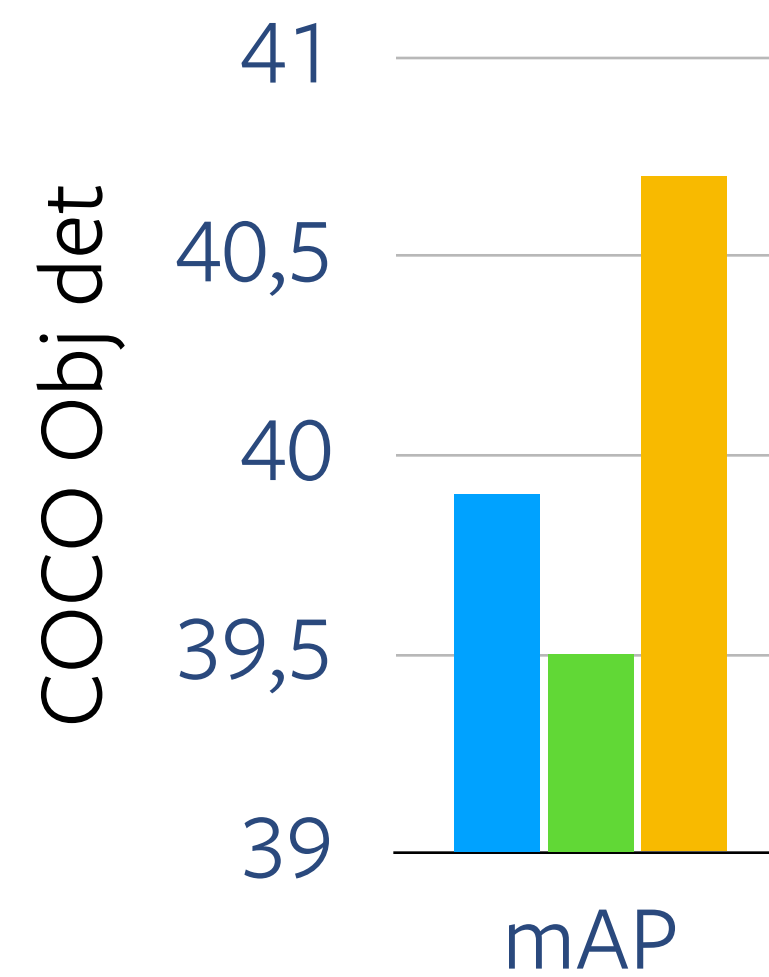
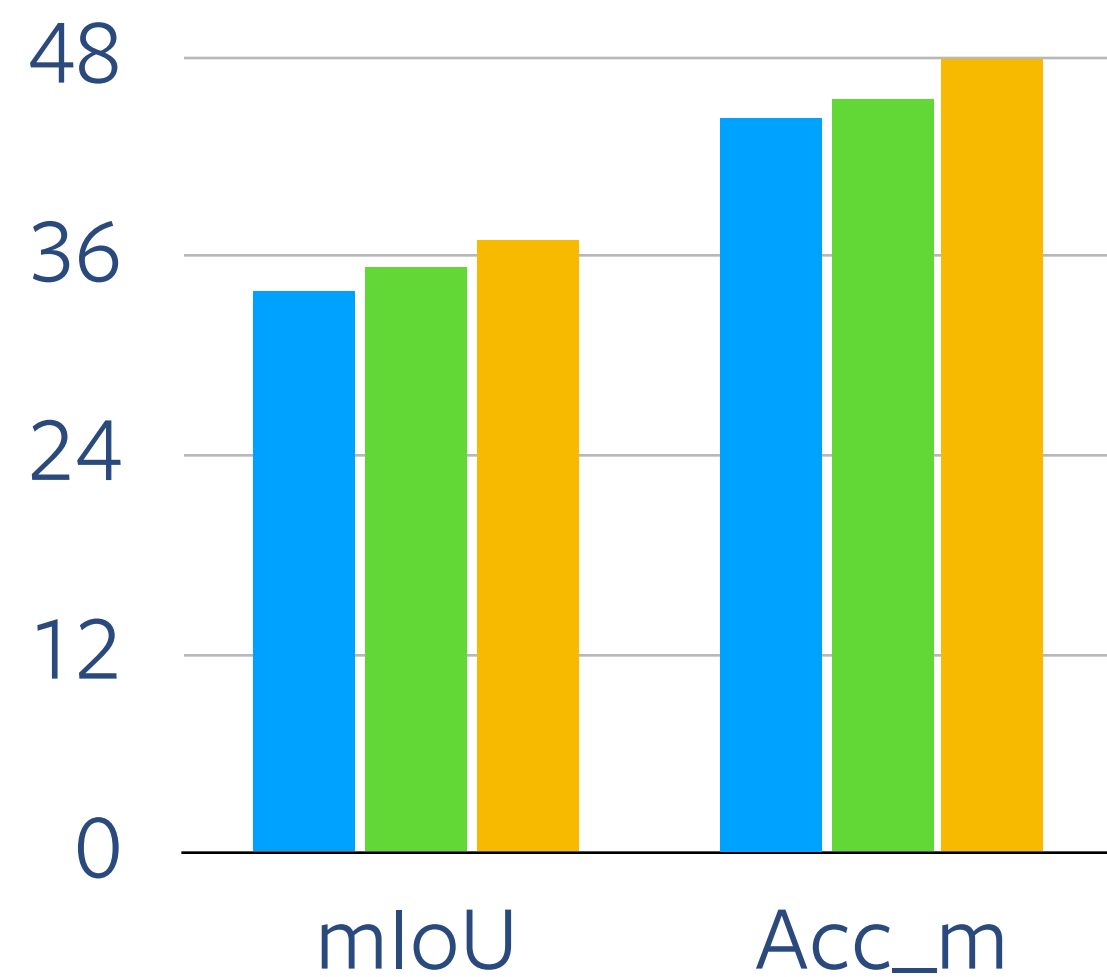
WT+ Dora: great match



But how does it compare against ImageNet pretraining?

■ DINO (IN-1k) ■ Dora (1 WT) ■ Dora (10 WT)

ADE20k Semantic Seg



Dora (1WT) ~ on par with DINO (IN-1k)
Dora (10WT) > DINO (IN-1k) everywhere

Summary

**CLUSTERING
WITH
SINKHORN-KNOPP**

[Asano et al. ICLR 2020]



**CLUSTERING
IN THE
SPATIAL DIMENSION**

[Ziegler & Asano. CVPR 2022]



**CLUSTERING
IN
SPACE AND TIME**

[Salehi et al. ICCV 2023]



**HAVING
MULTI-OBJECT
TRACKING EMERGE**

[Vekataraman et al. Arxiv 2023]

imgflip.com



1 TimeTuning:

DINO as init & use temporal info of videos.

How powerful is time without image-pretraining?

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

✓ Videos allow for strong self-supervised learning

VLMs?

Can we reduce the need for paired data?

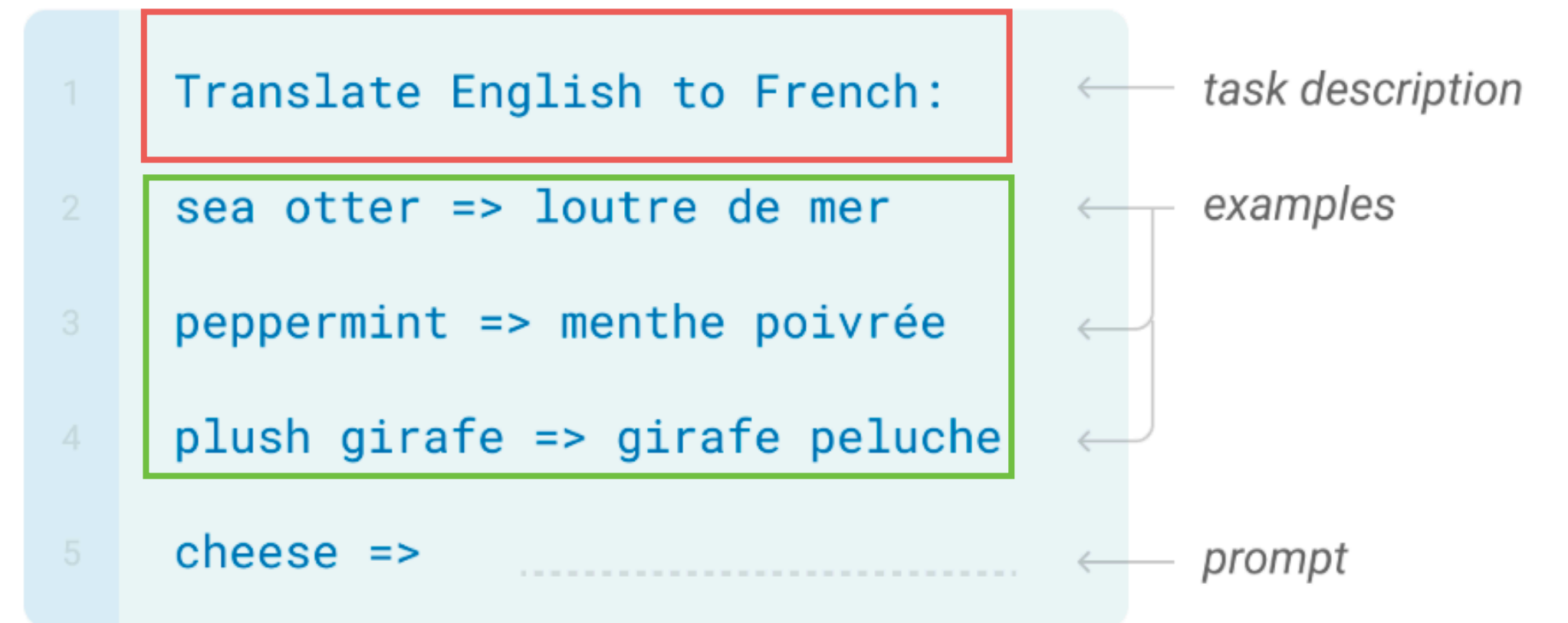
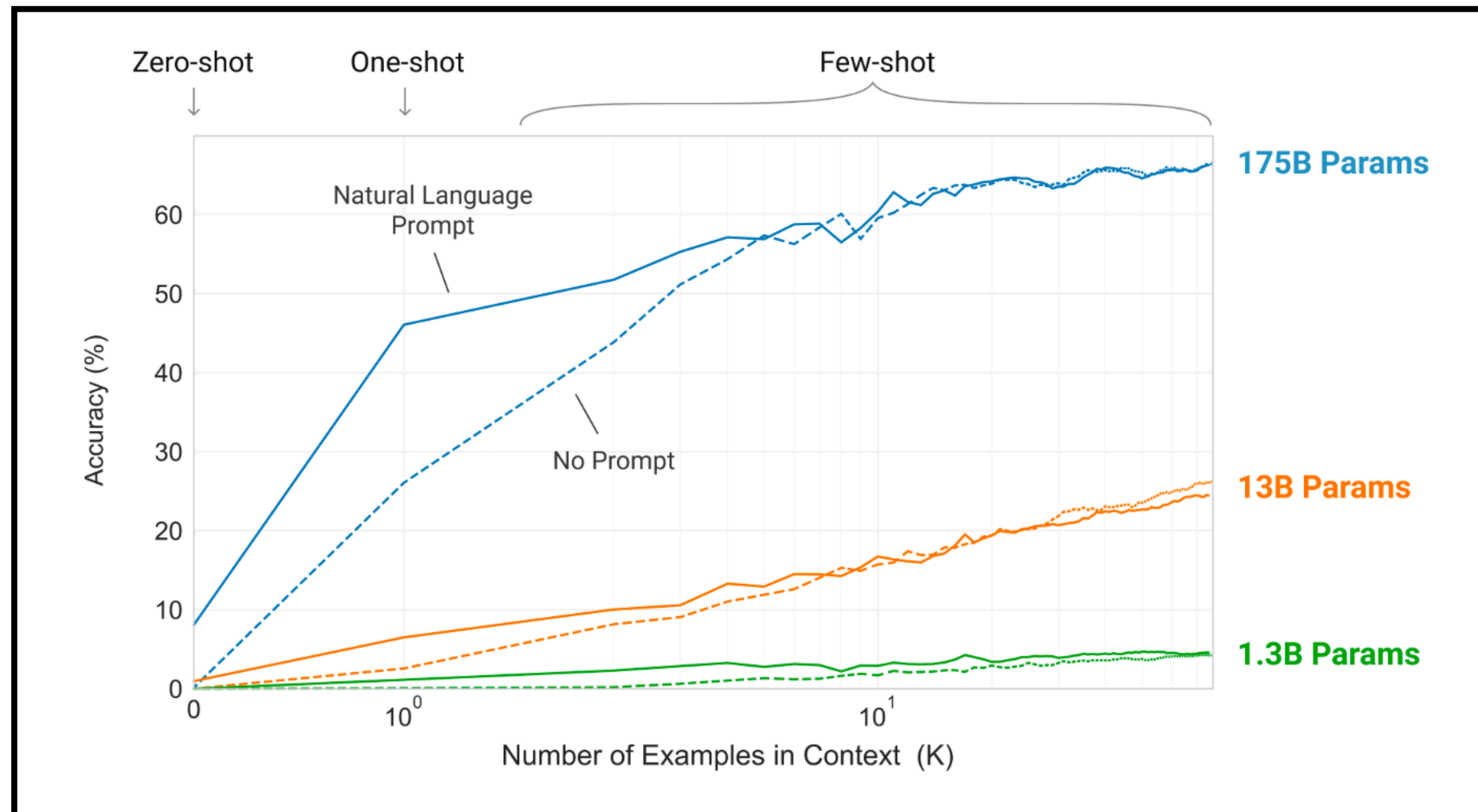
3 Use self-supervised features to create noisily paired data.



REDUCING THE NEED
FOR PAIRED
TEXT-IMAGE DATA

Dirk Jacobsz, *Painting a Portrait of His Wife*, 1550

Why are LLMs so sexy? a) scaling behavior, b) In-context Learning!

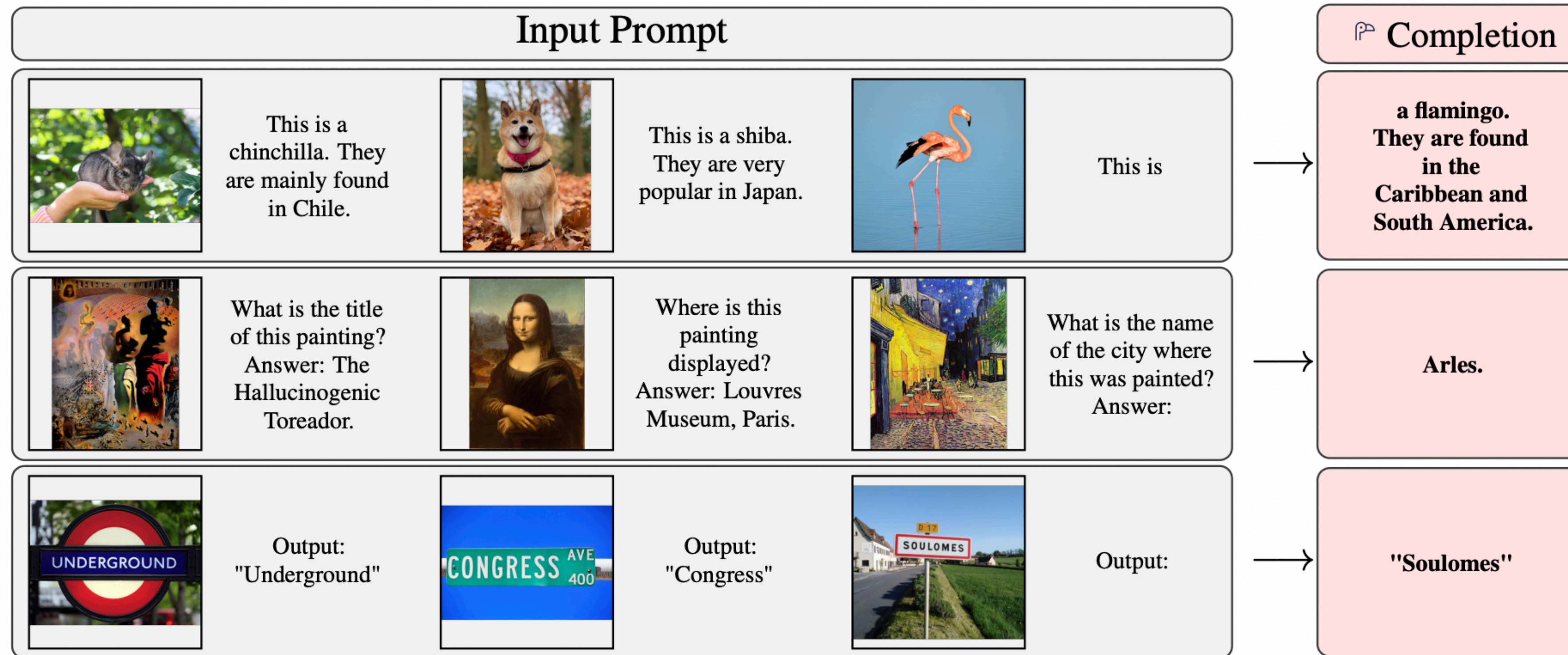


The possibility to

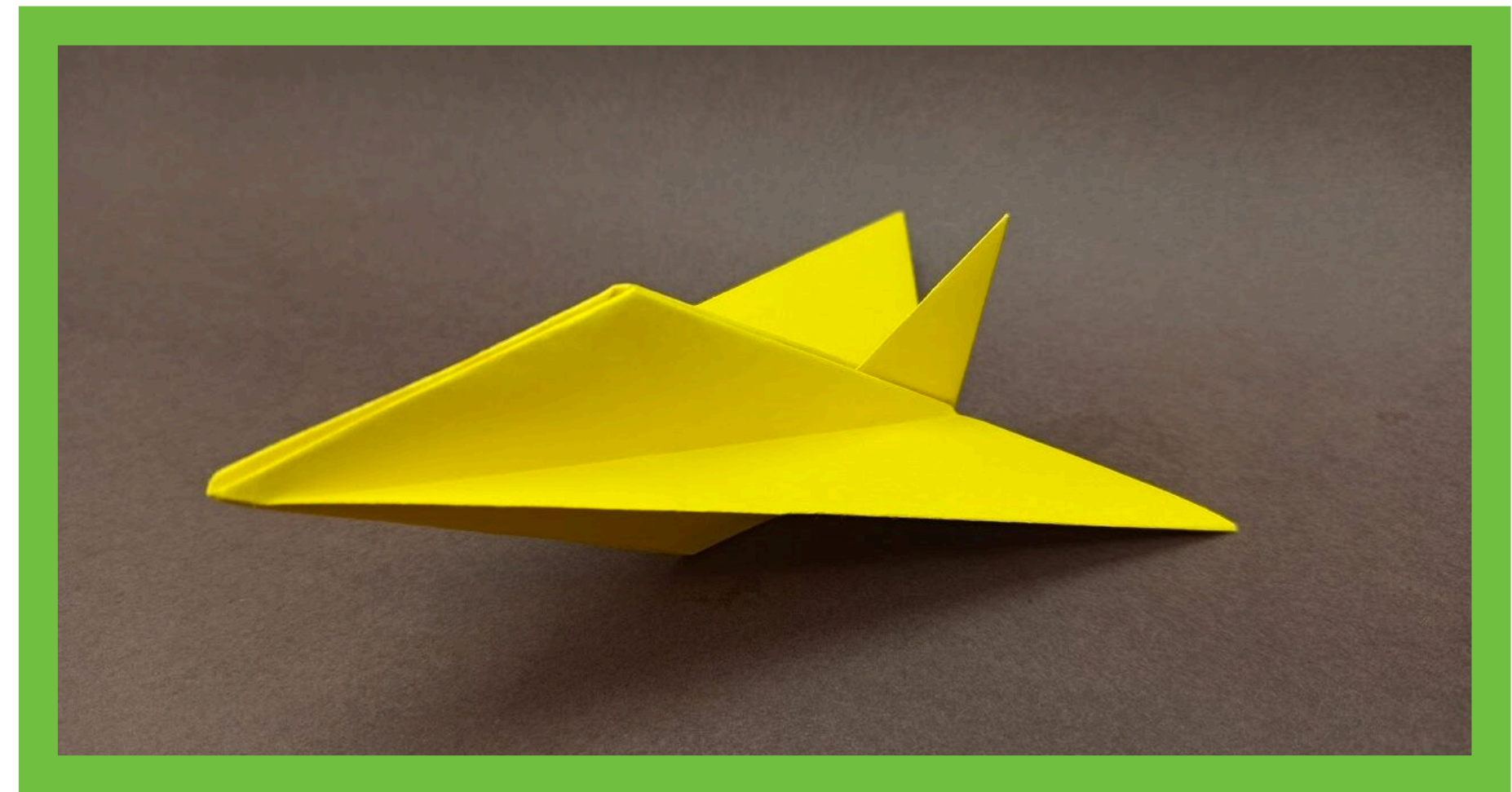
- *define a task* and the
- *learning-like* behaviour via few-shot examples
- with a *frozen model*

allows big scaling

In-context Learning emerges also for Visual Language Models

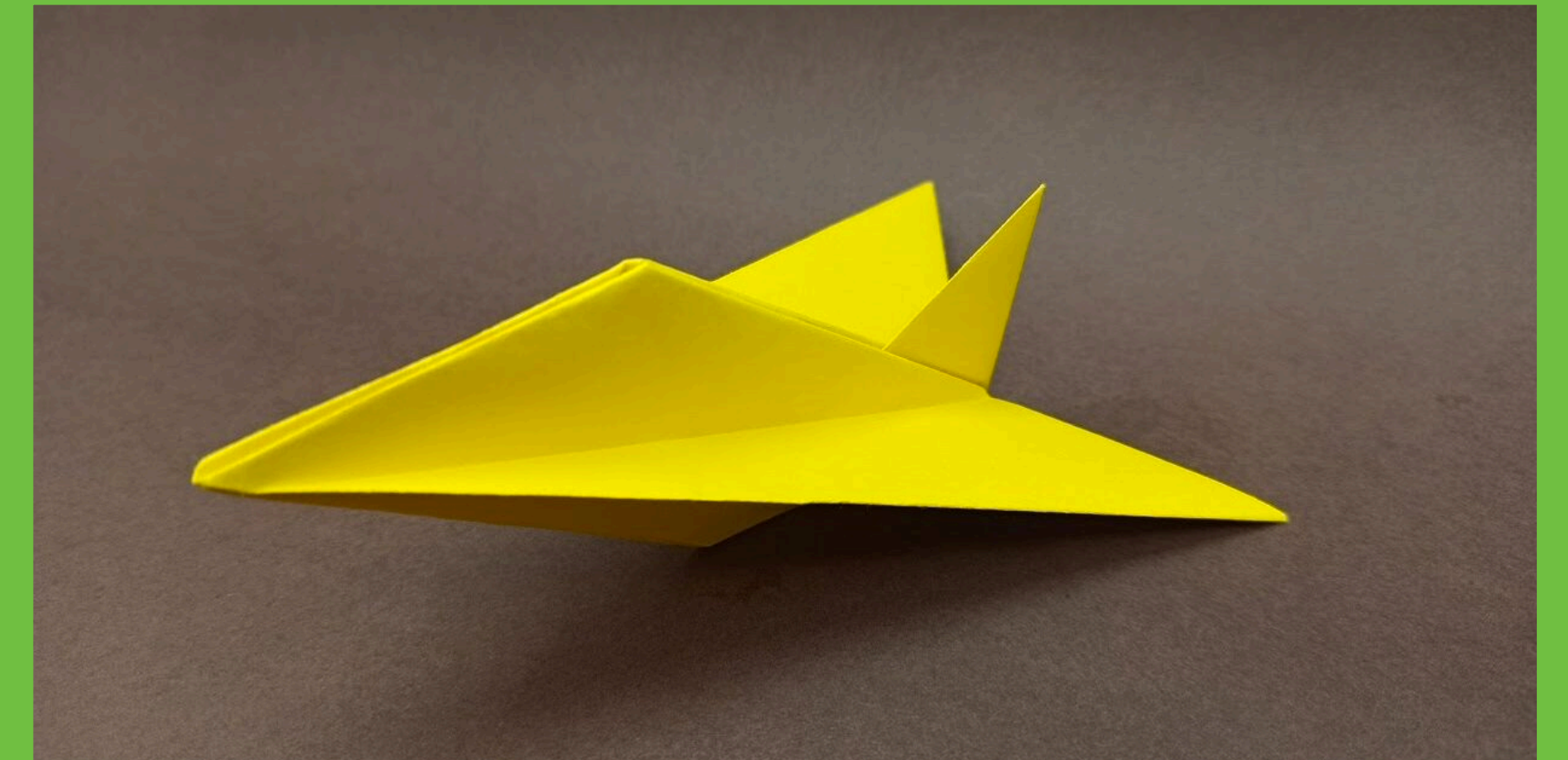


But just because it emerges for **6B+ sized models**, does it mean we cannot do this with more **light-weight ones**?



Our goal

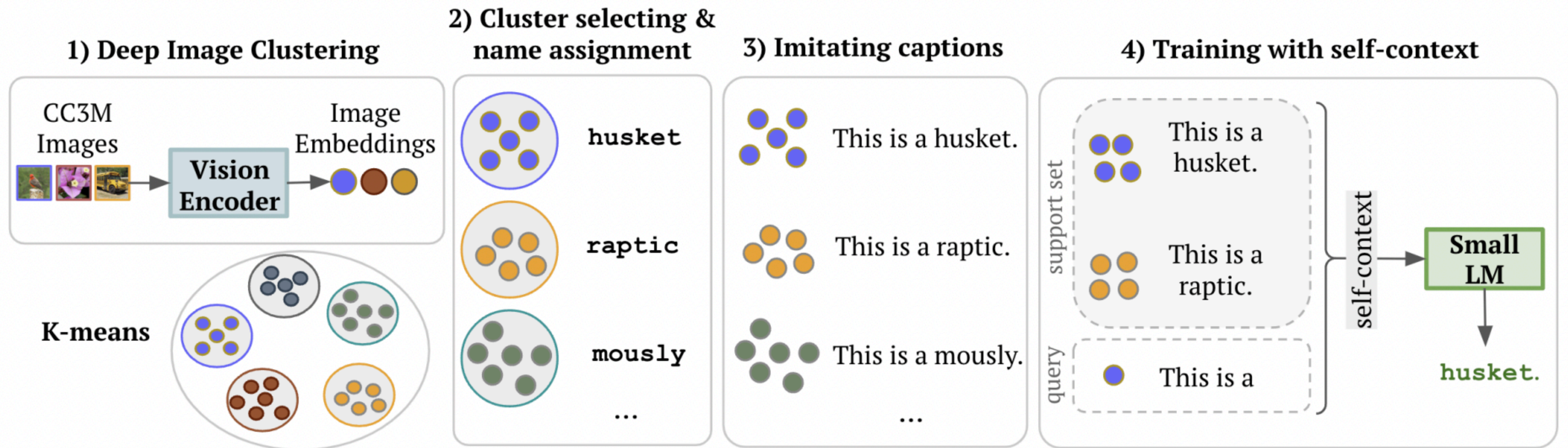
open-ended multi-modal ICL
with small VLMs
without using supervised data.



Why?

- Ultimately, paired data is rare
- ICL as algorithm: symbols replaceable
- As an existence-proof

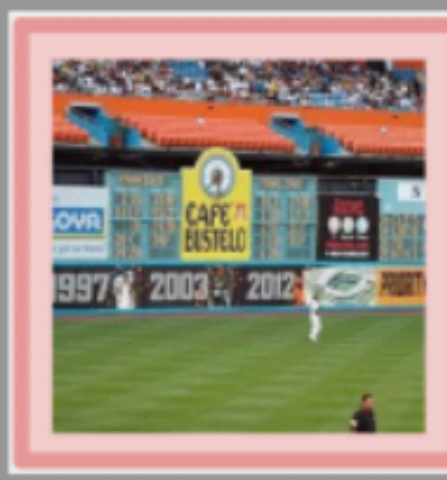
How? Our method simply *mimics* supervised data by using SSL



Where are we?

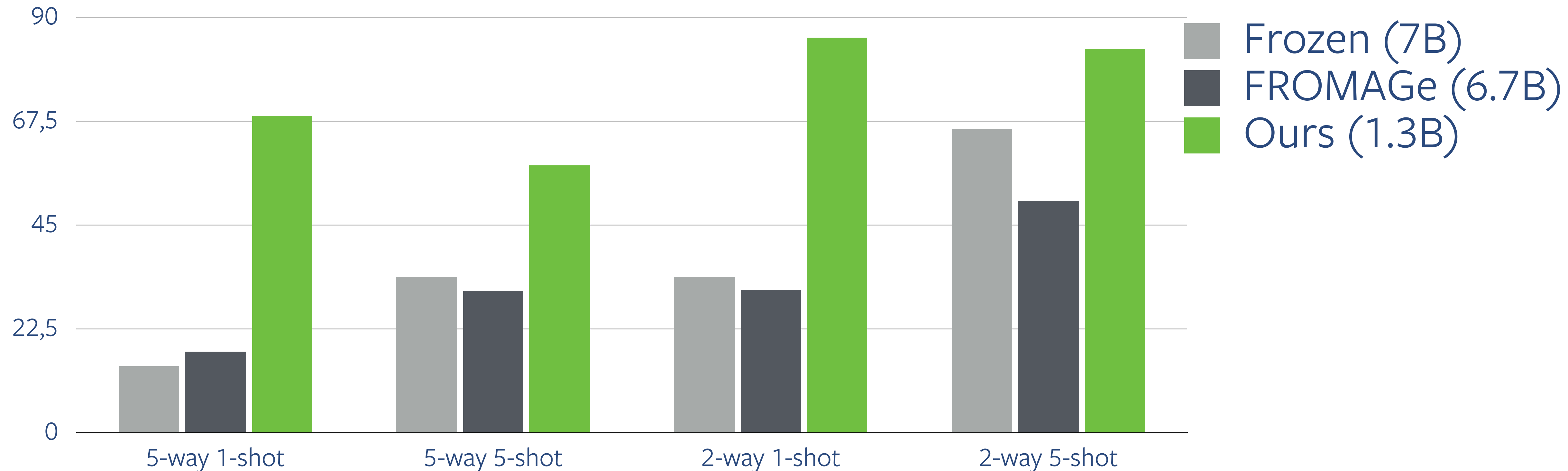


We train with fake names, but evaluation works just fine!

			<p>✗ <i>ClipCap</i>: distinguished from all other by its long slender torso. ✗ <i>FROMAGe</i>: school bus that is parked in the school yard. ✓ <i>SeCAt (Ours)</i>: school bus.</p>
<p>This is a scoreboard.</p>	<p>This is a school bus.</p>	<p>This is a <?></p>	

Simple-as-that; beats 6B-sized models on open-ended multi-modal classif.

Open-ended mini-ImageNet ICL evaluation



Team for the works presented

TimeTuning



Mohammadreza Salehi



Efstratios Gavves



Cees G. M. Snoek



Yuki M. Asano

WTour Dora



Shashanka Venkataramanan



Mamshad N Rizve



Joao Carreira

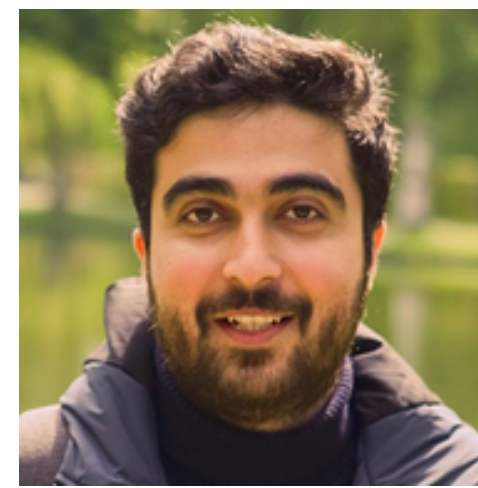


Yuki M. Asano*



Yannis Avrithis*

SeCA_t



Mohammad M Derakhshani¹



Ivona Najdenkoska¹



Cees G. M. Snoek*



Marcel Worring*



Yuki M. Asano*

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*. ArXiv 2023

Derakhshani, Najdenkoska, Snoek, Worring, Asano. *Self-Supervised Open-Ended Classification with Small Visual Language Models*. ArXiv 2023.

¹: co-first authors; *: co-last authors

Especially videos open exciting new directions

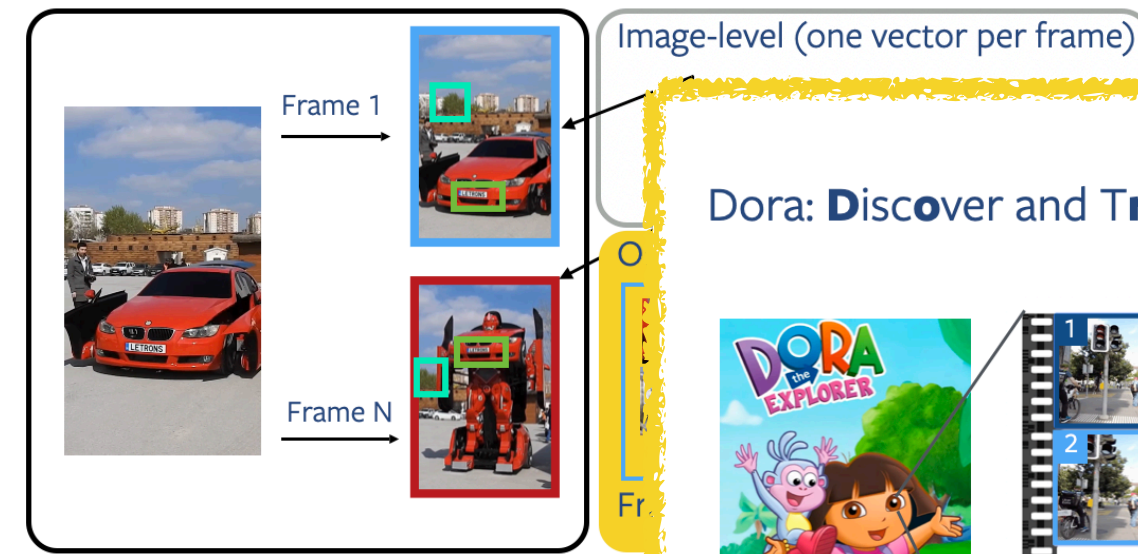


Visual development for AI

Bonus: insan

UNIVERSITY OF AMSTERDAM

Solution is obvious



UNIVERSITY OF AMSTERDAM

Salehi, Gavves, Snoek, Asano. Time does tell: self-supervised time-

Dora: Discover and Track



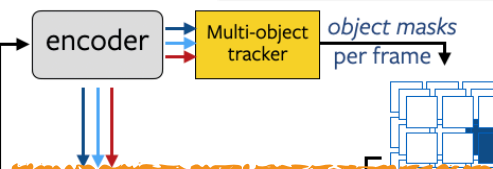
Much like Dora, we walk around and learn from what we see.

UNIVERSITY OF AMSTERDAM

Venkataramanan, Ribeiro, Carreira, Avrithis, Asano. In

High-level idea:

- 1) track multiple objects across time
- 2) enforce invariance of features across time



We train with fake names, but evaluation works just fine in-context

			<input checked="" type="checkbox"/> ClipCap : distinguished from all other by its long slender torso.
<input checked="" type="checkbox"/> FROMAGE : school bus that is parked in the school yard.			<input checked="" type="checkbox"/> SeCat (Ours) : school bus.
This is a scoreboard.	This is a school bus.	This is a <?>	
			<input checked="" type="checkbox"/> ClipCap : close up of a alpine sea holly.
<input checked="" type="checkbox"/> FROMAGE : bird that is native to the United States and Canada.			<input checked="" type="checkbox"/> SeCat (Ours) : gray kingbird.
This is a desert-rose	This is a gray kingbird.	This is a <?>	
			<input checked="" type="checkbox"/> ClipCap : daffodil in my garden.
<input checked="" type="checkbox"/> FROMAGE : russian blue cat.			<input checked="" type="checkbox"/> SeCat (Ours) : egyptian mau.
This is a russian blue.	This is a egyptian mau	This is a <?>	

UNIVERSITY OF AMSTERDAM

Derakhshani, Najdenkoska, Snoek, Worring, Asano. Small Visual Language Models can also be Open-Ended Few-Shot Learners. ArXiv 2023.

31

Future Foundation Models will be massively pretrained with videos. Current multi-modal training will become only the cherry on top.