



# Video Understanding

# Definitions...

**Ego**... a person's sense of self-esteem or self-importance



# Definitions...

**Ego**... a person's sense of self-esteem or self-importance

**Egocentric vision**... the wearer serves as the central reference point in the study of interesting entities: objects, actions, interactions and intentions



# In today's tutorial



Motivation and Datasets in  
Egocentric Video Understanding



Video Understanding  
Out of the Frame



Video Understanding:  
Data and Tasks



Teaser: The Wizard of Oz  
at the Sphere



Videos are Multimodal



Outlook into the Future of  
Egocentric Vision



Connected Videos of One's Life



Conclusion

# In today's tutorial



Motivation and Datasets in  
Egocentric Video Understanding



Video Understanding  
Out of the Frame



Video Understanding:  
Data and Tasks



Teaser: The Wizard of Oz  
at the Sphere



Videos are Multimodal



Outlook into the Future of  
Egocentric Vision



Connected Videos of One's Life

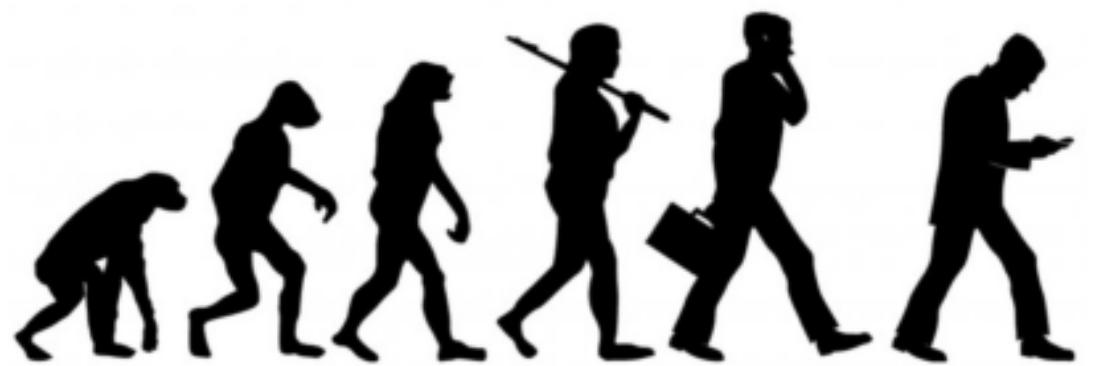


Conclusion

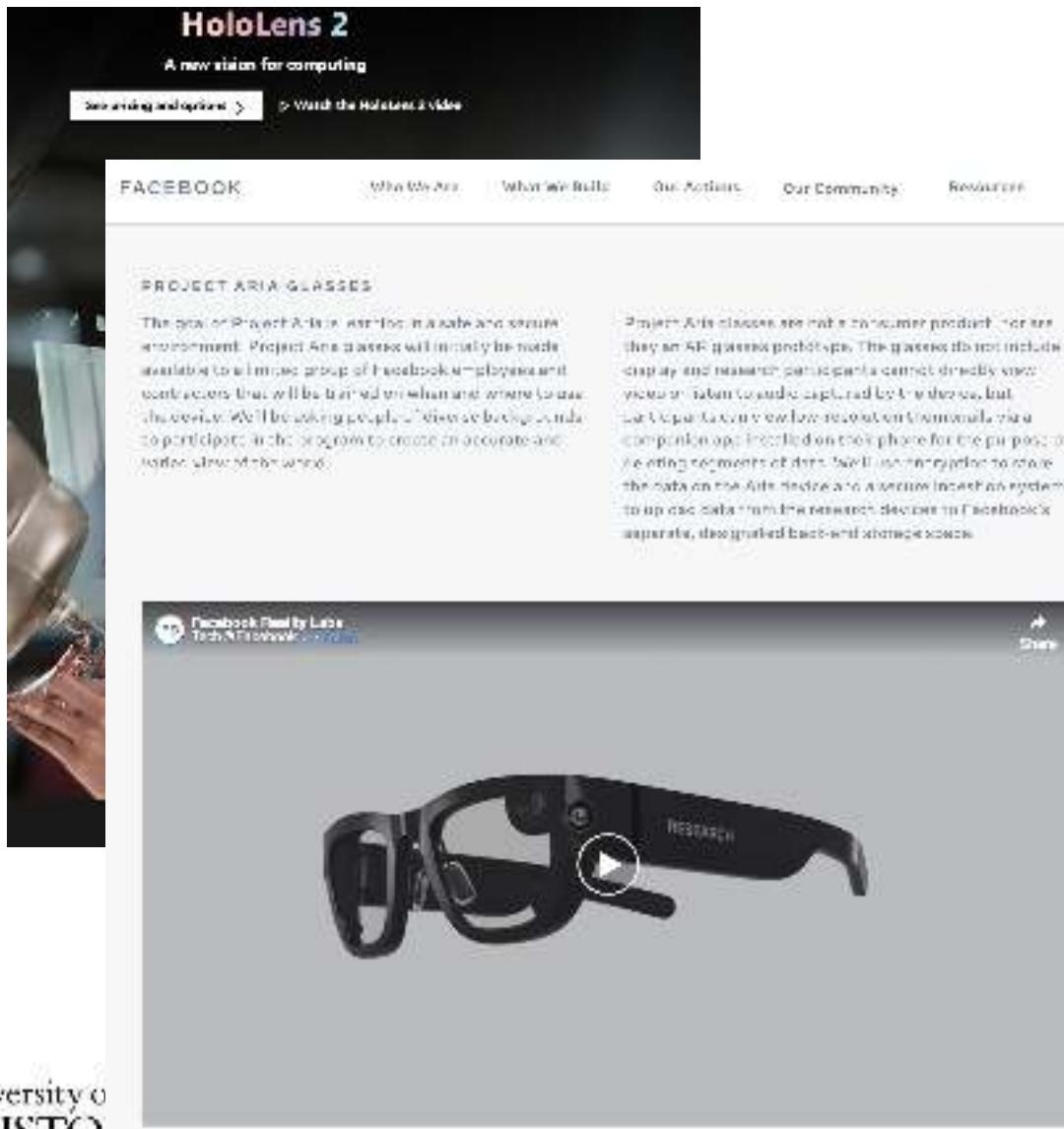
# The present...



Photo \*Illustration\* by Pelle Cass



# The future...



**HoloLens 2**  
A new vision for computing  
[See Inside and Optimize](#) > [Watch the HoloLens 2 video](#)

[FACEBOOK](#) [About Microsoft](#) [What we build](#) [Our Actions](#) [Our Community](#) [Resources](#)

**PROJECT ARIA GLASSES**

The goal of Project Aria is to test in a safe and secure environment. Project Aria glasses will initially be made available to a limited group of Facebook employees and contractors that will be trained on what and where to use the device. We'll be asking people to view us being asked to participate in the program to create an accurate and wider view of the world.

**Facebook Reality Labs**  
Tech Newsroom | [Facebook Reality Labs](#)



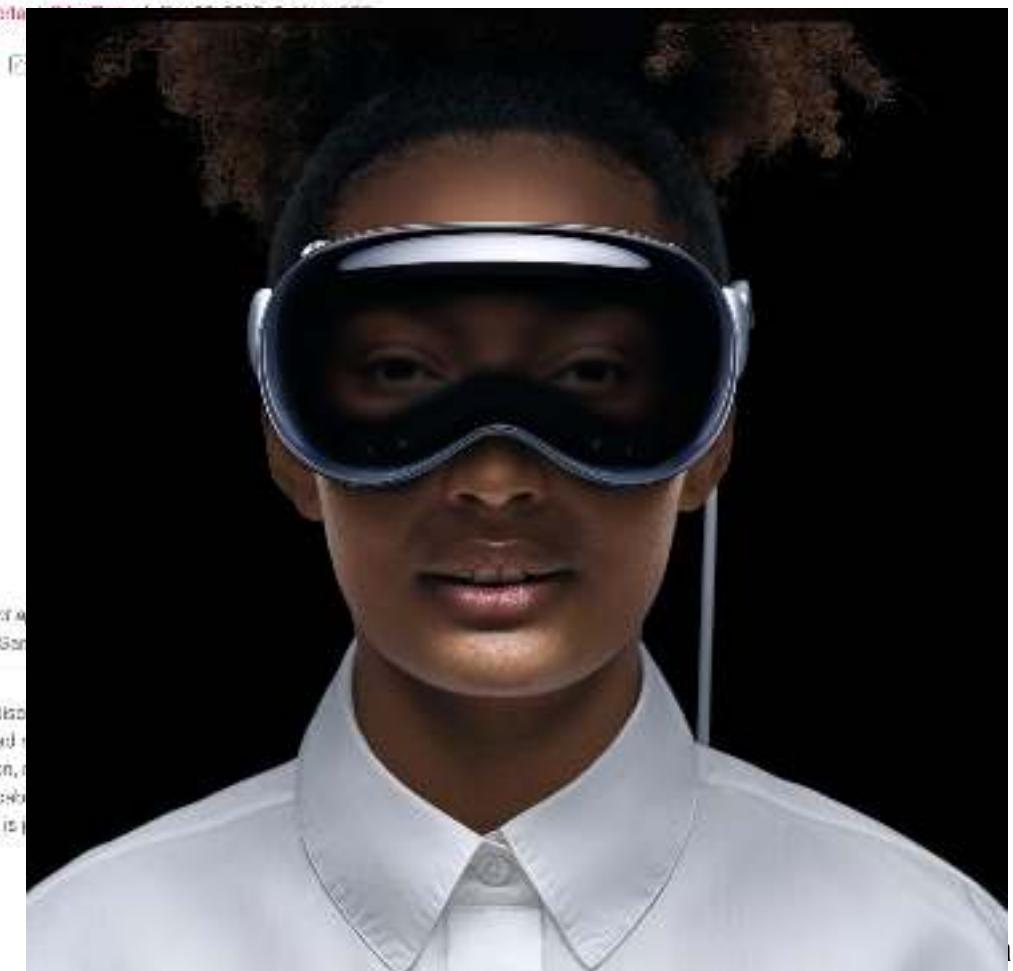
Share

he patent application was filed by Samsung Electronics Co. Ltd. on April 12, 2018, and was published on June 14, 2018. It describes a system for tracking the position of a user's head and eyes, and for displaying virtual content based on the tracked information. The system includes a head-mounted display (HMD) and a tracking unit that captures images of the user's eyes and head. The tracking unit also includes a camera and a light source. The system uses the captured images to track the user's gaze and head movements, and to generate a 3D model of the user's face. The system then displays virtual content on the HMD based on the tracked information. The system can be used for various applications, such as gaming, entertainment, and education.

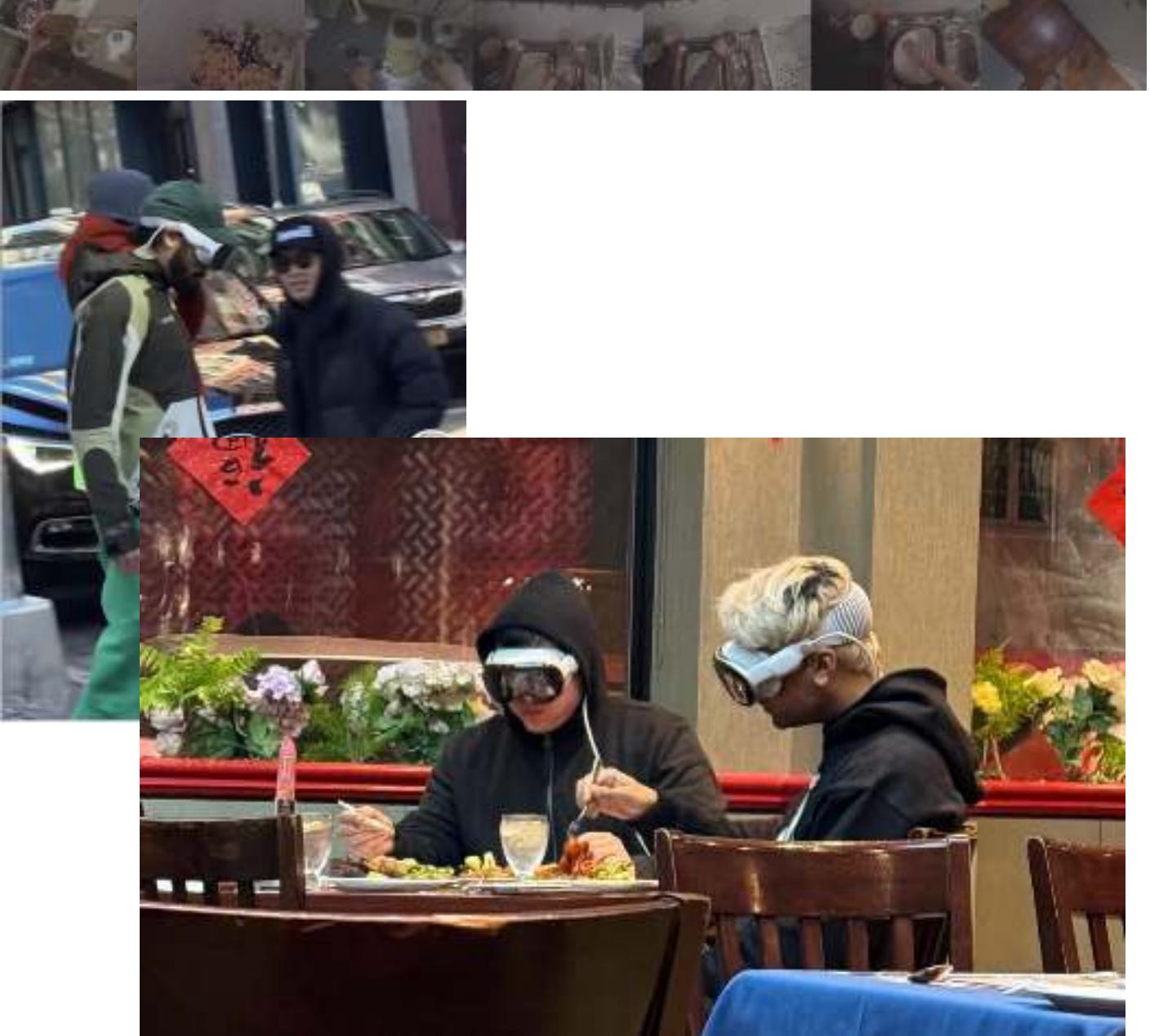
**University of BRISTOL**

## Samsung patent application reveals augmented reality headset design

*It comes as the Gear VR slowly fades away*



# The future is here...



# The future can be imagined...



# Egocentric Videos?



# Egocentric Videos?



Damen  
S2025

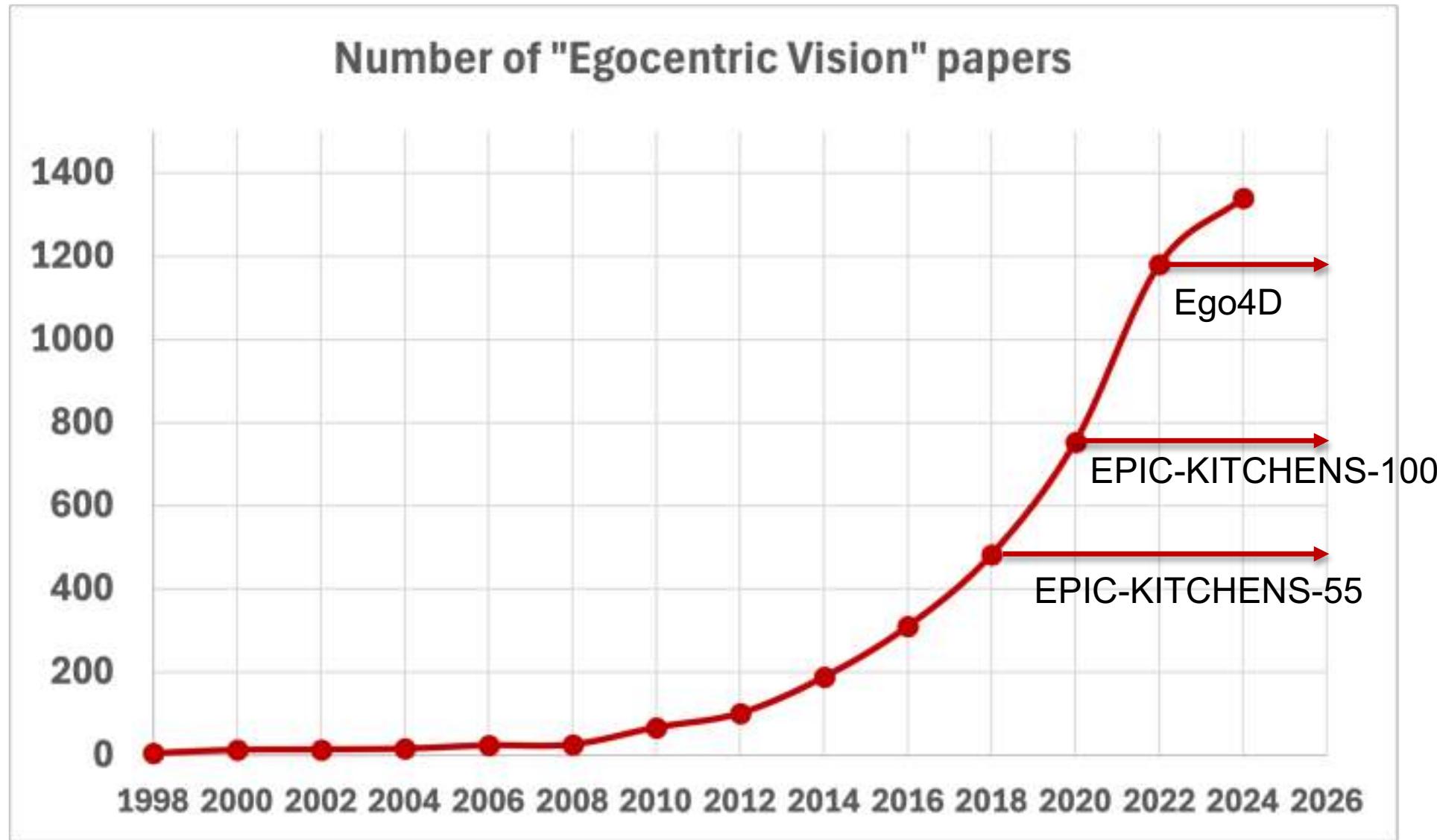
# Machine Learning in Practice



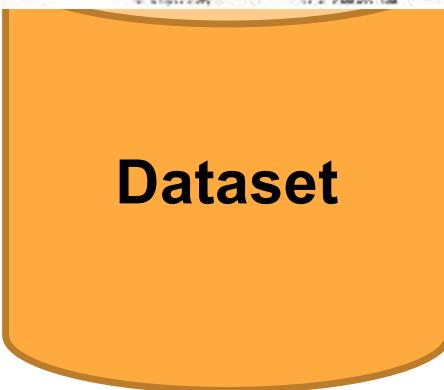
**no data  
no machine  
learning research**



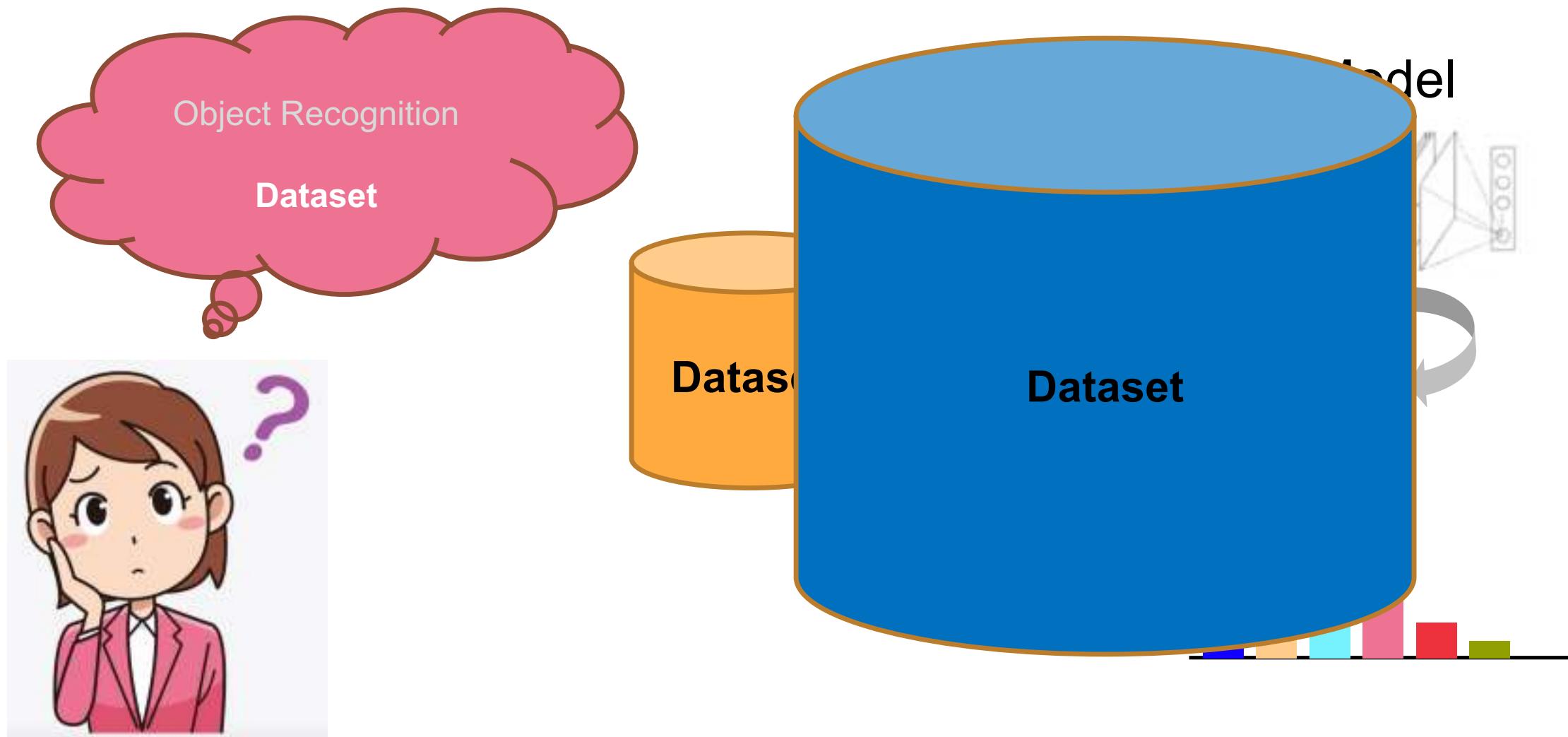
# Is research catching up?



# Machine Learning in Practice



# Machine Learning in Practice



# Machine Learning in Practice

- Applies to Most ML research at the moment
    - Object Recognition (Pascal, ImageNet, Places, ...)
    - Action Recognition (Kinetics-400, -600, -700, AVA, SS, ...)
    - ...
  - Datasets:
    - Methods overfit to the dataset
    - useful for **one** task
    - unnaturally balanced (or nearly balanced) – unrelated to priors outside the dataset itself
- 

# Machine Learning in Practice

- Autonomous Driving...

## Welcome to the KITTI Vision Benchmark Suite!

We take advantage of our autonomous driving platform [Annieway](#) to develop novel challenging real-world computer vision benchmarks. Our tasks of interest are: stereo, optical flow, visual odometry, 3D object detection and 3D tracking. For this purpose, we equipped a standard station wagon with two high-resolution color and grayscale video cameras. Accurate ground truth is provided by a Velodyne laser scanner and a GPS localization system. Our datasets are captured by driving around the mid-size city of [Karlsruhe](#), in rural areas and on highways. Up to 15 cars and 30 pedestrians are visible per image. Besides providing all data in raw format, we extract benchmarks for each task. For each of our benchmarks, we also provide an evaluation metric and this evaluation website. Preliminary experiments show that methods ranking high on established benchmarks such as [Middlebury](#) perform below average when being moved outside the laboratory to the real world. Our goal is to reduce this bias and complement existing benchmarks by providing real-world benchmarks with novel difficulties to the community.

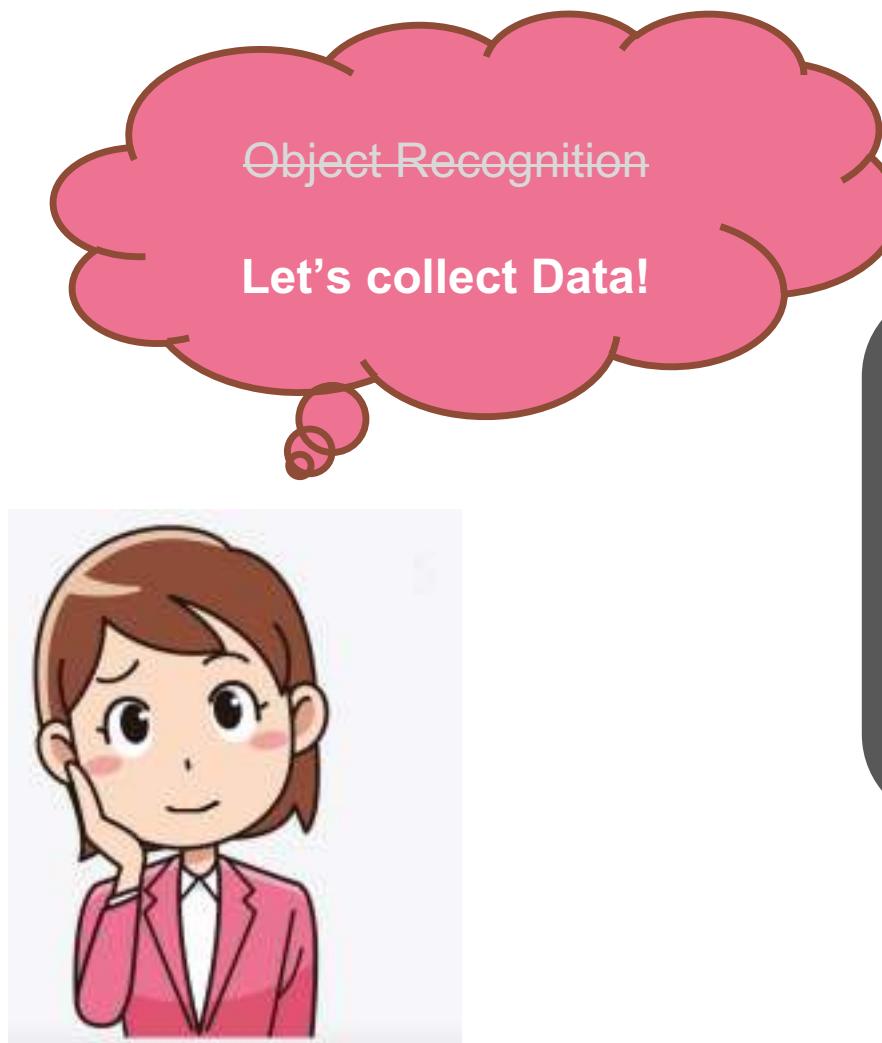
 Share



To get started, grab a cup of your favorite beverage and watch our video trailer (5 minutes):

stereo flow sceneflow depth odometry object tracking road semantics raw data

# Machine Learning in Practice



# EPIC-KITCHENS



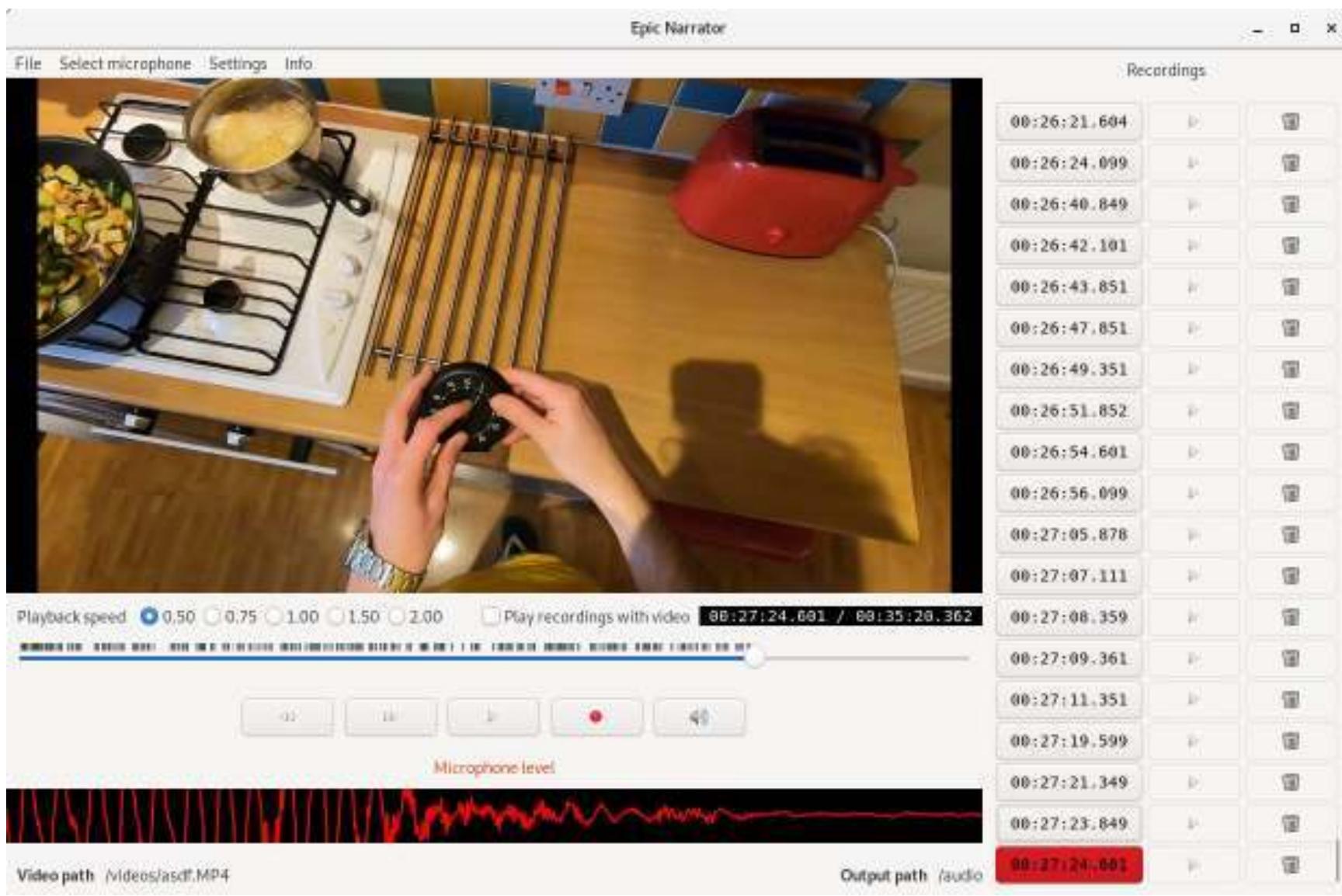
# EPIC-KITCHENS



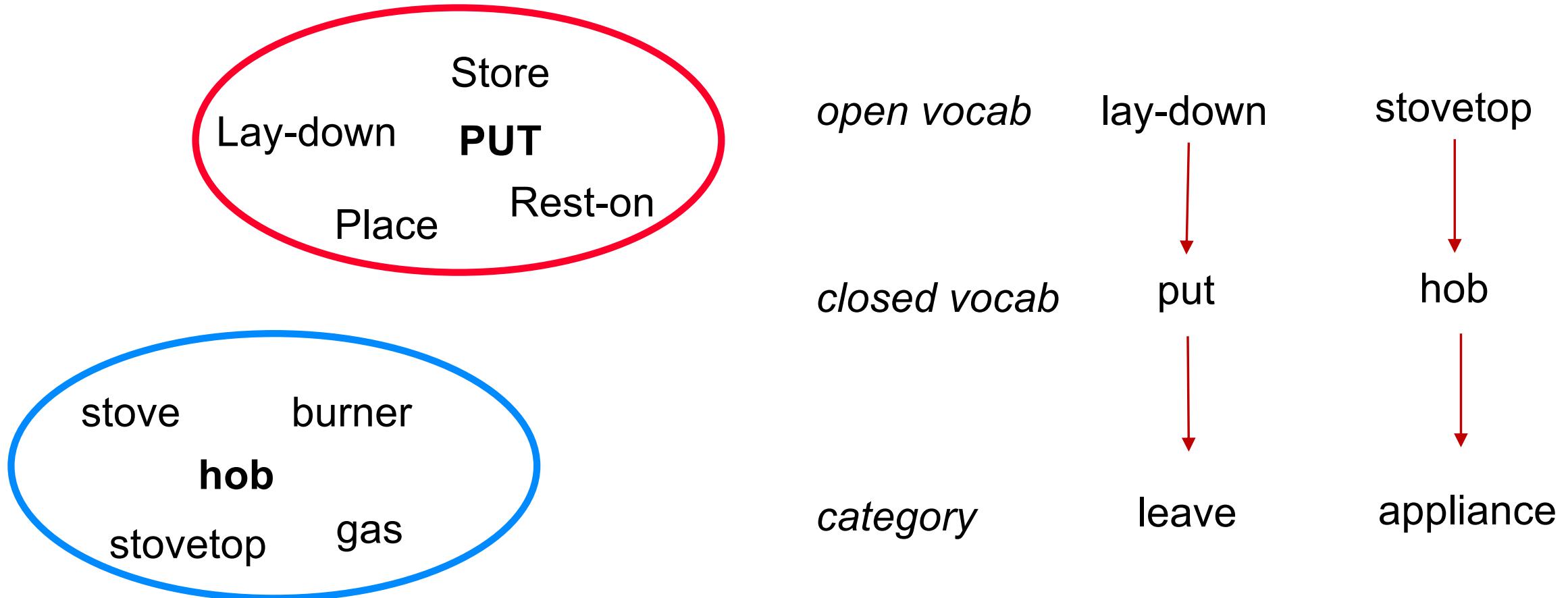
EPIC  
KITCHENS



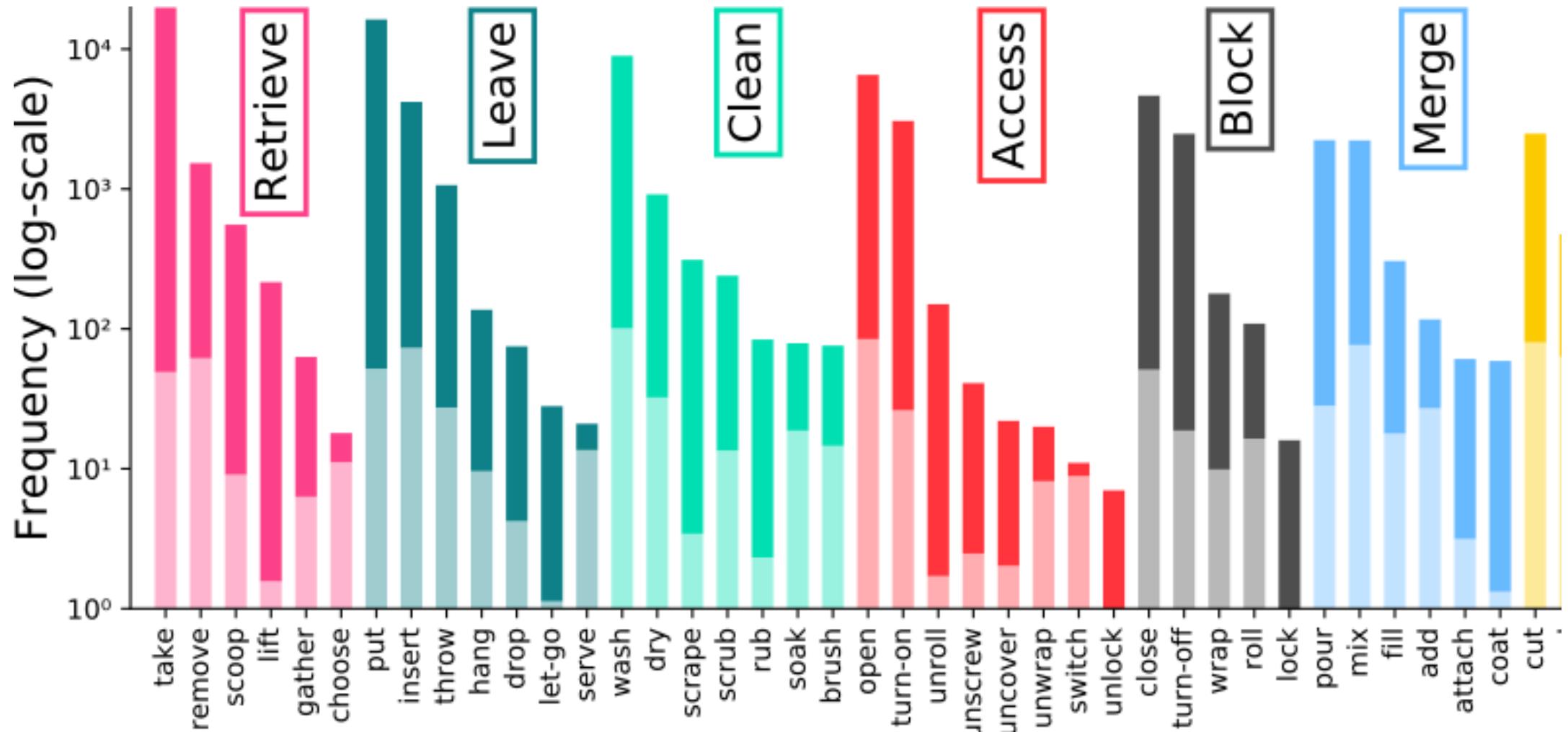
# EPIC-KITCHENS



# EPIC-KITCHENS and Ego4D



# EPIC-KITCHENS-100 Statistics



## Narration

C: camera wearer

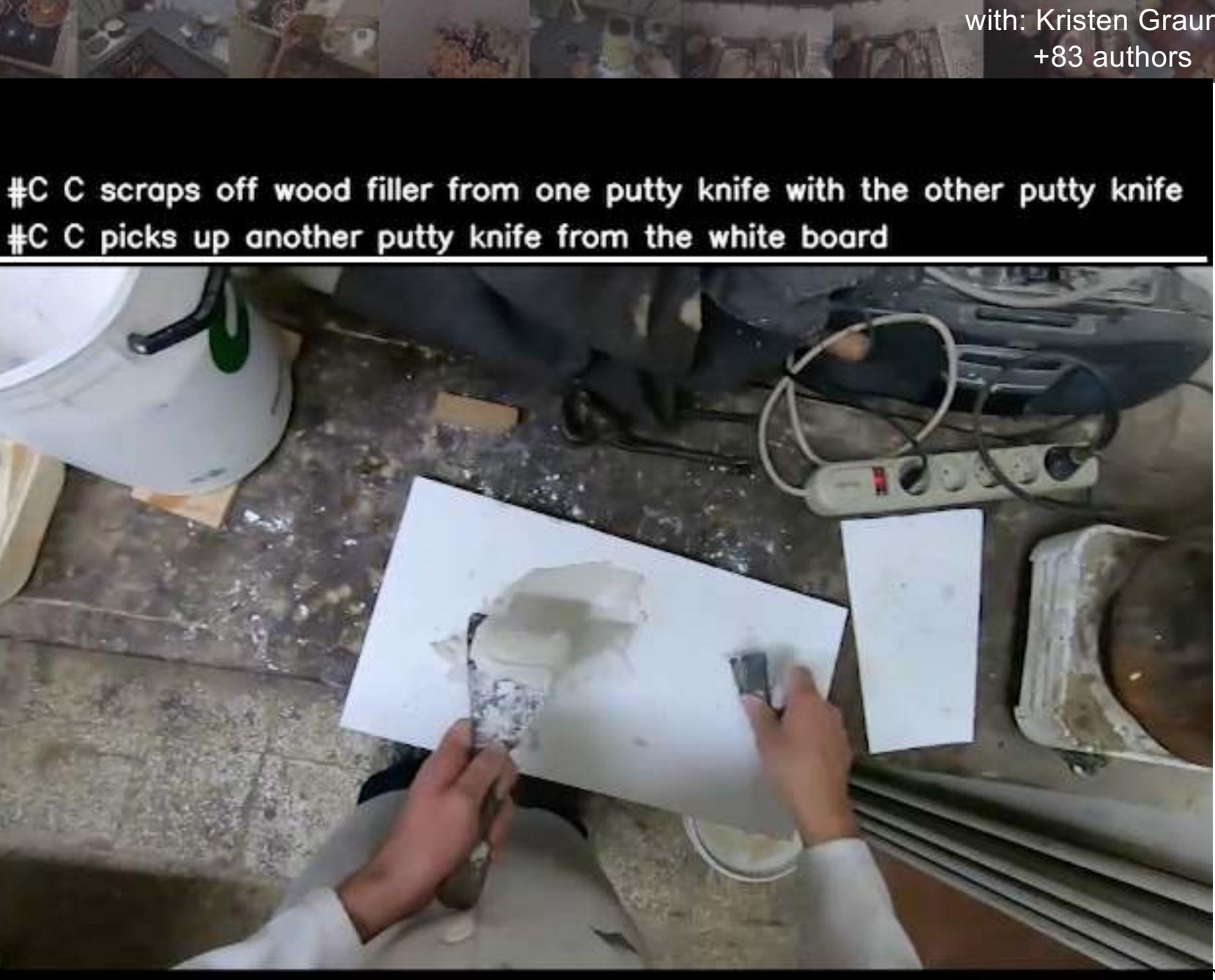
13.2 sentences/min  
3.8 M sentences

1,772 verbs

remove place open  
adjust pull push  
**put pick move**  
take hit turn fix  
carry fold cut lift  
clean up cover eat  
hold drop

4,336 nouns

bowl spoon knife  
card cloth  
bottle plank  
**hand** brush  
paper tray  
container screw cover



# Data Collection Exercise



2017 - now

100 hours  
45 kitchens  
4 countries  
Long-term recording  
Kitchen-based activities



2020 - now

13 universities  
3670 hours  
923 participants  
74 locations  
9 countries  
Short-term recording  
All daily activities

# Data Collection Exercise



2024 - now

1286 hours  
740 camera wearers  
Skilled activities



2025 - now

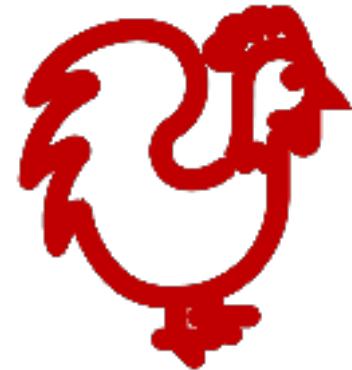
Validation dataset  
41 hours  
9 participants  
Highly-Detailed  
Digital Twin

# Data Collection Exercise



**Labels**

Pascal VOC  
ImageNet  
Kinetics  
Something-Something



**Data**

EPIC-KITCHENS  
Ego4D  
Ego-Exo4D  
HD-EPIC  
...  
KITTI

# The chicken or the egg...

## Data



Naturally unbalanced

Harder to label (exposes ambiguity)

Closer to application

Multiple tasks

## Labels



Unnaturally balanced (or nearly)

Easier to label (hides ambiguity)

Can be expanded

Single task

# In today's tutorial



Motivation and Datasets in  
Egocentric Video Understanding



Video Understanding  
Out of the Frame



Video Understanding:  
Data and Tasks



Teaser: The Wizard of Oz  
at the Sphere



Videos are Multimodal



Outlook into the Future of  
Egocentric Vision



Connected Videos of One's Life



Conclusion

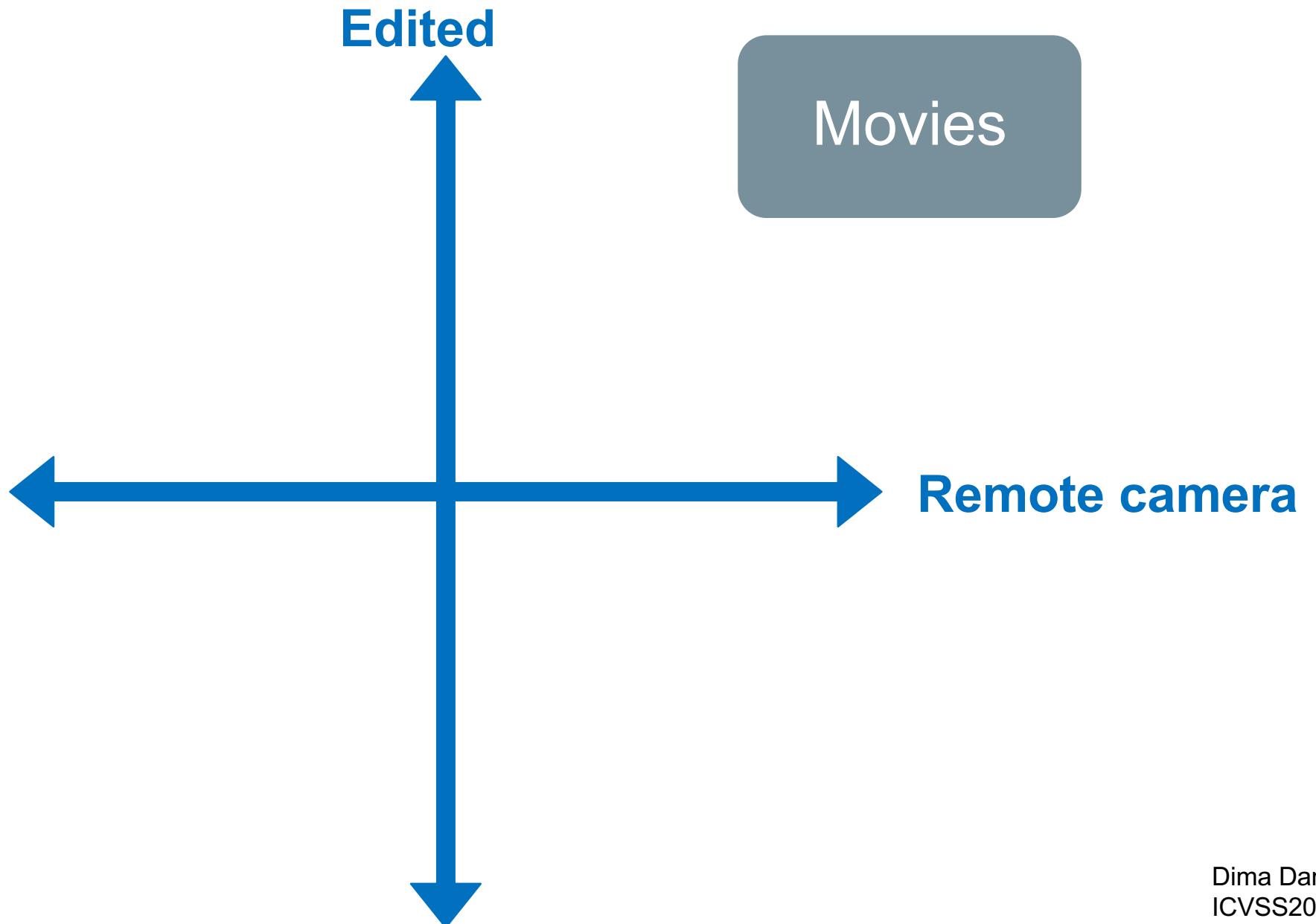
# The history of VIDEO



# The history of VIDEO



# The history of **VIDEO** understanding



# The history of **VIDEO** understanding



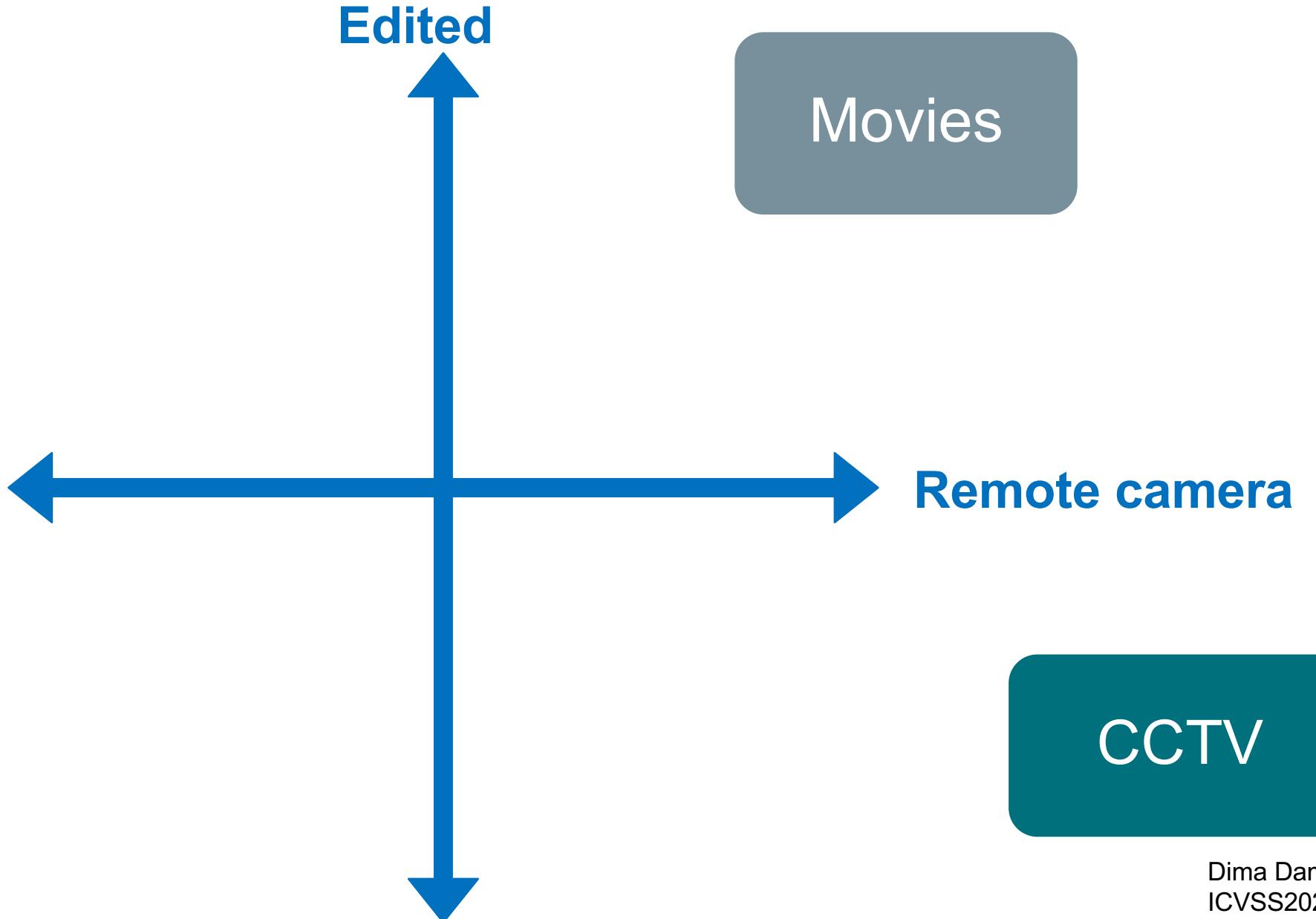
Figure 1. Examples of two action classes (drinking and smoking) from the movie “Coffee and Cigarettes”. Note the high within-

Laptev and Perez (2007)

# The history of **VIDEO** understanding



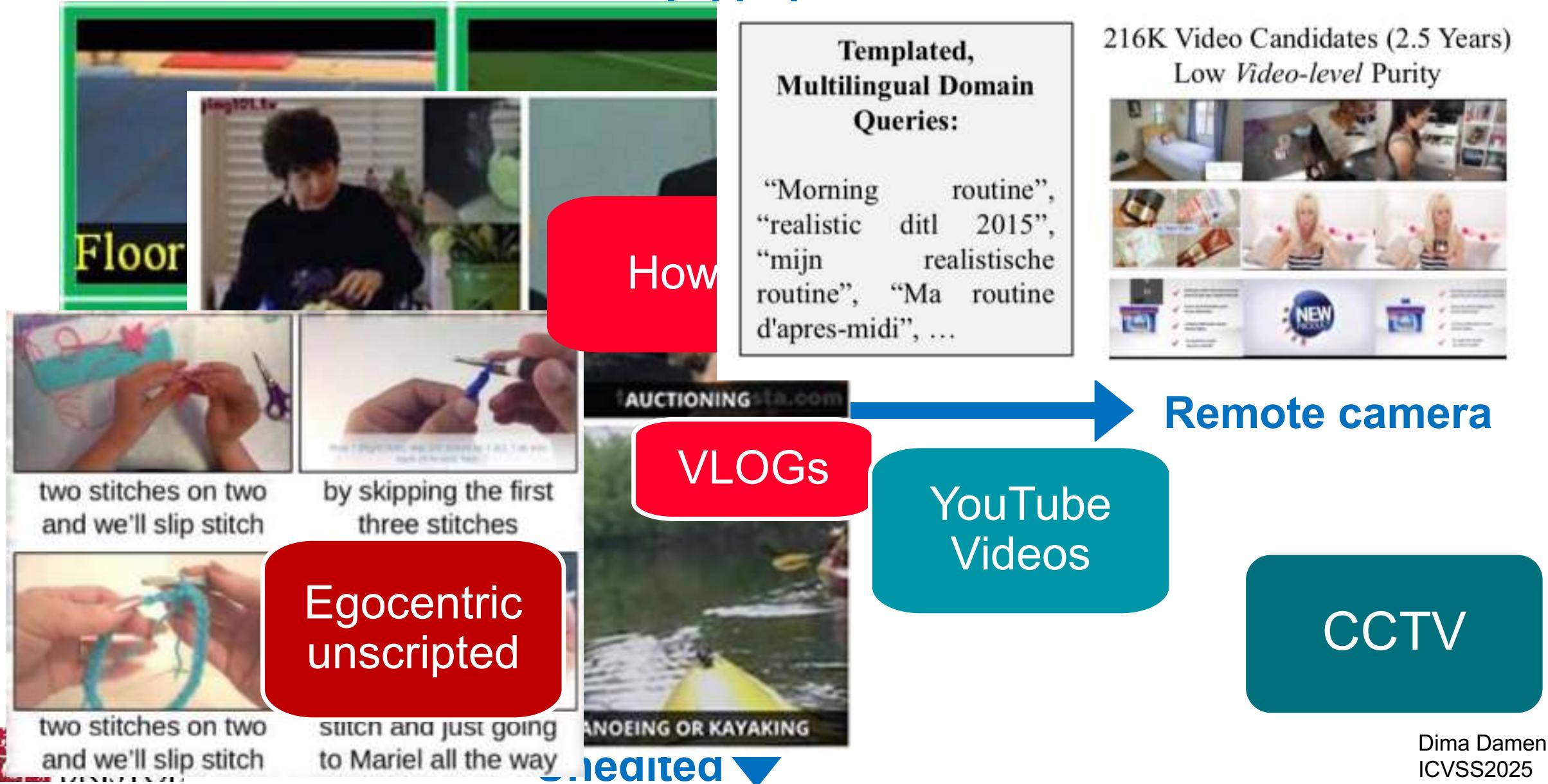
# The history of **VIDEO** Understanding



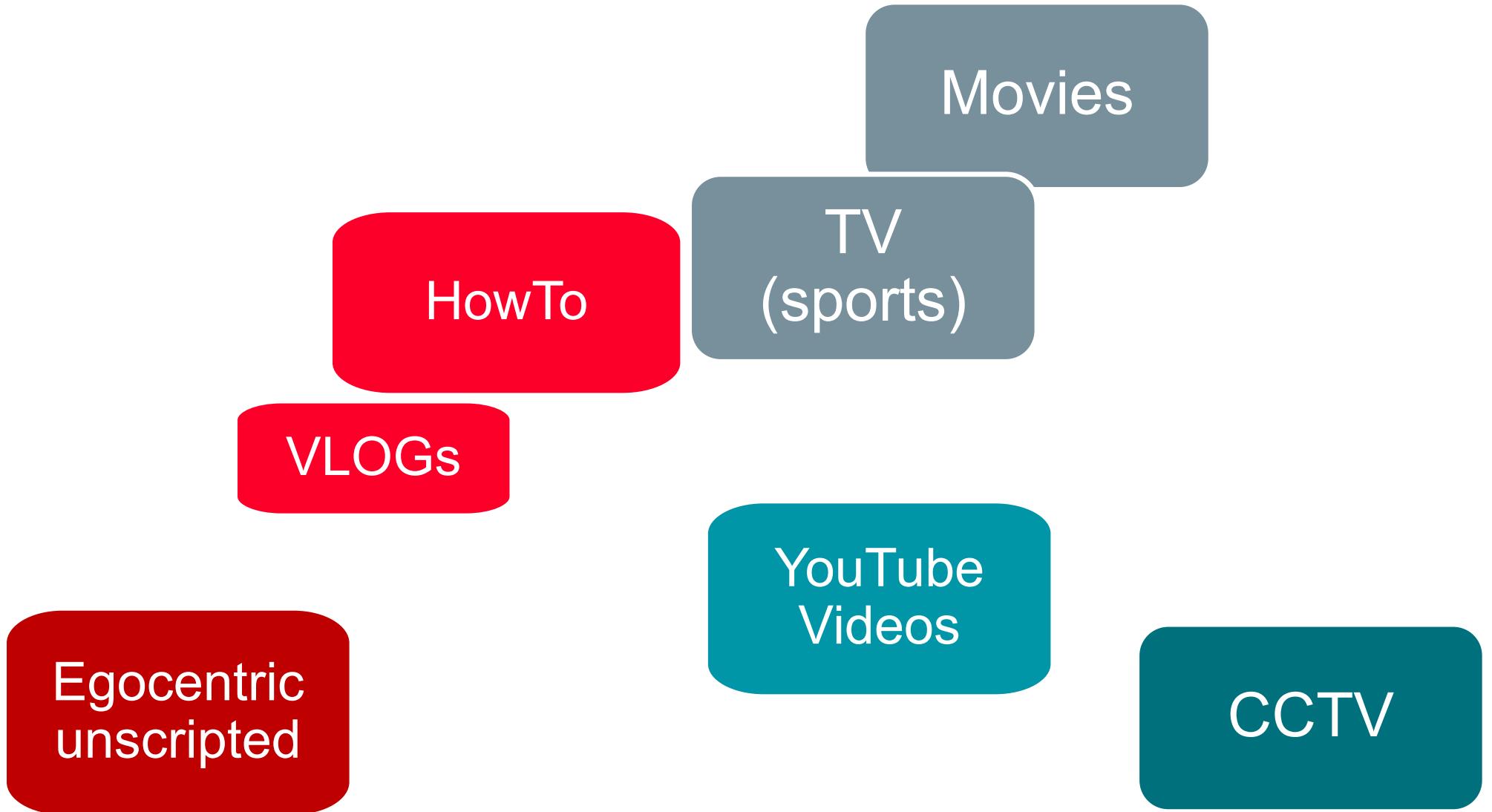
# The history of **VIDEO** understanding



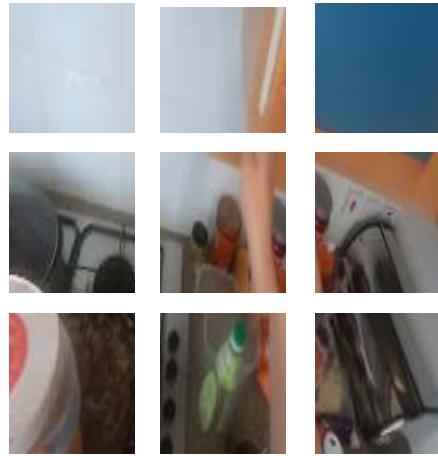
# The history of **VIDEO** understanding



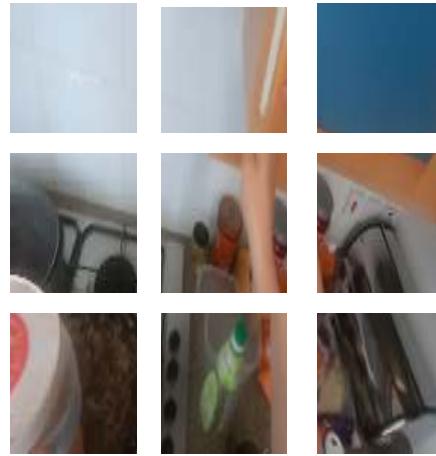
# The history of **VIDEO** understanding



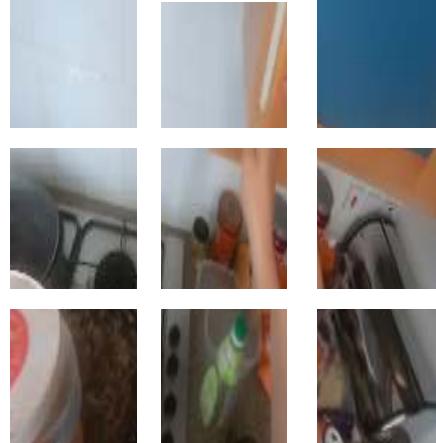
# How to Patch-ify a Video?



$H \times W \times C$

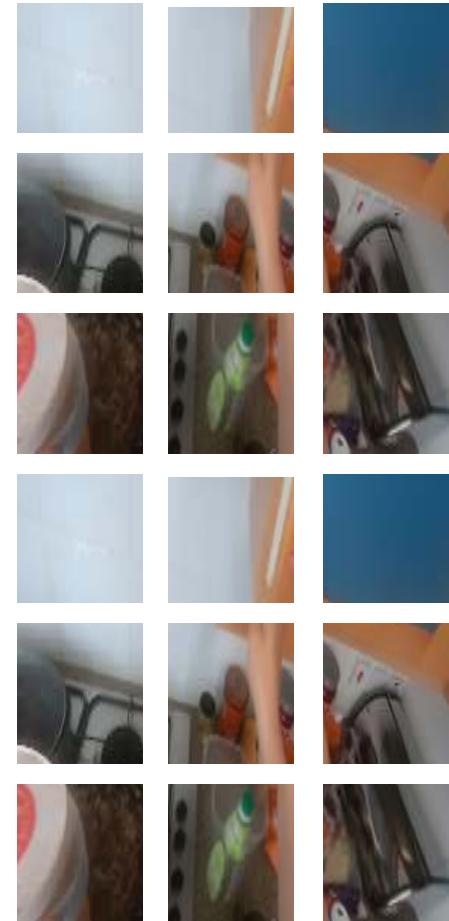


$H \times W \times [CT]$



$[TH] \times W \times C$

# How to Patch-ify a Video?

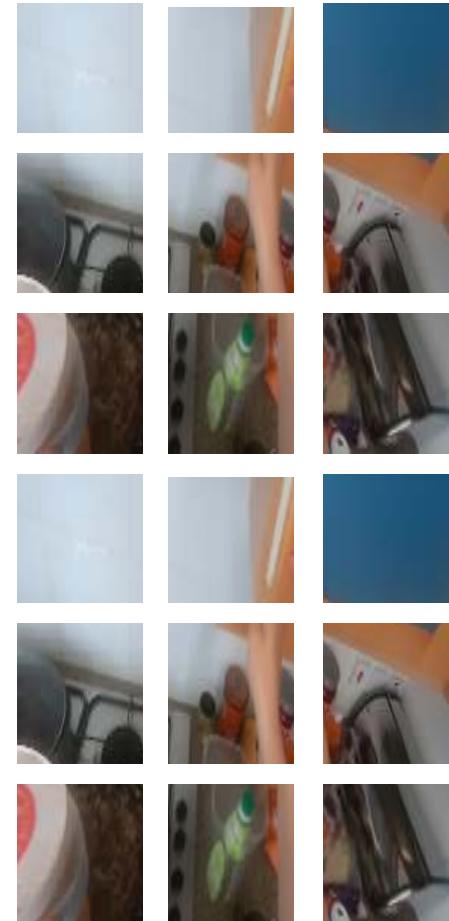


Flatten

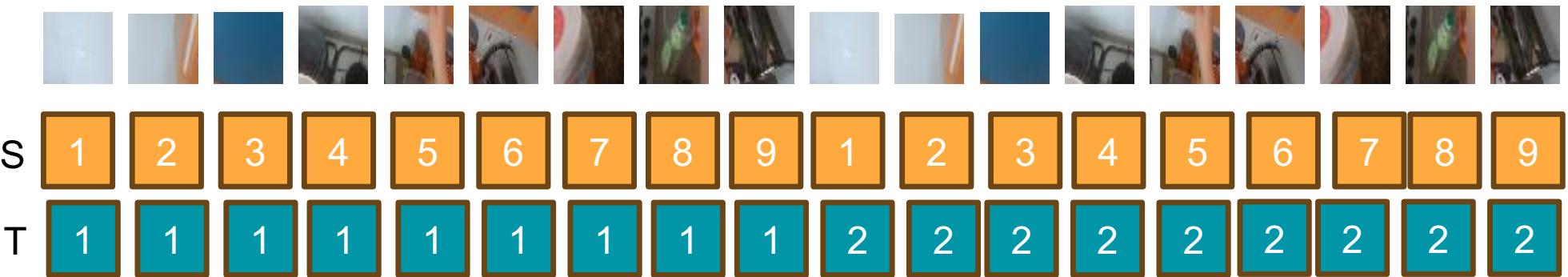


$[T \cdot H] \times W \times C$

# How to Patch-ify a Video?

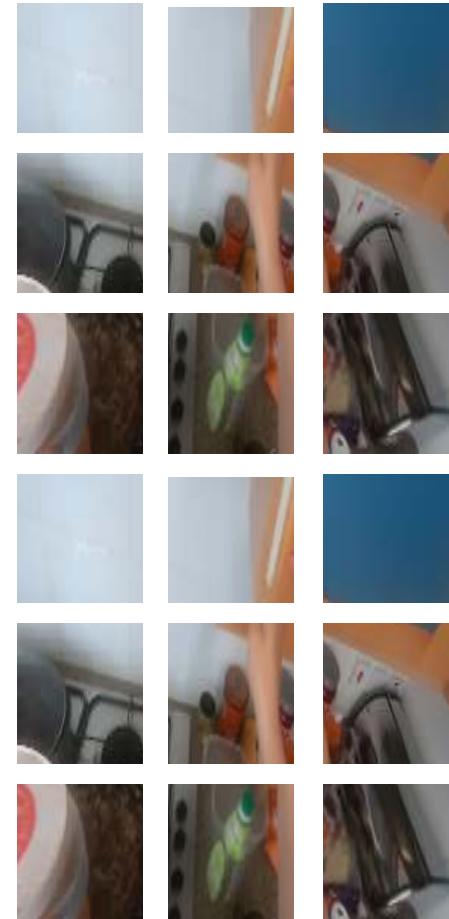


Flatten



$[TH] \times W \times C$

# How to Patch-ify a Video?



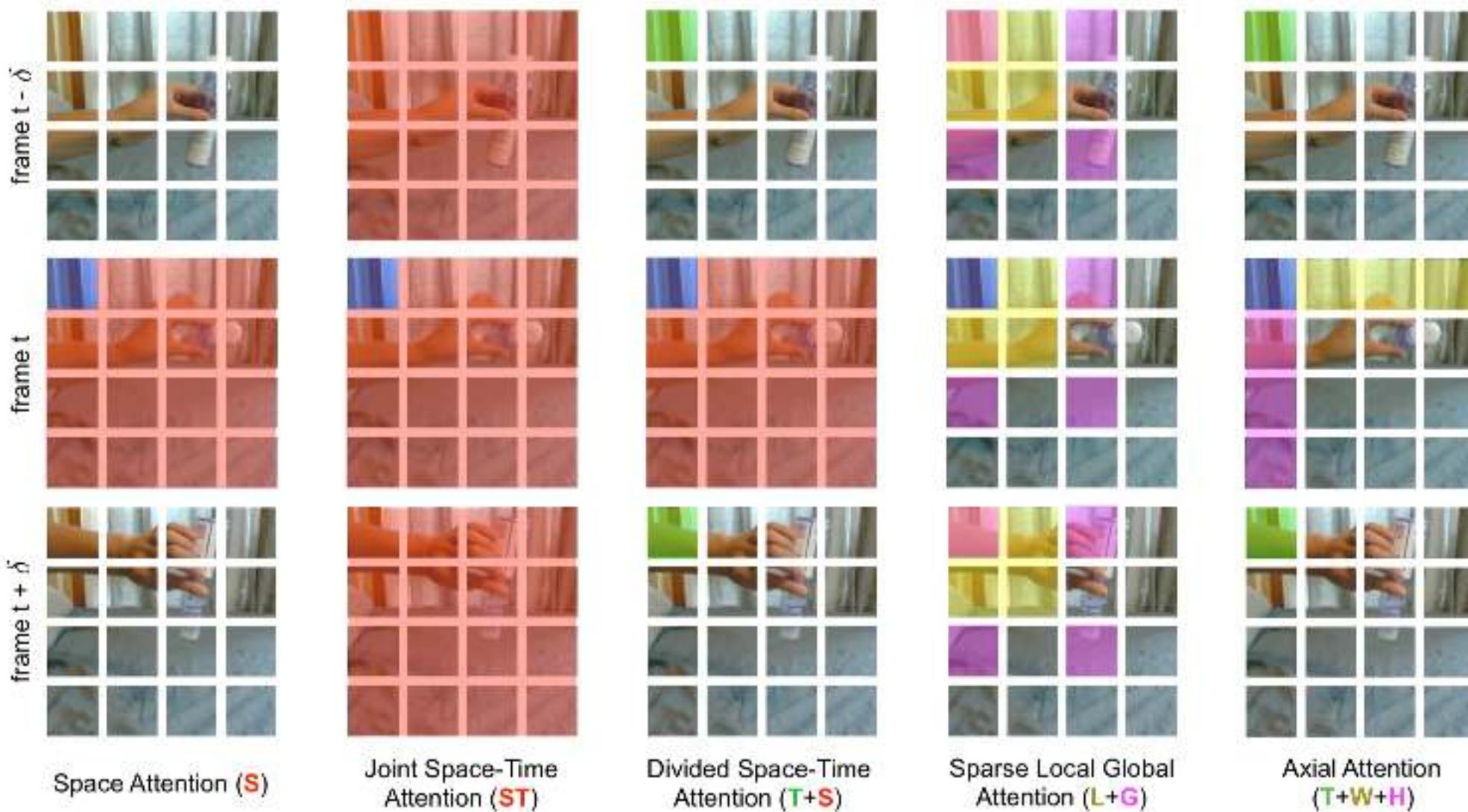
Transformer



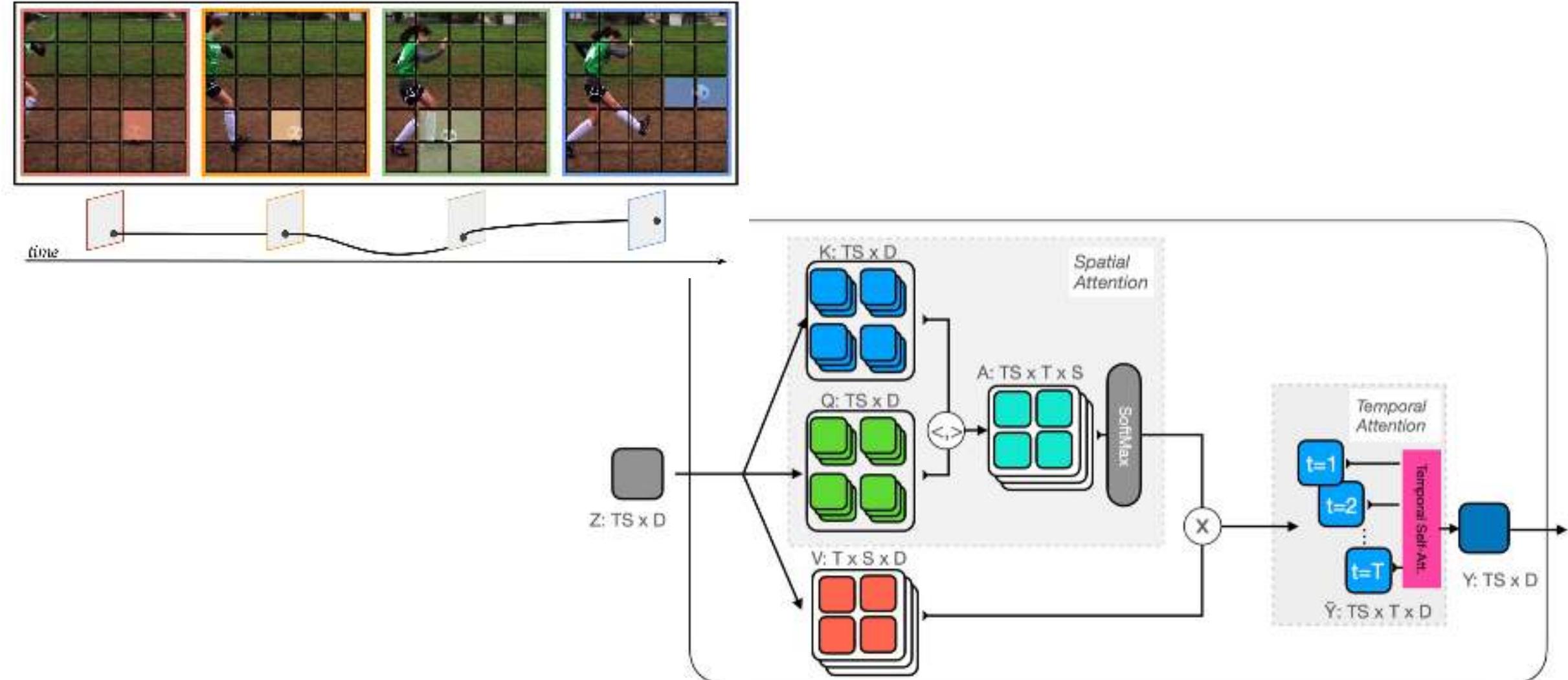
S	1	2	3	4	5	6	7	8	9	1	2	3	4	5	6	7	8	9
T	1	1	1	1	1	1	1	1	1	2	2	2	2	2	2	2	2	2

[TH] x W x C

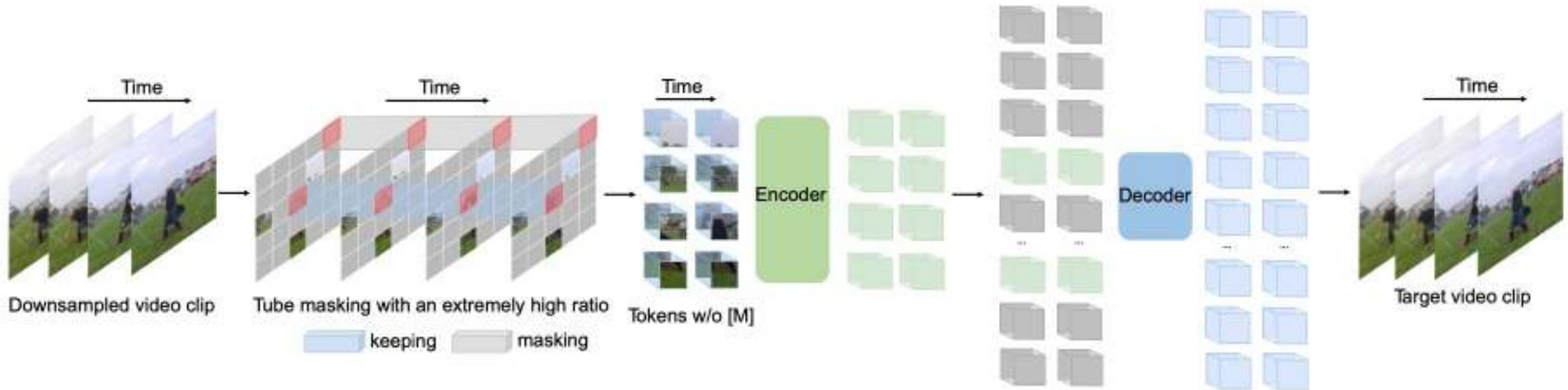
# TimeSFormer



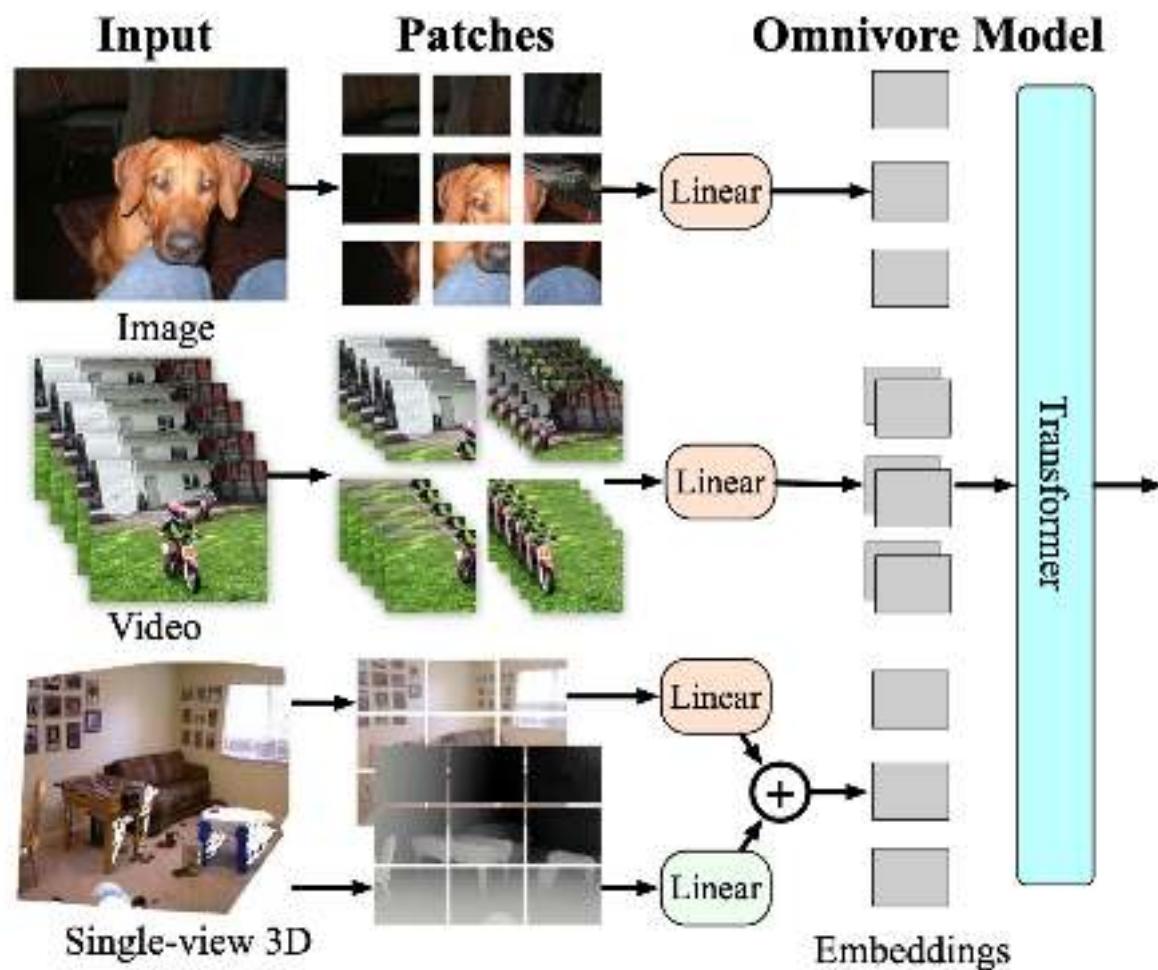
# MotionFormer



# VideoMAE

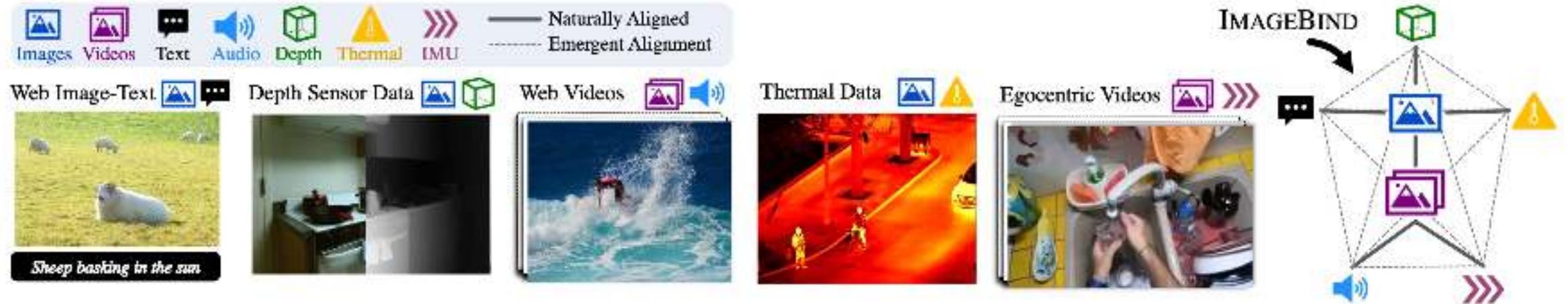


# OmniVore



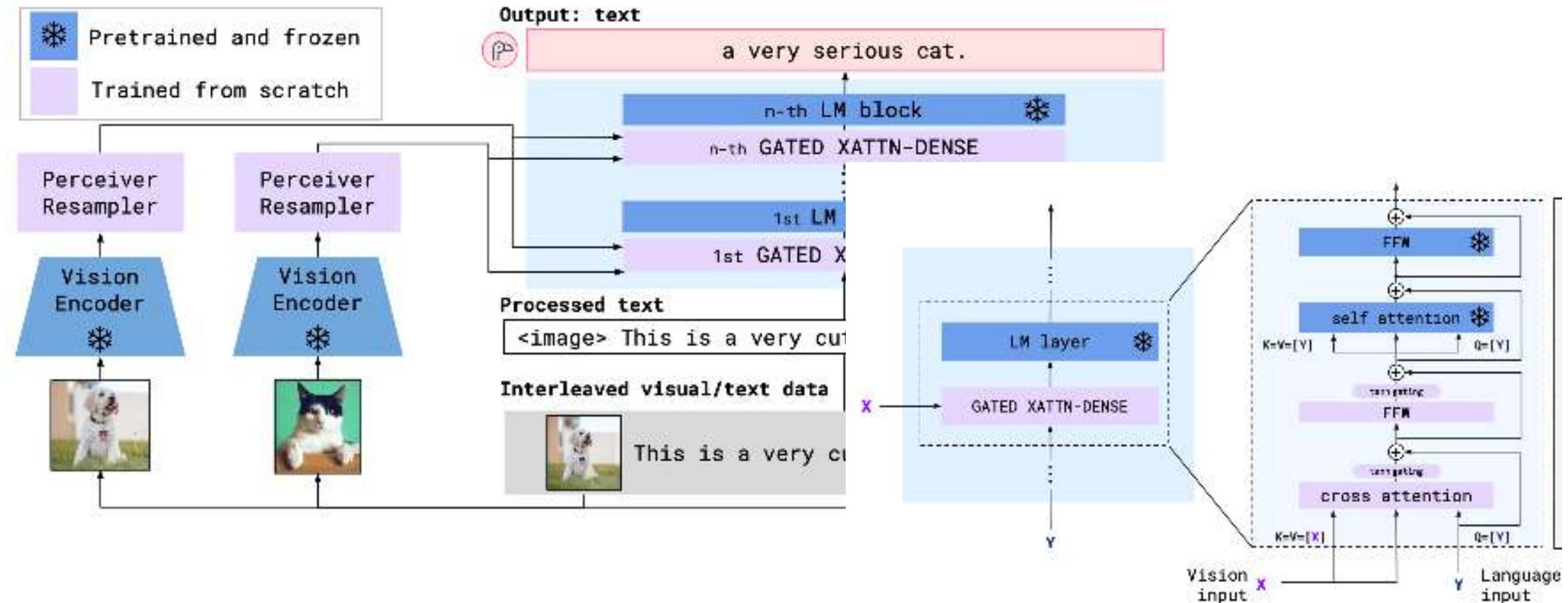
**Figure 2. Multiple visual modalities in the OMNIVORE model.**

# ImageBind



$$L_{\mathcal{I}, \mathcal{M}} = -\log \frac{\exp(\mathbf{q}_i^\top \mathbf{k}_i / \tau)}{\exp(\mathbf{q}_i^\top \mathbf{k}_i / \tau) + \sum_{j \neq i} \exp(\mathbf{q}_i^\top \mathbf{k}_j / \tau)}$$

# Flamingo



# InternVideo

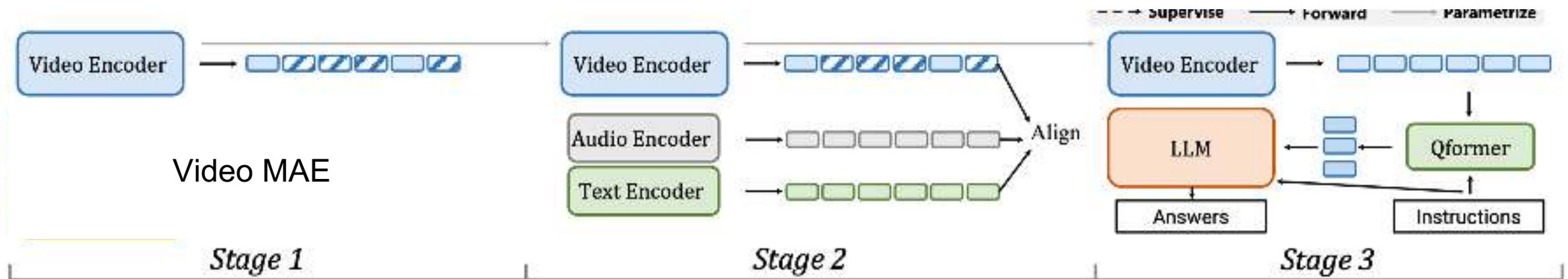
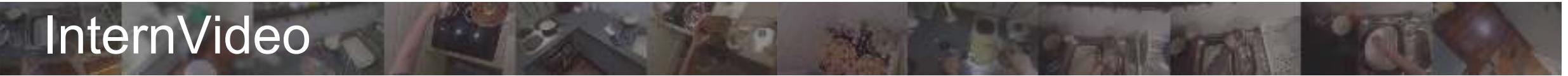


Figure 2: Framework of **InternVideo2**. It consists of three consecutive training phases: unmasked video token reconstruction, multimodal contrastive learning, and next token prediction. In stage 1, the video encoder is trained from scratch, while in stages 2 and 3, it is initialized from the version used in the previous stage.

# InternVideo

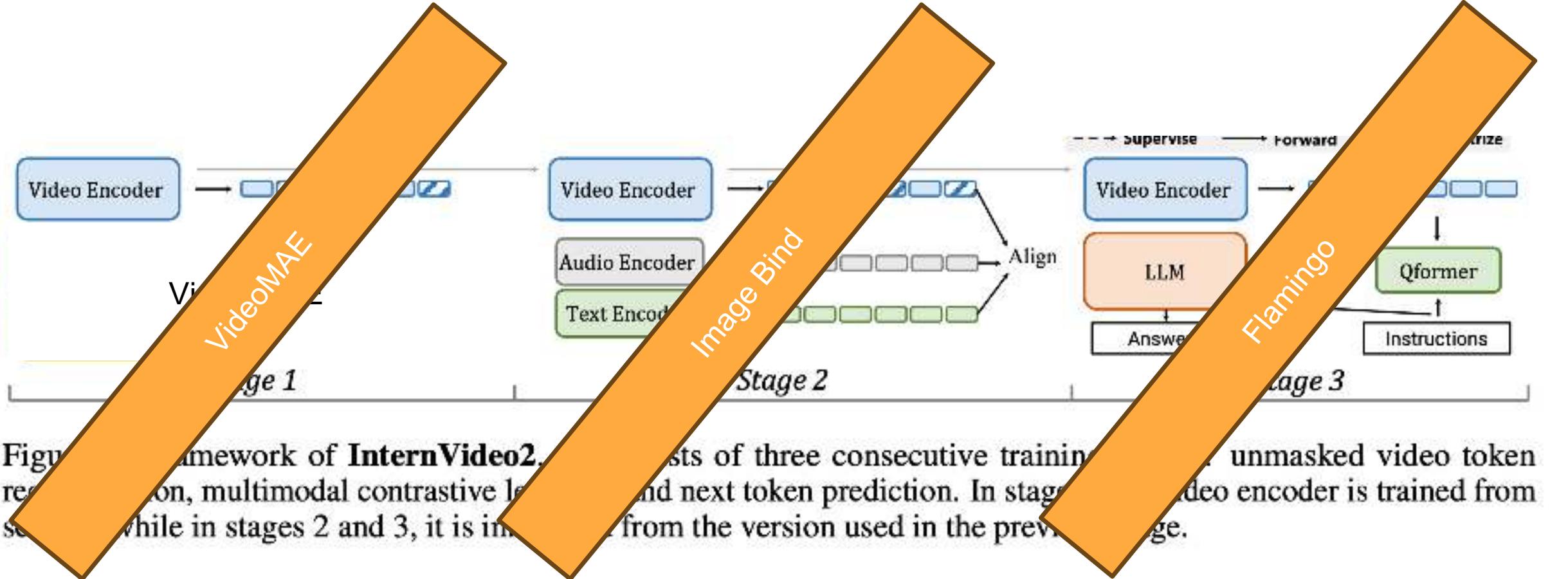


Figure 2: The framework of **InternVideo2**. The framework consists of three consecutive training stages. In stage 1, the unmasked video token prediction, multimodal contrastive learning, and next token prediction. In stage 2, the video encoder is trained from scratch, while in stages 2 and 3, it is initialized from the version used in the previous stage.

# Two types of video understanding tasks

## Video Understanding Tasks

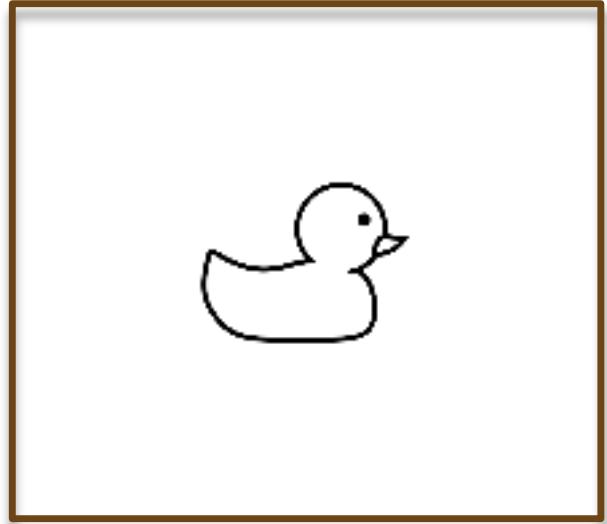
Analogous to Image-based Tasks

Novel Video Tasks

# Analogous Tasks

## Image

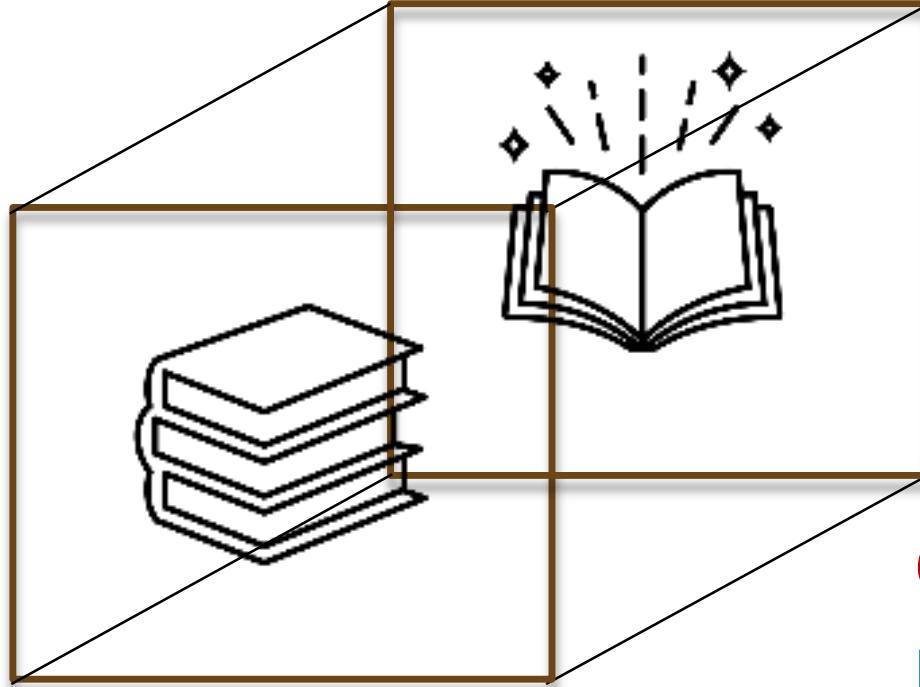
- Object Recognition



Duck

## Video

- Action Recognition

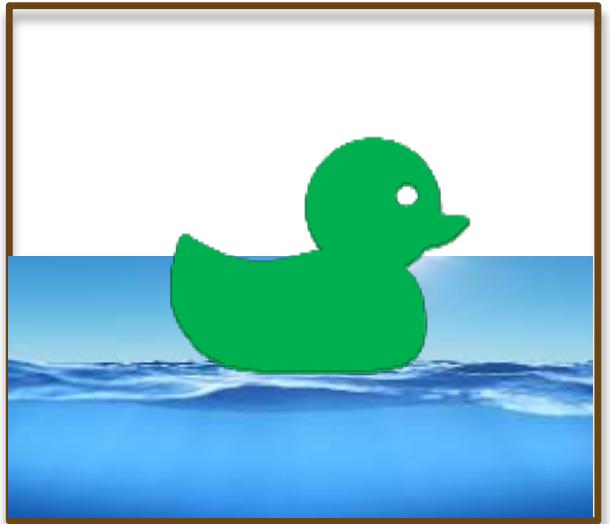


Open  
Book

# Analogous Tasks

## Image

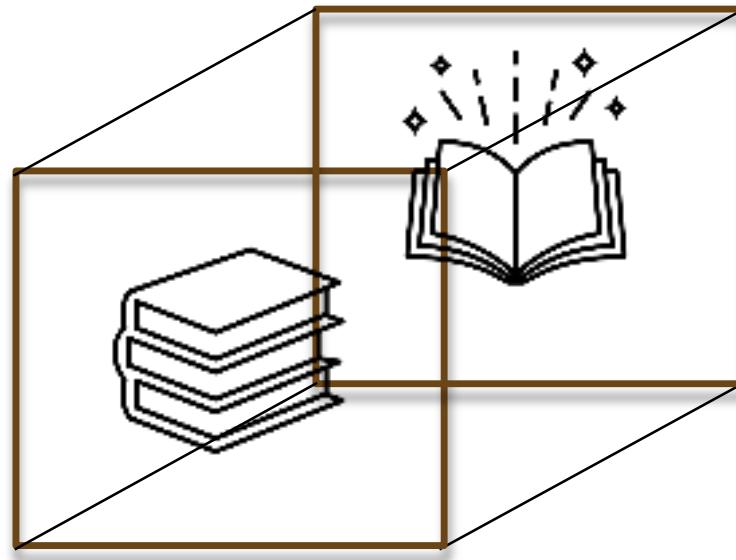
- Image Captioning



A green duck swimming  
In clear water

## Video

- Video Captioning

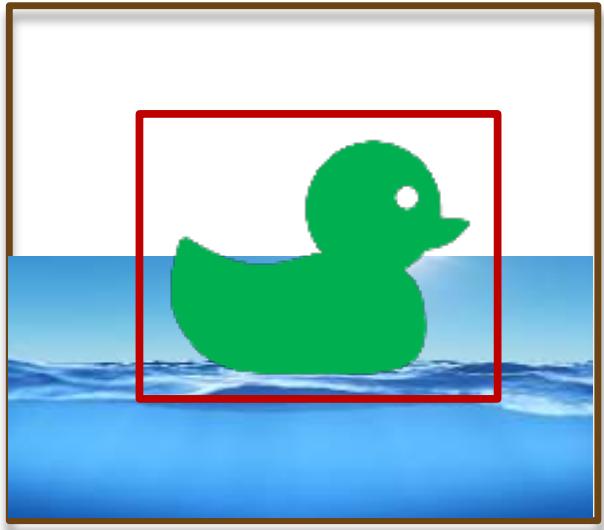


A book picked from top of the pile  
and opened to a page in the middle

# Analogous Tasks

## Image

- Object Detection



Duck

## Video

- Action Detection

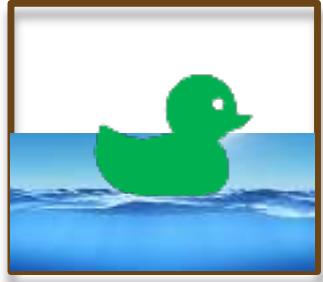


Open Book

# Analogous Tasks

## Image

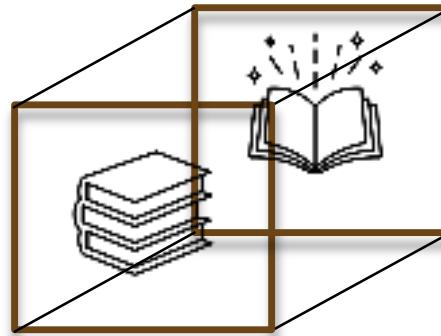
- Image Retrieval



Duck

## Video

- Video Retrieval

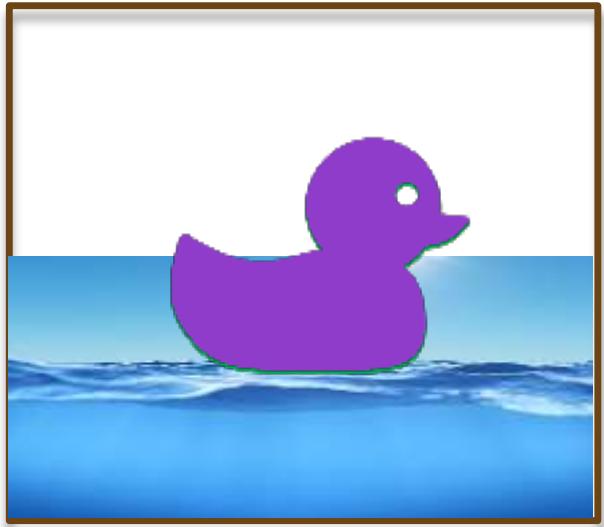


Open Book

# Analogous Tasks

## Image

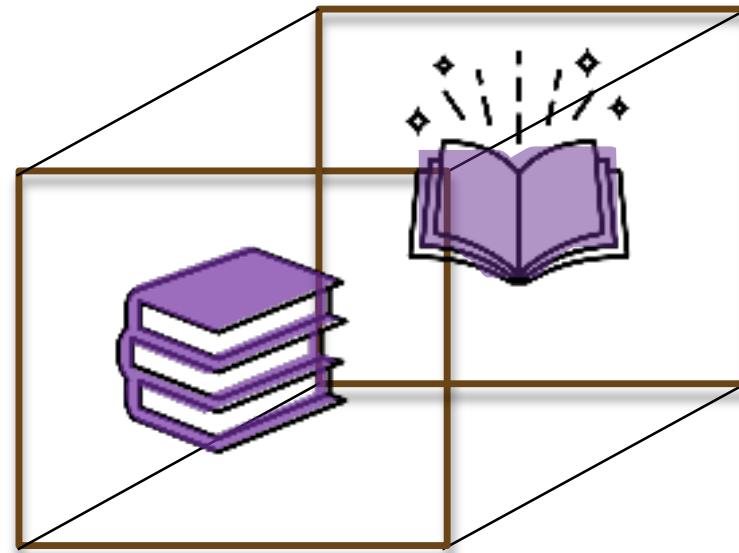
- Object Segmentation



Duck

## Video

- Video Object Segmentation



Book

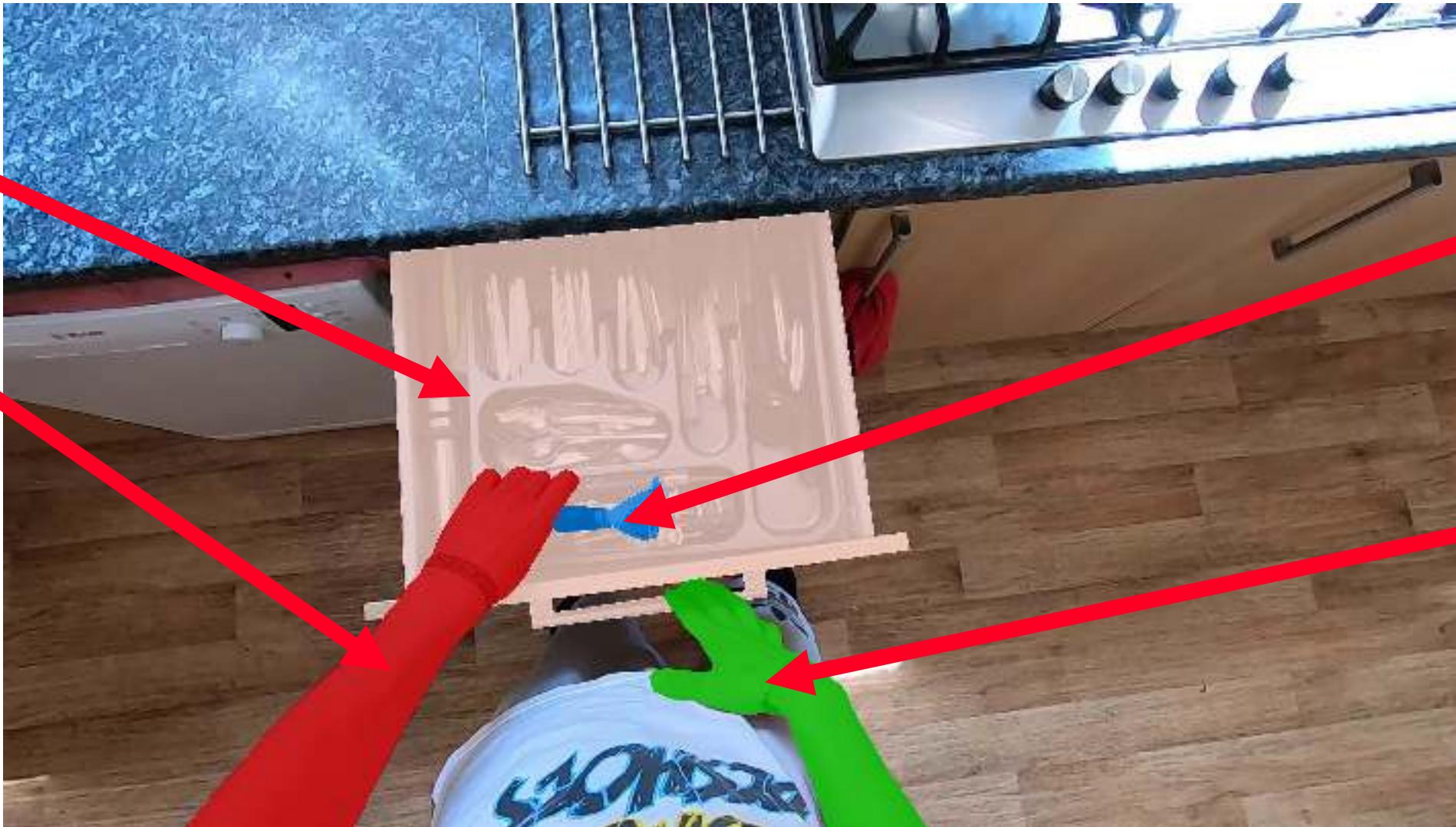
# EPIC-KITCHENS VISOR

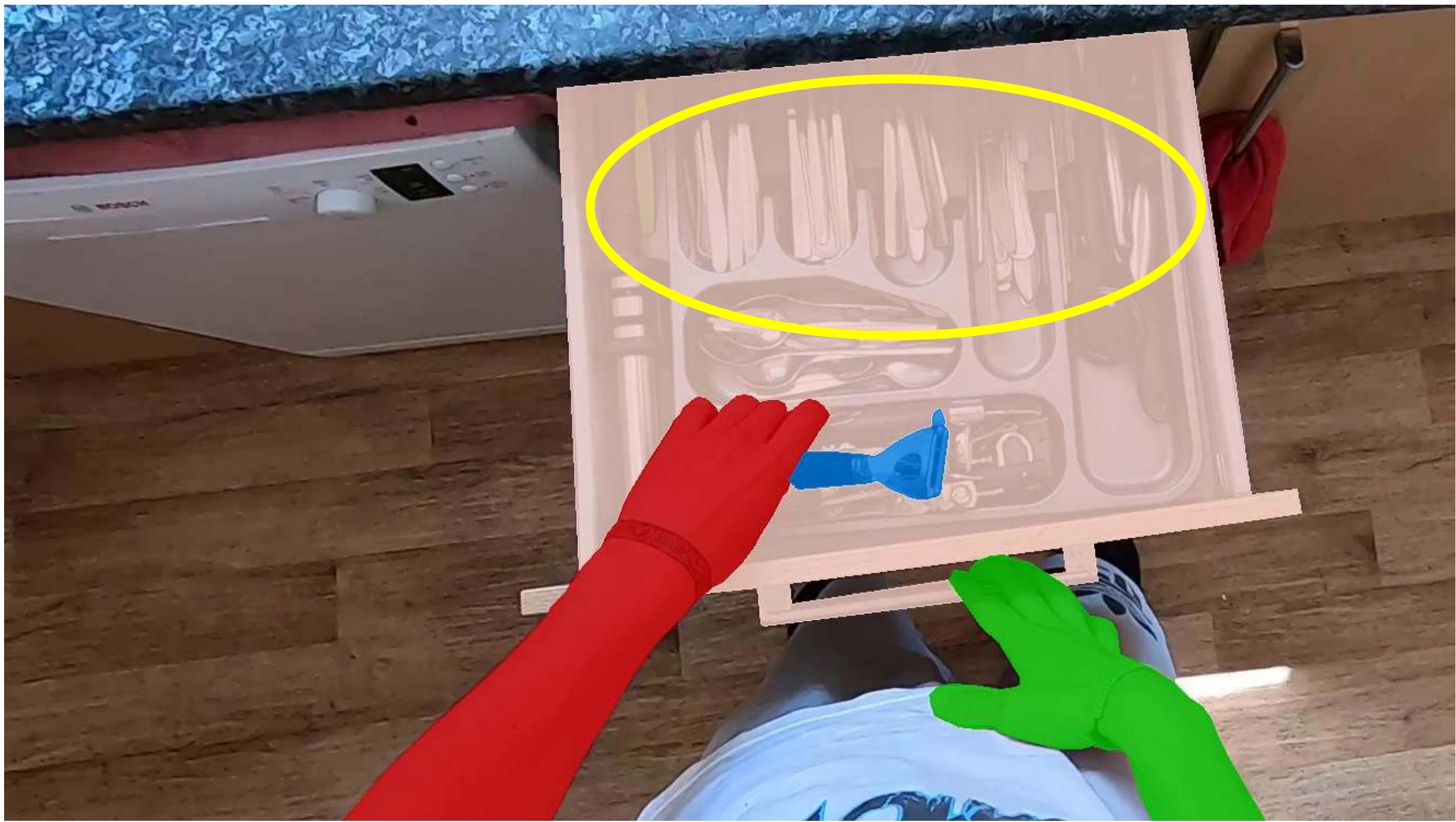
with: Ahmad Darkhalil, Dandan Shan, Bin Zhu, Jian Ma,  
Amlan Kar, Richard Higgins, David Fouhey, Sanja Fidler



# EPIC-KITCHENS VISOR

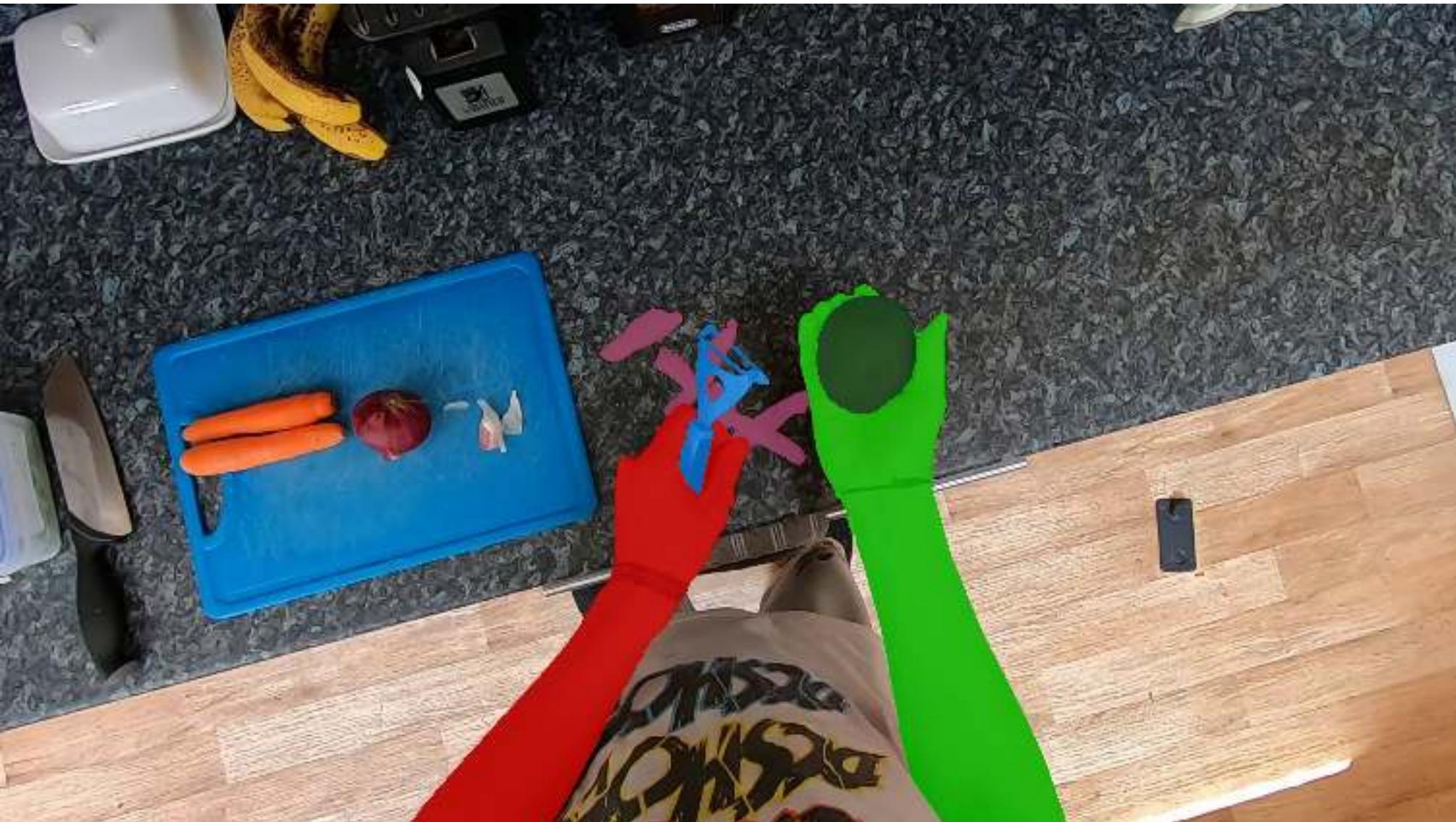
with: Ahmad Darkhalil, Dandan Shan, Bin Zhu, Jian Ma,  
Amlan Kar, Richard Higgins, David Fouhey, Sanja Fidler





# EPIC-KITCHENS VISOR

with: Ahmad Darkhalil, Dandan Shan, Bin Zhu, Jian Ma,  
Amlan Kar, Richard Higgins, David Fouhey, Sanja Fidler





# EPIC-KITCHENS VISOR

with: Ahmad Darkhalil, Dandan Shan, Bin Zhu, Jian Ma,  
Amlan Kar, Richard Higgins, David Fouhey, Sanja Fidler



# EPIC-KITCHENS VISOR

with: Ahmad Darkhalil, Dandan Shan, Bin Zhu, Jian Ma,  
Amlan Kar, Richard Higgins, David Fouhey, Sanja Fidler



# EPIC-KITCHENS VISOR

with: Ahmad Darkhalil, Dandan Shan, Bin Zhu, Jian Ma,  
Amlan Kar, Richard Higgins, David Fouhey, Sanja Fidler



# EPIC-KITCHENS VISOR

with: Ahmad Darkhalil, Dandan Shan, Bin Zhu, Jian Ma,  
Amlan Kar, Richard Higgins, David Fouhey, Sanja Fidler



# EPIC-KITCHENS VISOR

with: Ahmad Darkhalil, Dandan Shan, Bin Zhu, Jian Ma,  
Amlan Kar, Richard Higgins, David Fouhey, Sanja Fidler



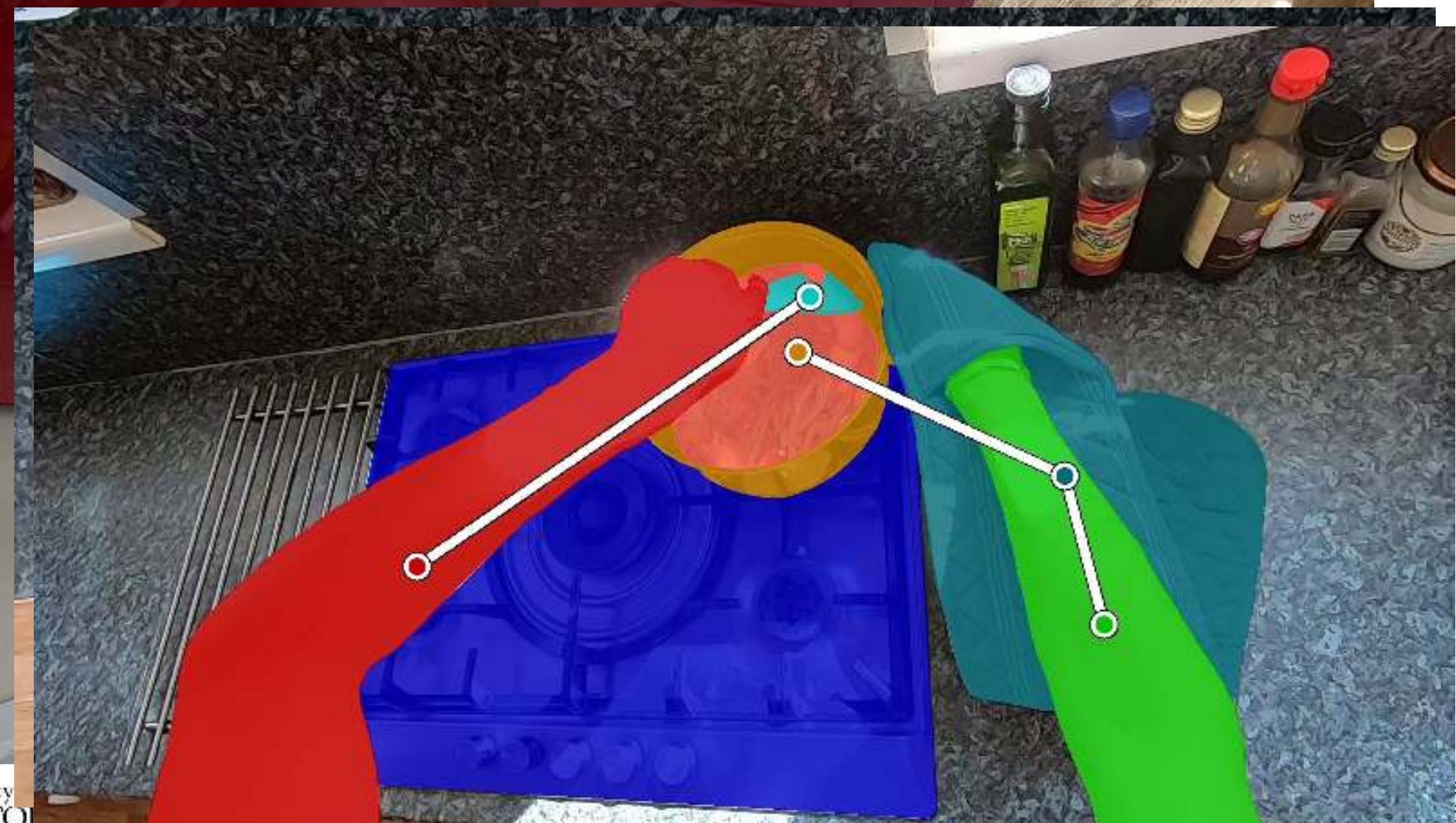
# EPIC-KITCHENS VISOR

with: Ahmad Darkhalil, Dandan Shan, Bin Zhu, Jian Ma,  
Amlan Kar, Richard Higgins, David Fouhey, Sanja Fidler



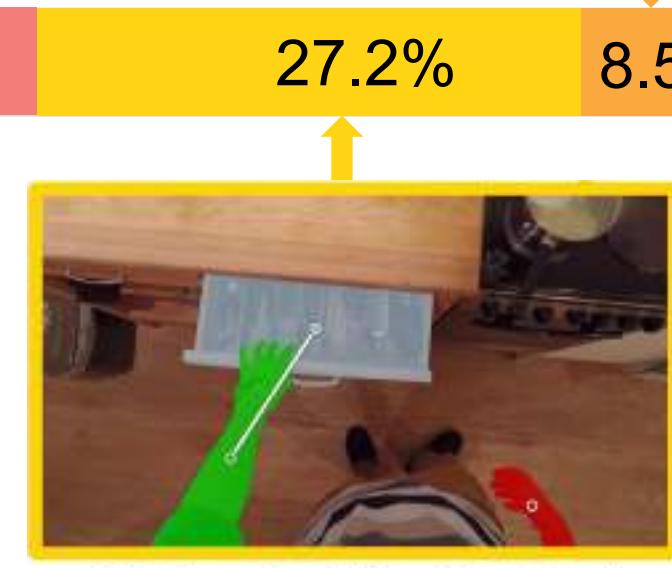
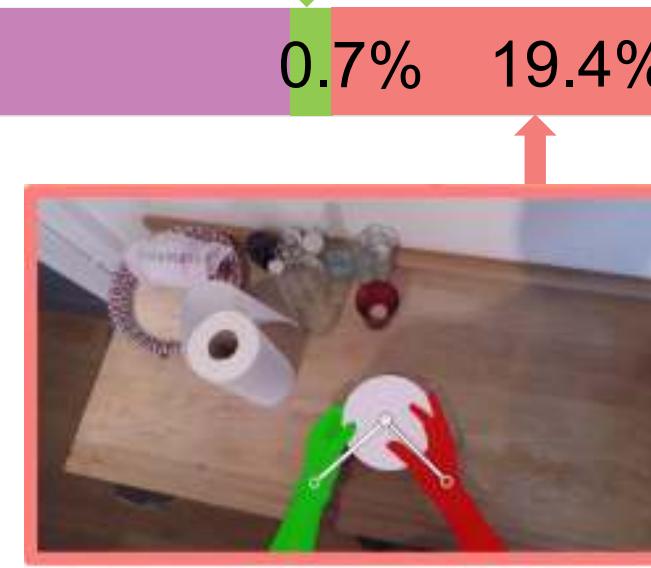
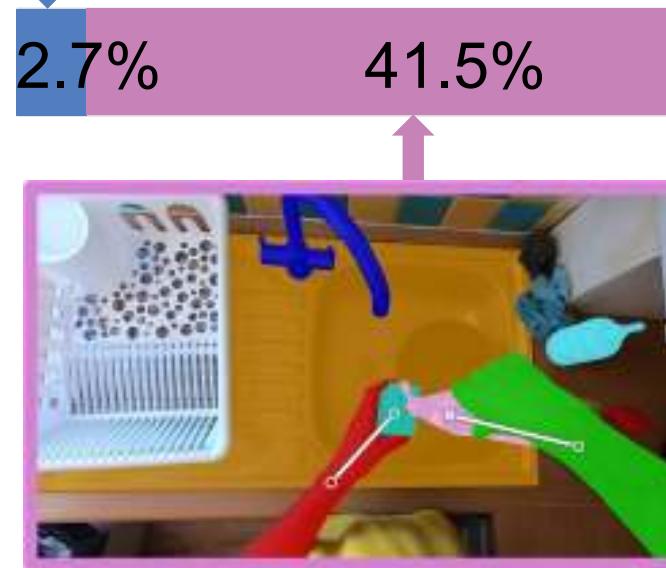
# VISOR Relations

with: Ahmad Darkhalil, Dandan Shan, Bin Zhu, Jian Ma, Amlan Kar,  
Richard Higgins, David Fouhey, Sanja Fidler, Dima Damen



# Object relation stats

with: Ahmad Darkhalil, Dandan Shan, Bin Zhu, Jian Ma, Amlan Kar,  
Richard Higgins, David Fouhey, Sanja Fidler, Dima Damen



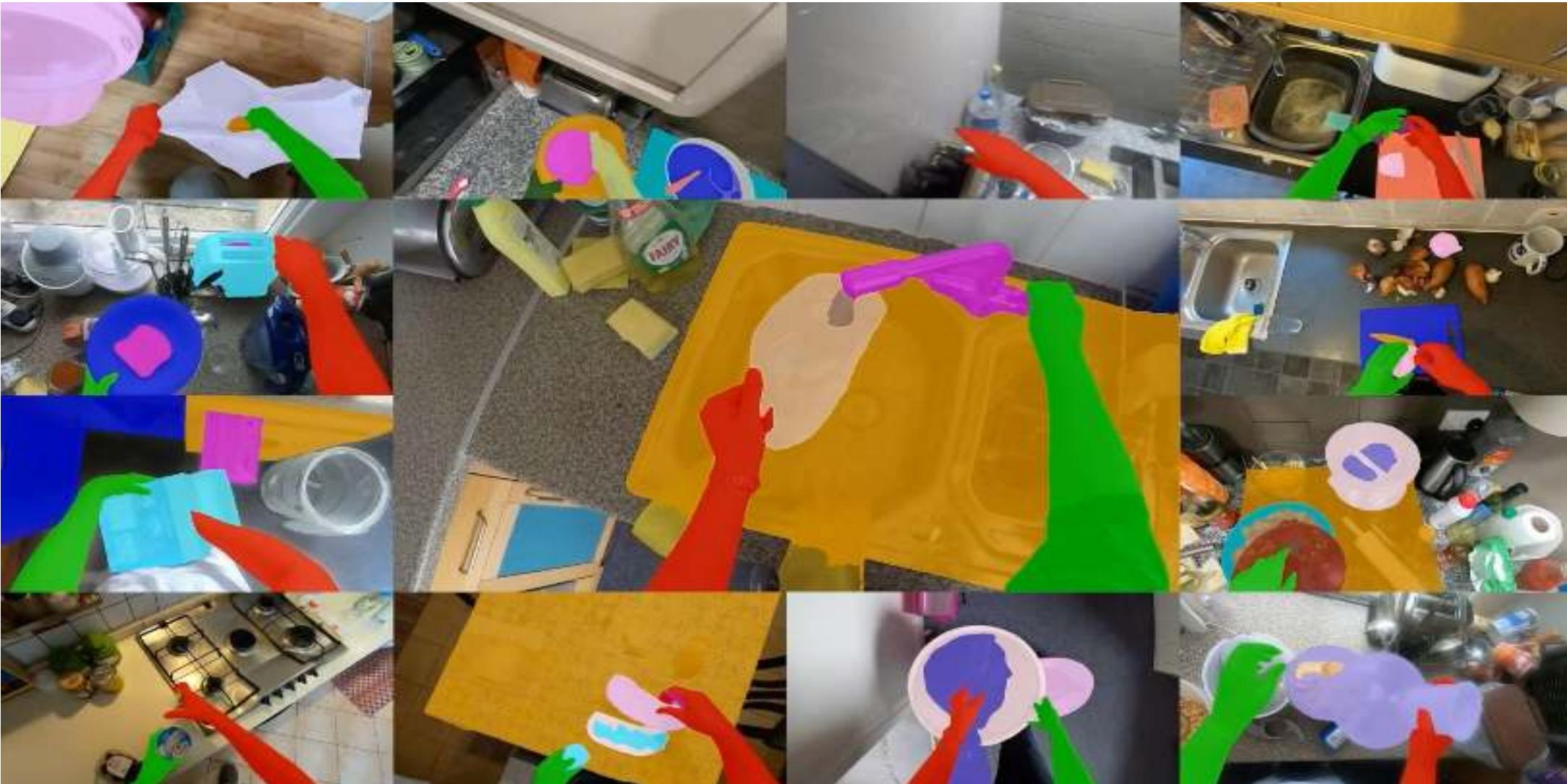
2.7% 41.5%

0.7% 19.4%

27.2% 8.5%

# EPIC-KITCHENS VISOR

with: Ahmad Darkhalil, Dandan Shan, Bin Zhu, Jian Ma,  
Amlan Kar, Richard Higgins, David Fouhey, Sanja Fidler



# Analogous Tasks

## Image

- Object Counting

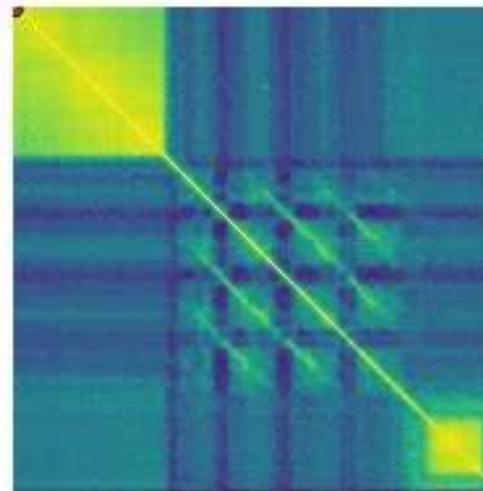
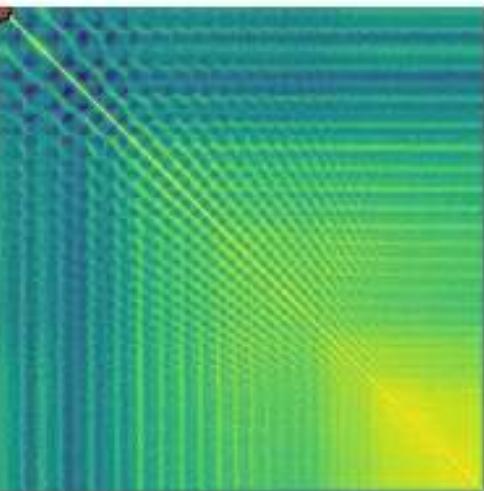
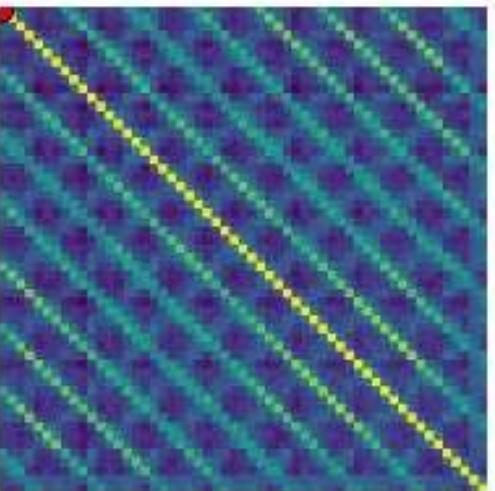
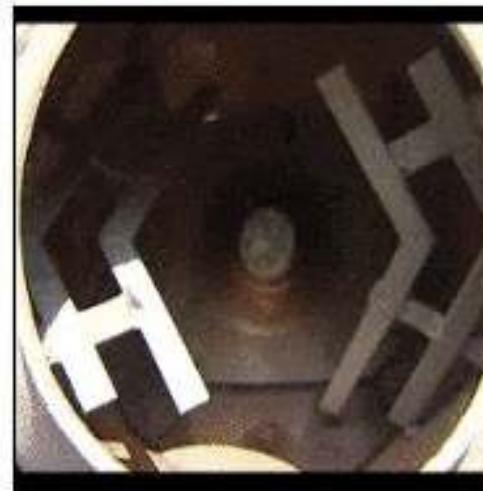


## Video

- Action Counting



# Countix



# Every Shot Counts

with: Saptarshi Sinha  
Alexandros Stergiou



# Analogous Tasks

## Image

- Text-to-image Generation



Stable Diffusion

## Video

- Text-to-Video Generation



SORA

# Text-to-Video Generation



# Text-to-Video Generation



Prompt: A grandmother with neatly combed grey hair stands behind a colorful birthday cake with numerous candles at a wood dining room table, expression is one of pure joy and happiness, with a happy glow in her eye. She leans forward and blows out...



Gemini Veo 3 – May 2025

Dima Damen  
ICVSS2025



Gemini Veo 3 – May 2025

Dima Damen  
ICVSS2025

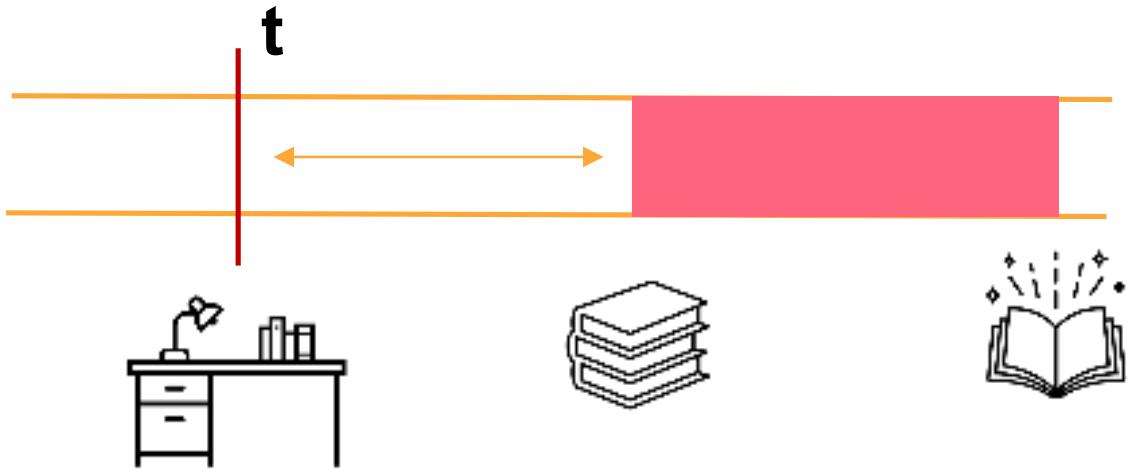
# Non-Analogous Tasks

Image



Video

- Action Anticipation  
What will happen after 1 second?



Open Book

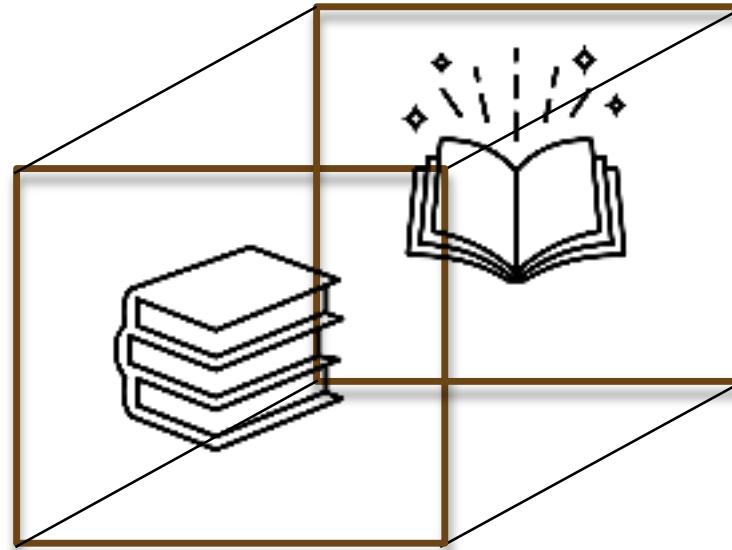
# Non-Analogous Tasks

Image



Video

- Manner of action  
How did you open the book?



# How?

with: Hazel Doughty  
Ivan Laptev  
Walterio Mayol-Cuevas



... if you **turn** the bowl upside down **slowly** they won't come out ...



... mix it well until it is **completely dissolved** ...



... you want to make sure you **fill** it up **partially** ...



... you want to **dice** it **finely**...

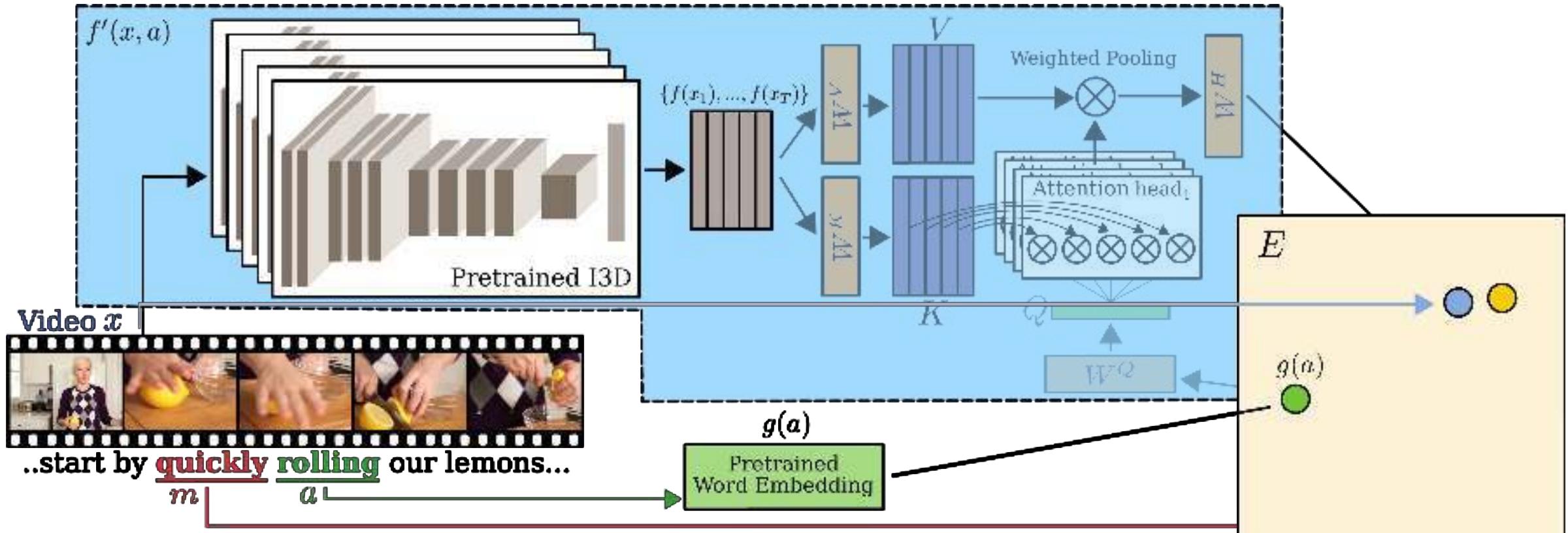
-10 seconds

timestamp

+10 seconds

# How?

with: Hazel Doughty  
Ivan Laptev  
Walterio Mayol-Cuevas



# How?

with: Hazel Doughty  
Ivan Laptev  
Walterio Mayol-Cuevas



... we're going to **mix** these up real **quick**...

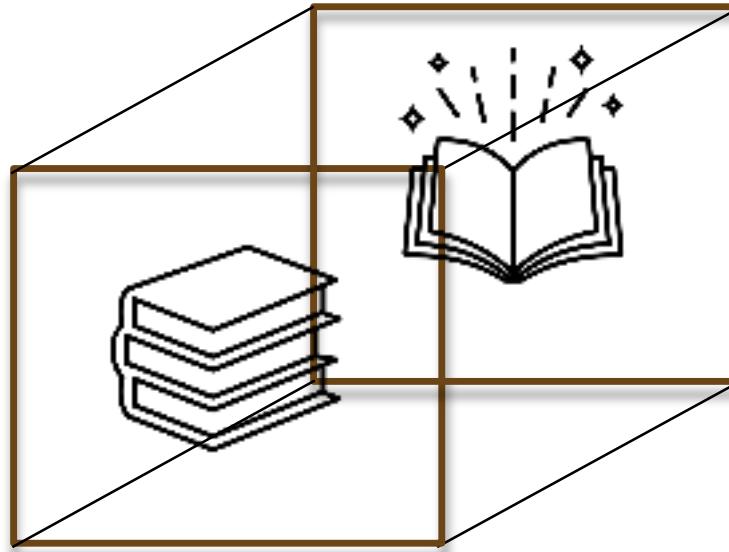
# Non-Analogous Tasks

Image

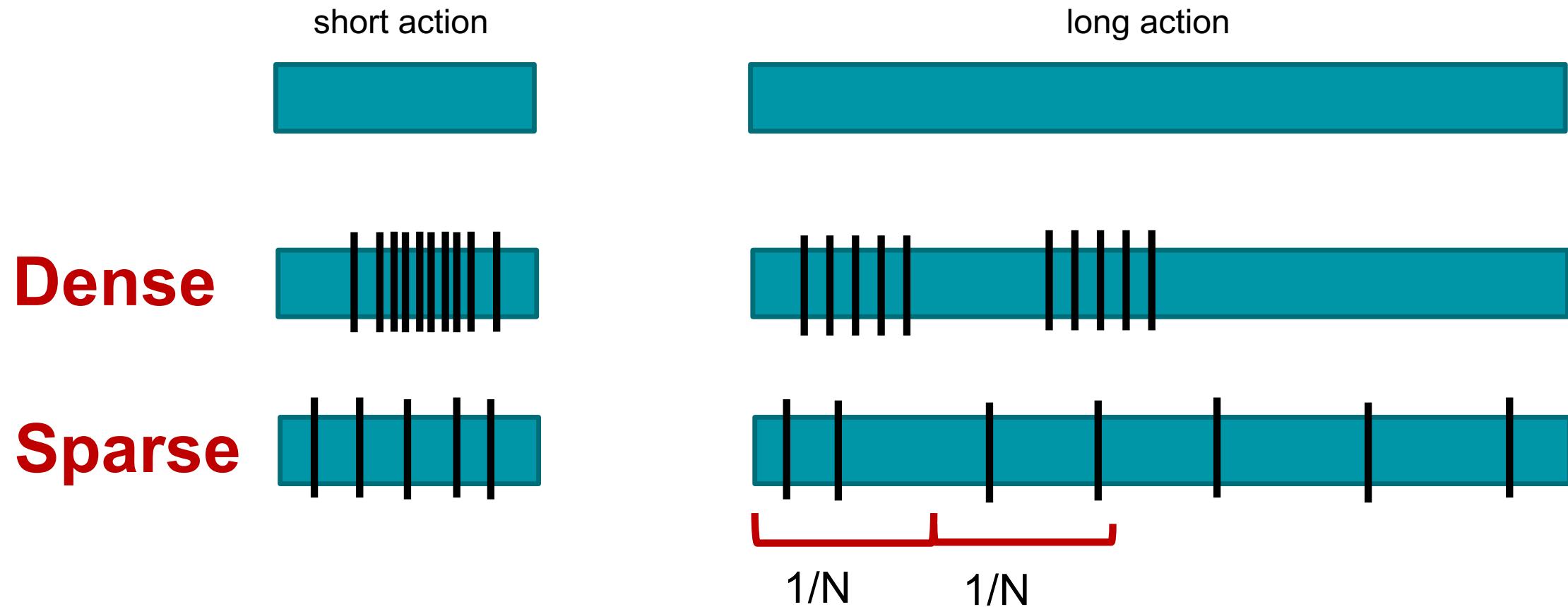


Video

- Frame Selection and Attribution



# Two sampling approaches

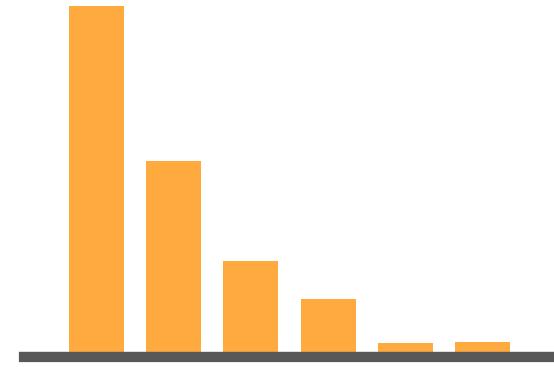


# Frame Attributions in Video Models

with: Will Price



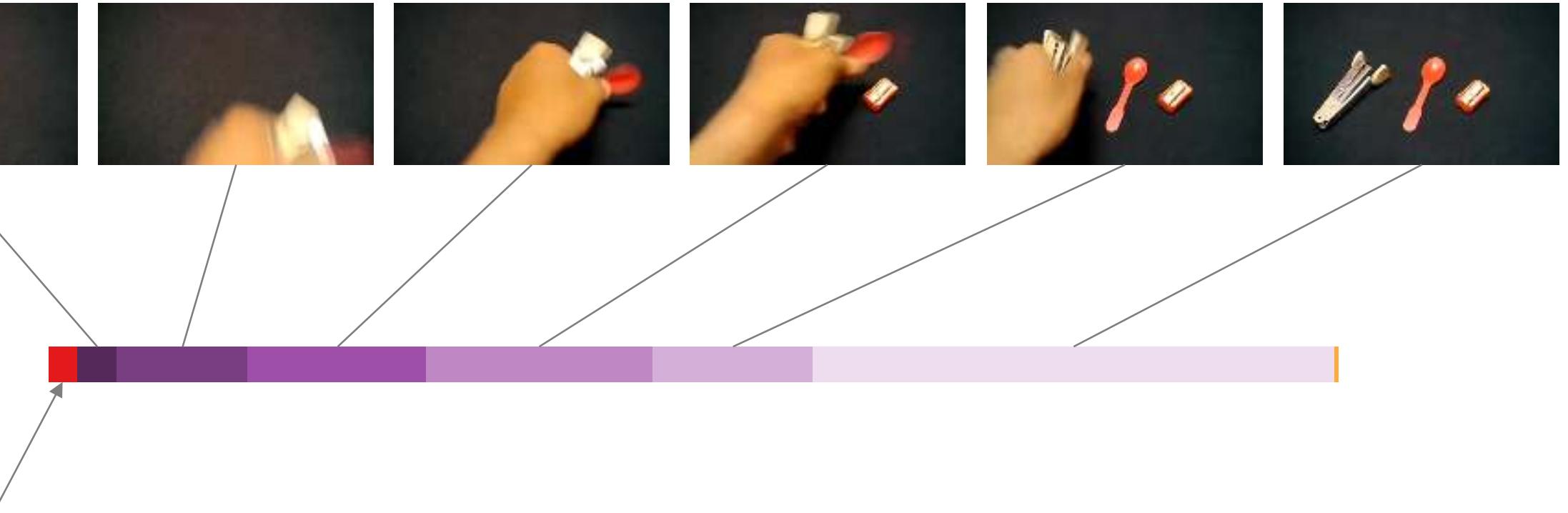
MODEL



Putting ?, ? and ?

# Frame Attributions in Video Models

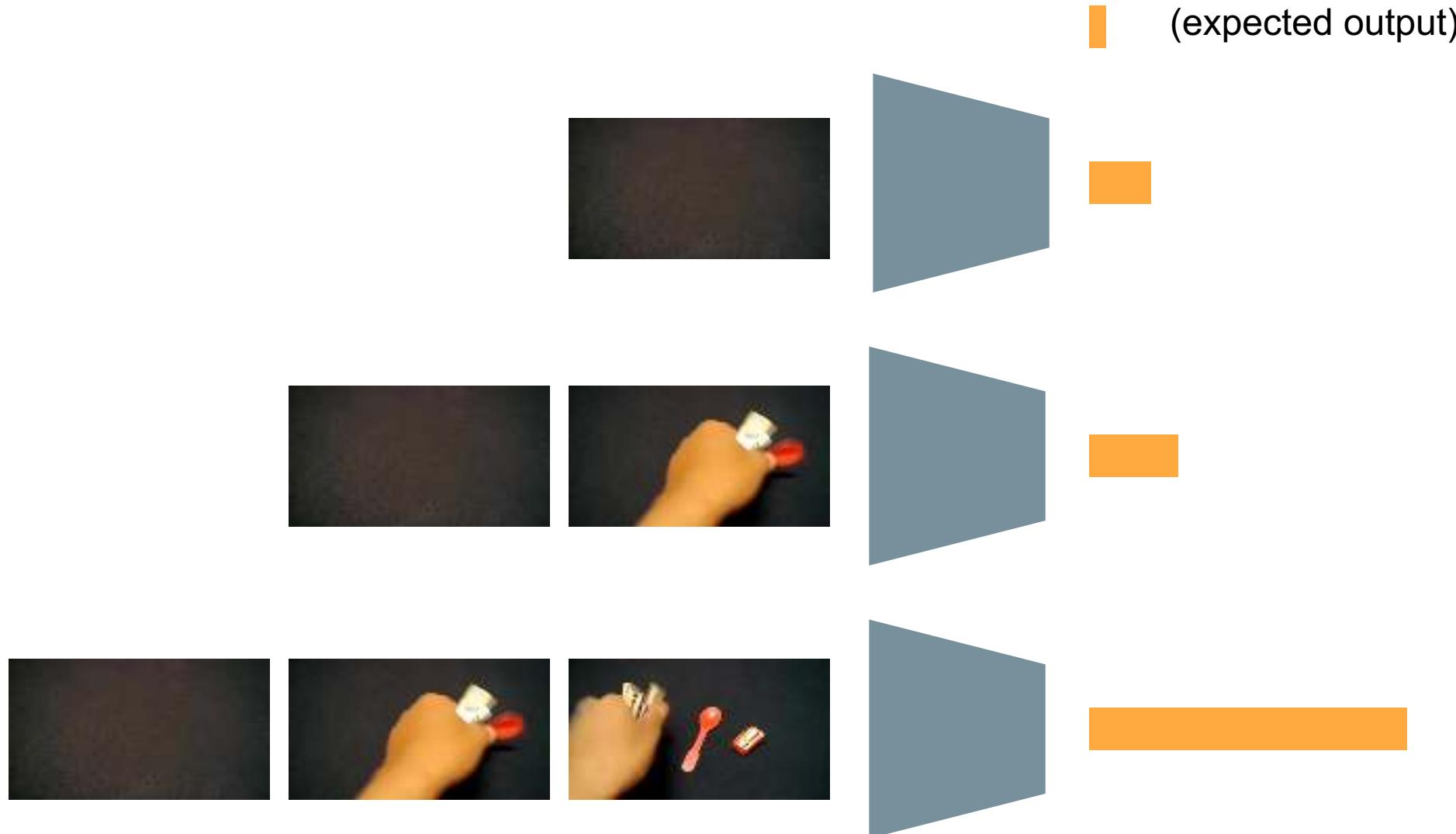
with: Will Price



Expected output  
(Prior probability for  
classification model)

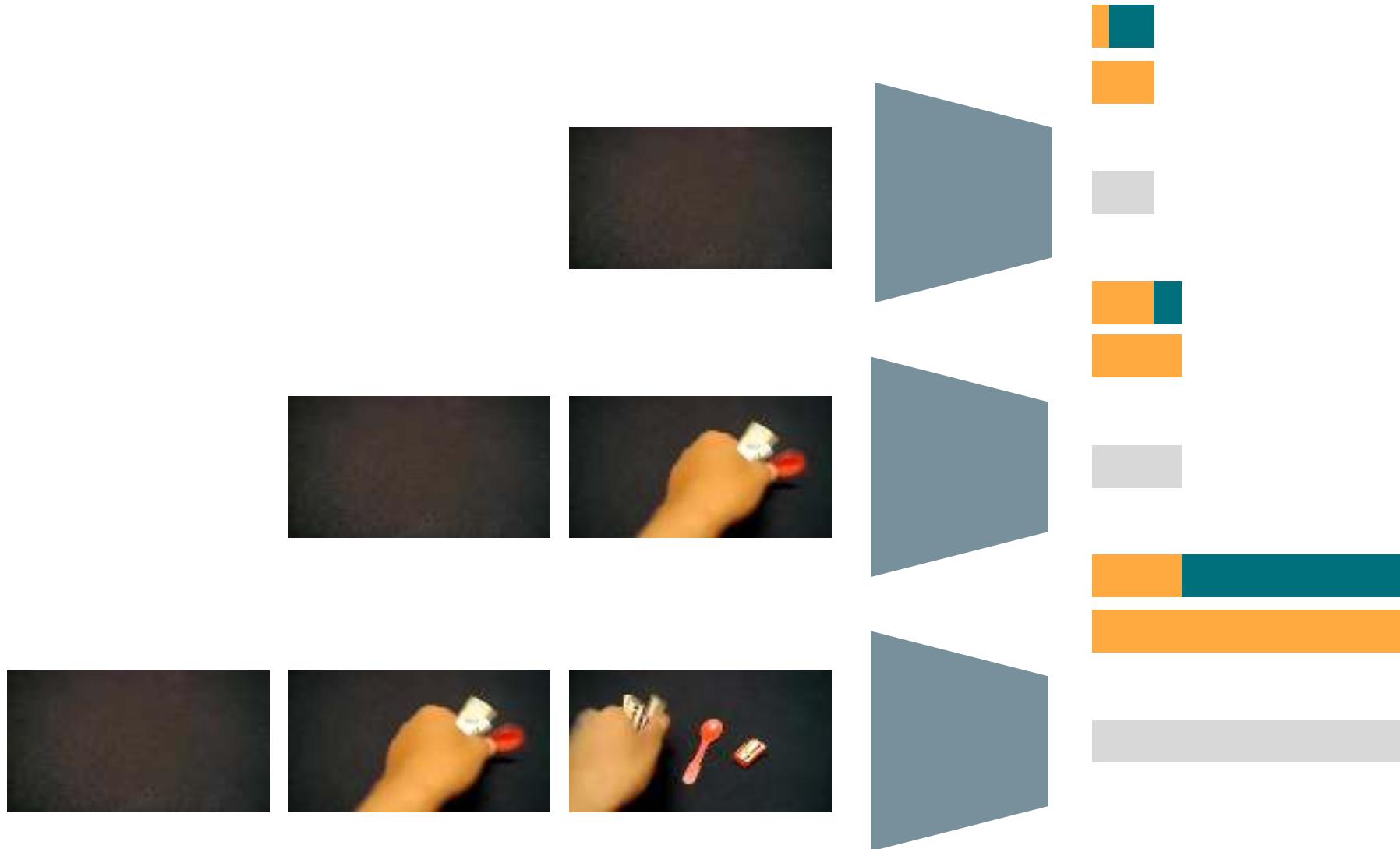
# Frame Attributions in Video Models

with: Will Price



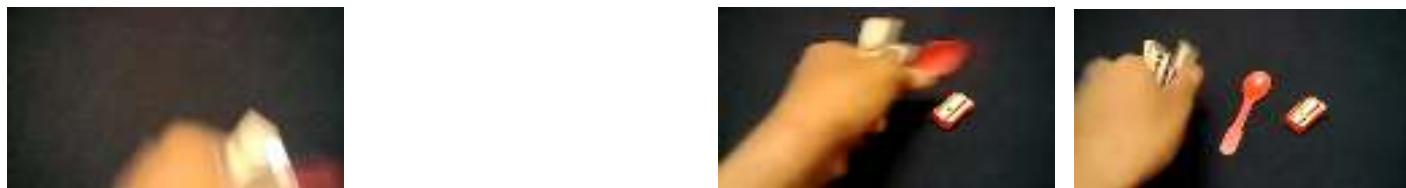
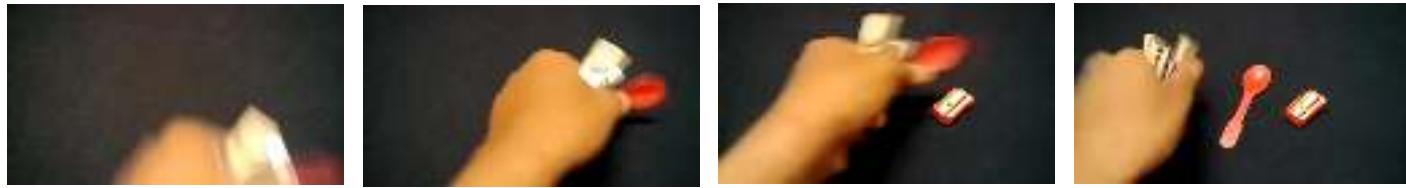
# Frame Attributions in Video Models

with: Will Price



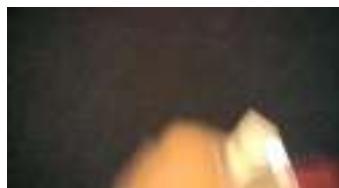
# Frame Attributions in Video Models

with: Will Price



# Frame Attributions in Video Models

with: Will Price



# Frame Attributions in Video Models

with: Will Price

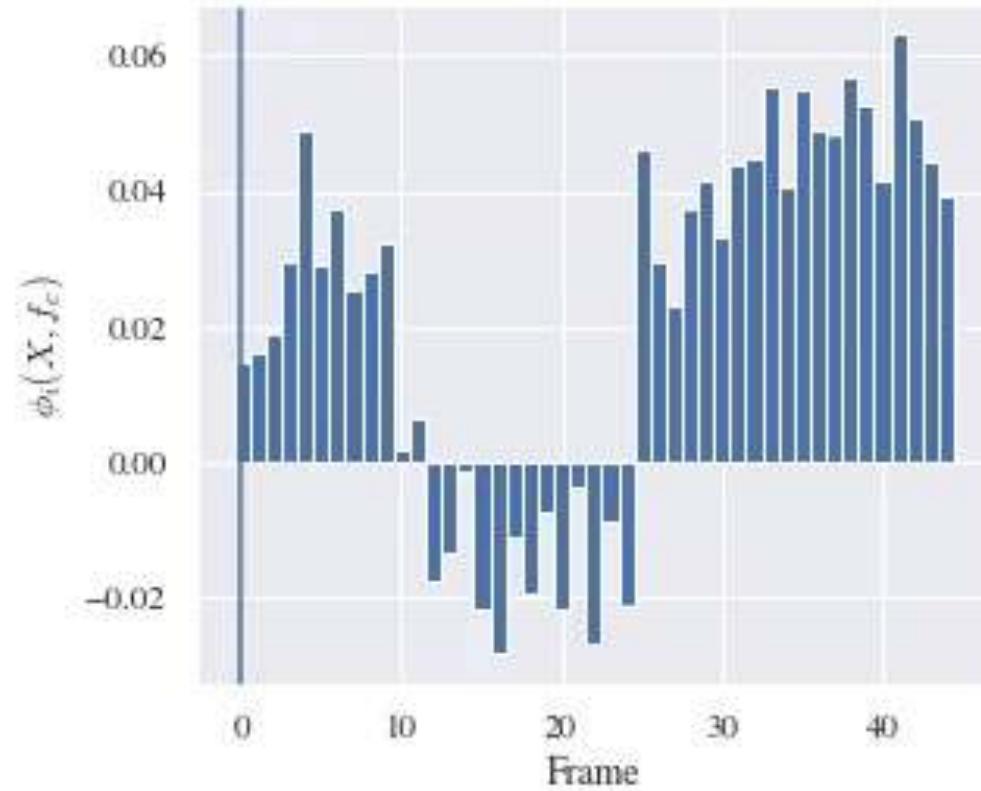


$$\Delta_3(\{1,2,4,5\}) = -2$$



# Frame Attributions in Video Models

with: Will Price



Showing that something is empty



# Frame Attributions in Video Models

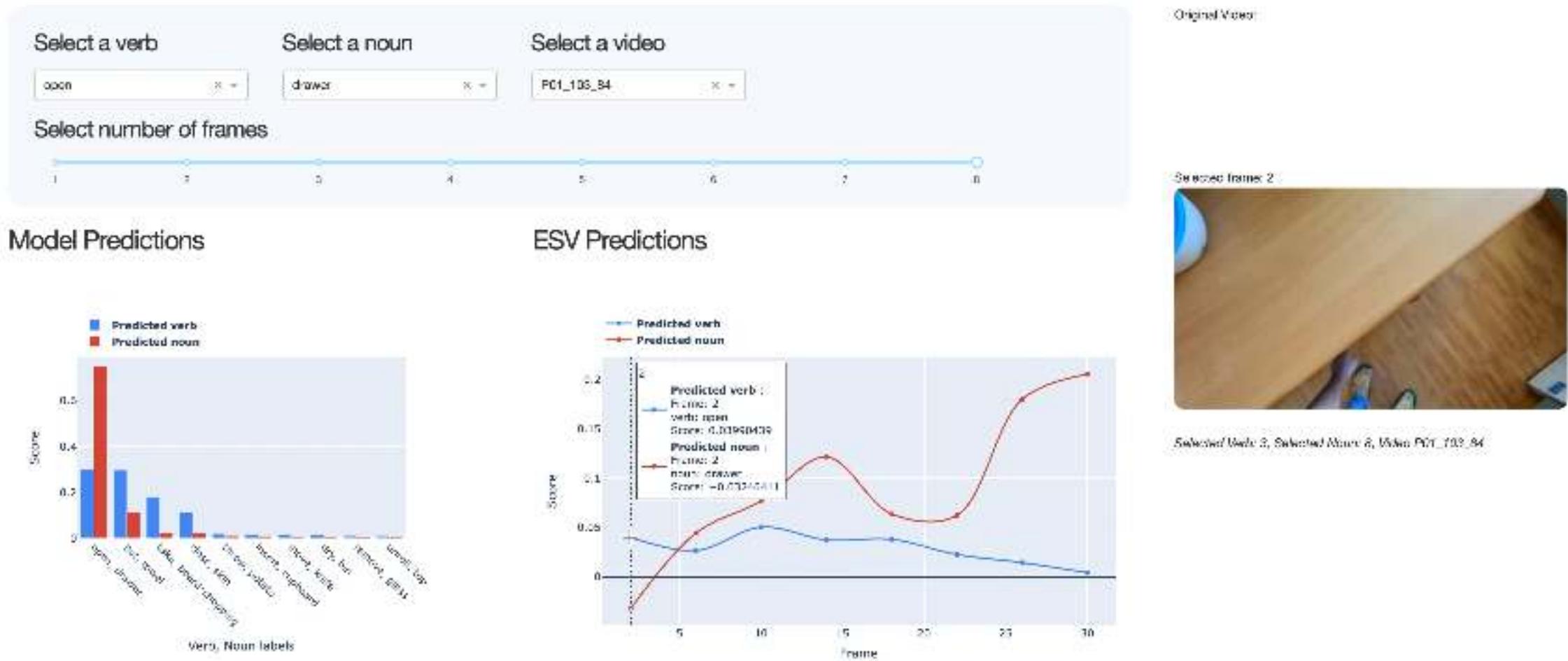
with: Will Price

# Dashboard

# Frame Attributions in Video Models

with: Will Price  
Tom Stark

## ESVs Dashboard for Epic



# Frame Attributions in Video Models

with: Will Price  
Tom Stark

## ESVs Dashboard for Epic

The figure shows a user interface for video analysis. At the top, three input fields allow selecting a verb ('cut'), a noun ('tomato'), and a video file ('P01\_17\_126'). Below these is a slider for 'Select number of frames' ranging from 1 to 8, with frame 8 selected. The next section, 'Model Predictions', displays a bar chart comparing predicted verbs (blue) and nouns (red) across various actions. The final section, 'ESV Predictions', shows a line graph of predicted verb and noun scores over time (Frame 0 to 800). A thumbnail image on the right shows a person chopping tomatoes.

# In today's tutorial



Motivation and Datasets in  
Egocentric Video Understanding



Video Understanding  
Out of the Frame



Video Understanding:  
Data and Tasks



Teaser: The Wizard of Oz  
at the Sphere



Videos are Multimodal



Outlook into the Future of  
Egocentric Vision



Connected Videos of One's Life



Conclusion



# Multi-modal learning...

with: Vangelis Kazakos  
Arsha Nagrani.  
Andrew Zisserman

Jaesung Huh  
Jacob Chalk

- The magic of audio-visual understanding...
- Object-Object interactions



# Multi-modal learning...

with: Vangelis Kazakos  
Arsha Nagrani.  
Andrew Zisserman

Jaesung Huh  
Jacob Chalk

- The magic of audio-visual understanding...
- Object-Object interactions
- Material sounds



# Multi-modal learning...

with: Vangelis Kazakos  
Arsha Nagrani.  
Andrew Zisserman

Jaesung Huh  
Jacob Chalk

- The magic of audio-visual understanding...
- Object-Object interactions
- Material sounds
- Sound-emitting objects





with: Jaesung Huh\* & Jacob Chalk\*  
Vangelis Kazakos Andrew Zisserman



# EPIC-Sounds: A Large-scale Dataset of Actions That Sound

Jaesung Huh\*, Jacob Chalk\*, Evangelos Kazakos, Dima Damen, Andrew Zisserman

\* : Equal contribution

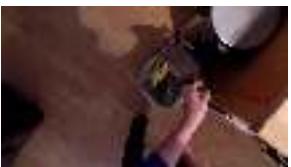


Dima Damen  
ICVSS2025

# Motivation

with: Jaesung Huh\* & Jacob Chalk\*  
Vangelis Kazakos Andrew Zisserman

Video

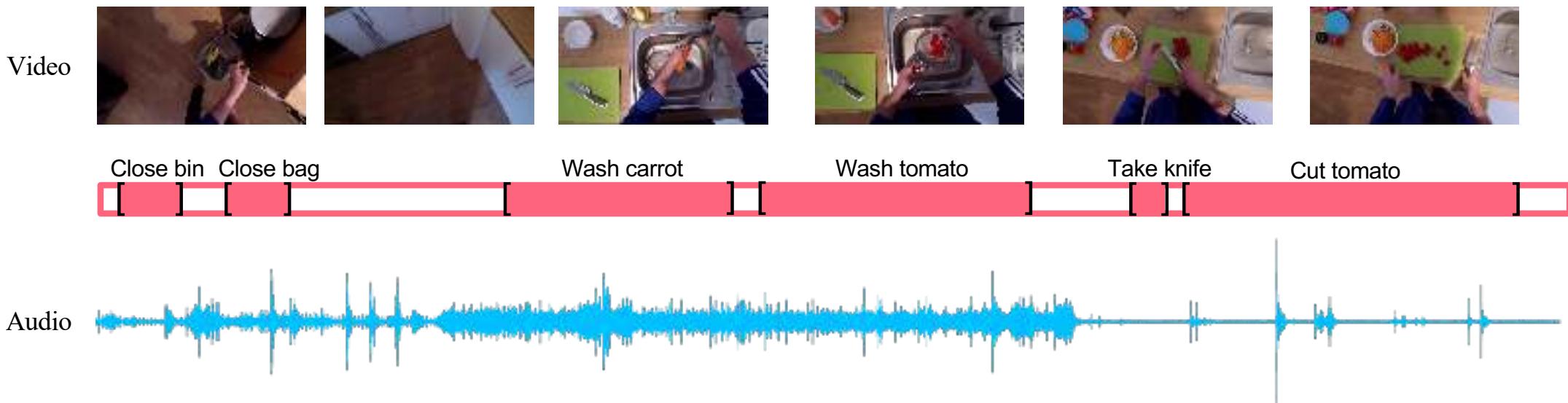


Audio



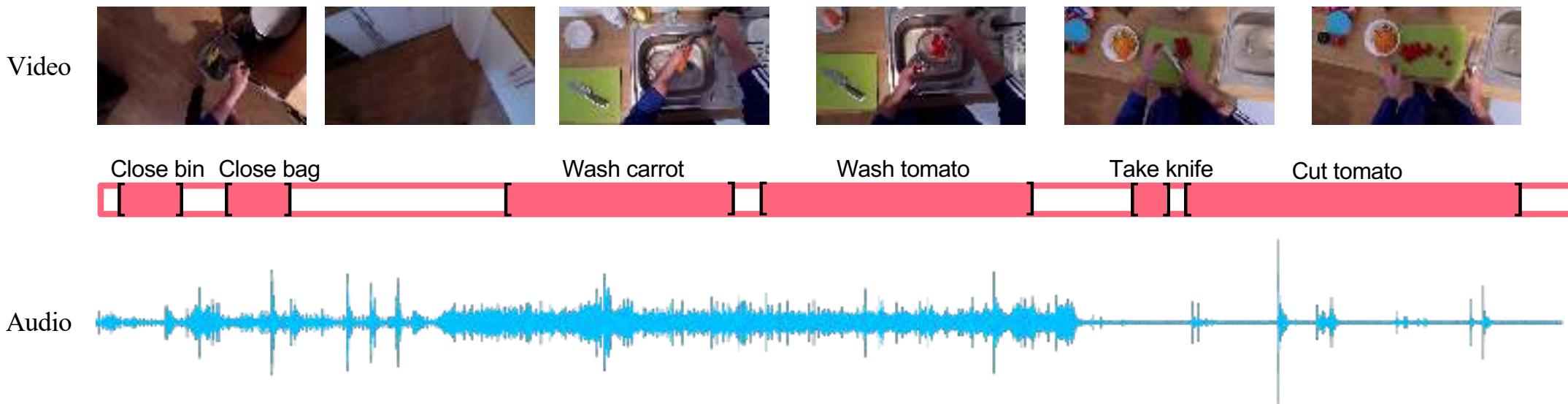
# Motivation

with: Jaesung Huh\* & Jacob Chalk\*  
Vangelis Kazakos Andrew Zisserman



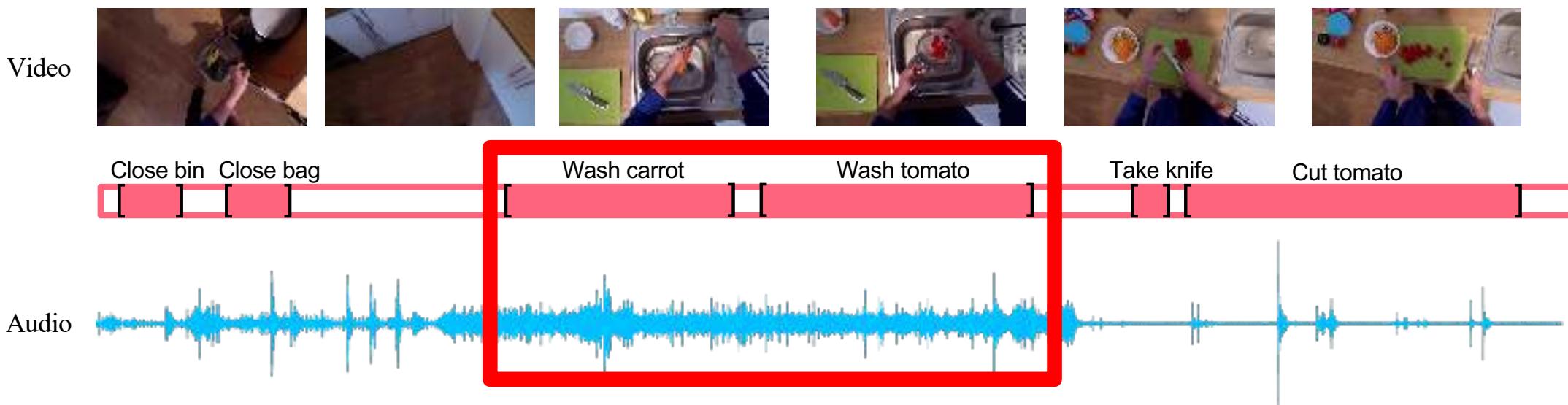
# Motivation

with: Jaesung Huh\* & Jacob Chalk\*  
Vangelis Kazakos Andrew Zisserman



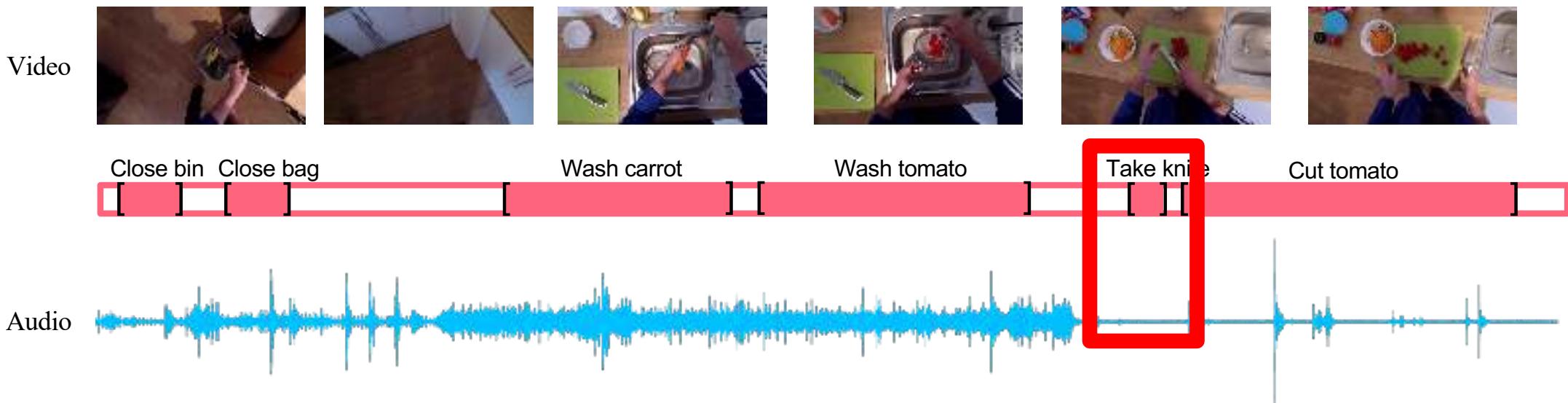
# Motivation

with: Jaesung Huh\* & Jacob Chalk\*  
Vangelis Kazakos Andrew Zisserman



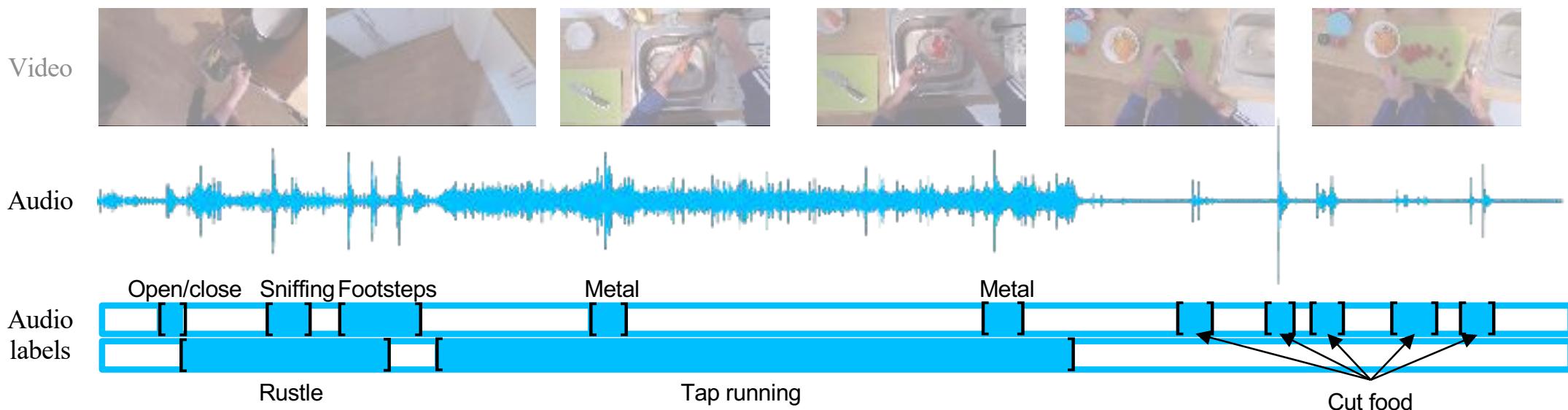
# Motivation

with: Jaesung Huh\* & Jacob Chalk\*  
Vangelis Kazakos Andrew Zisserman



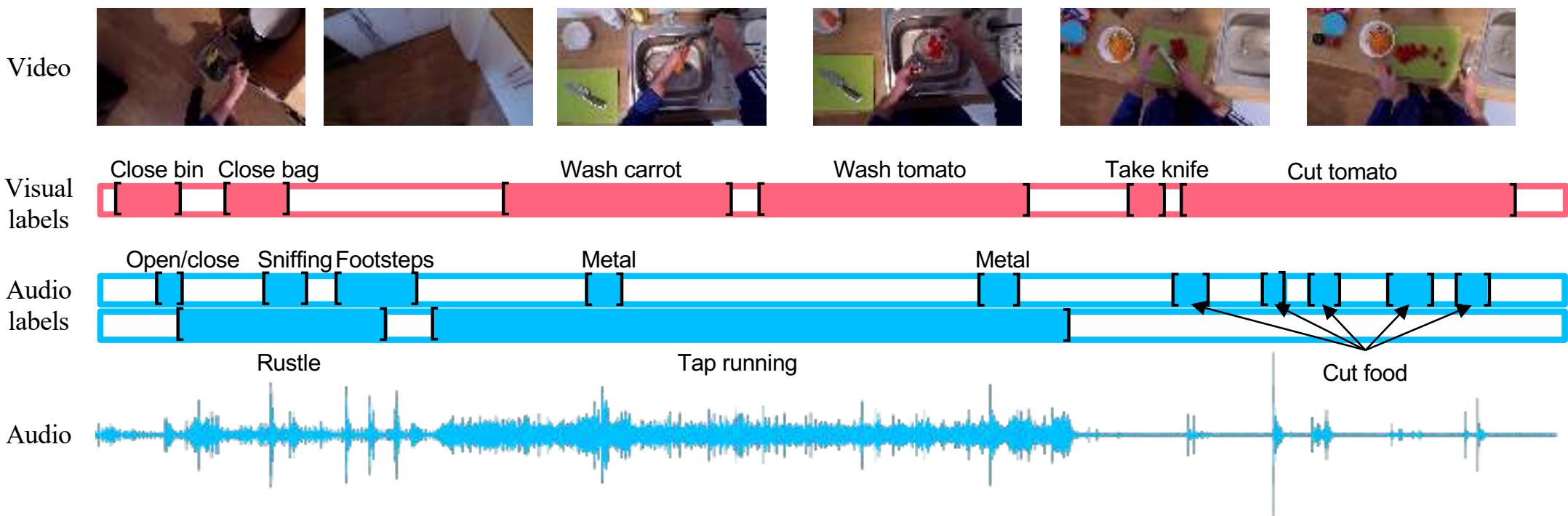
# Motivation

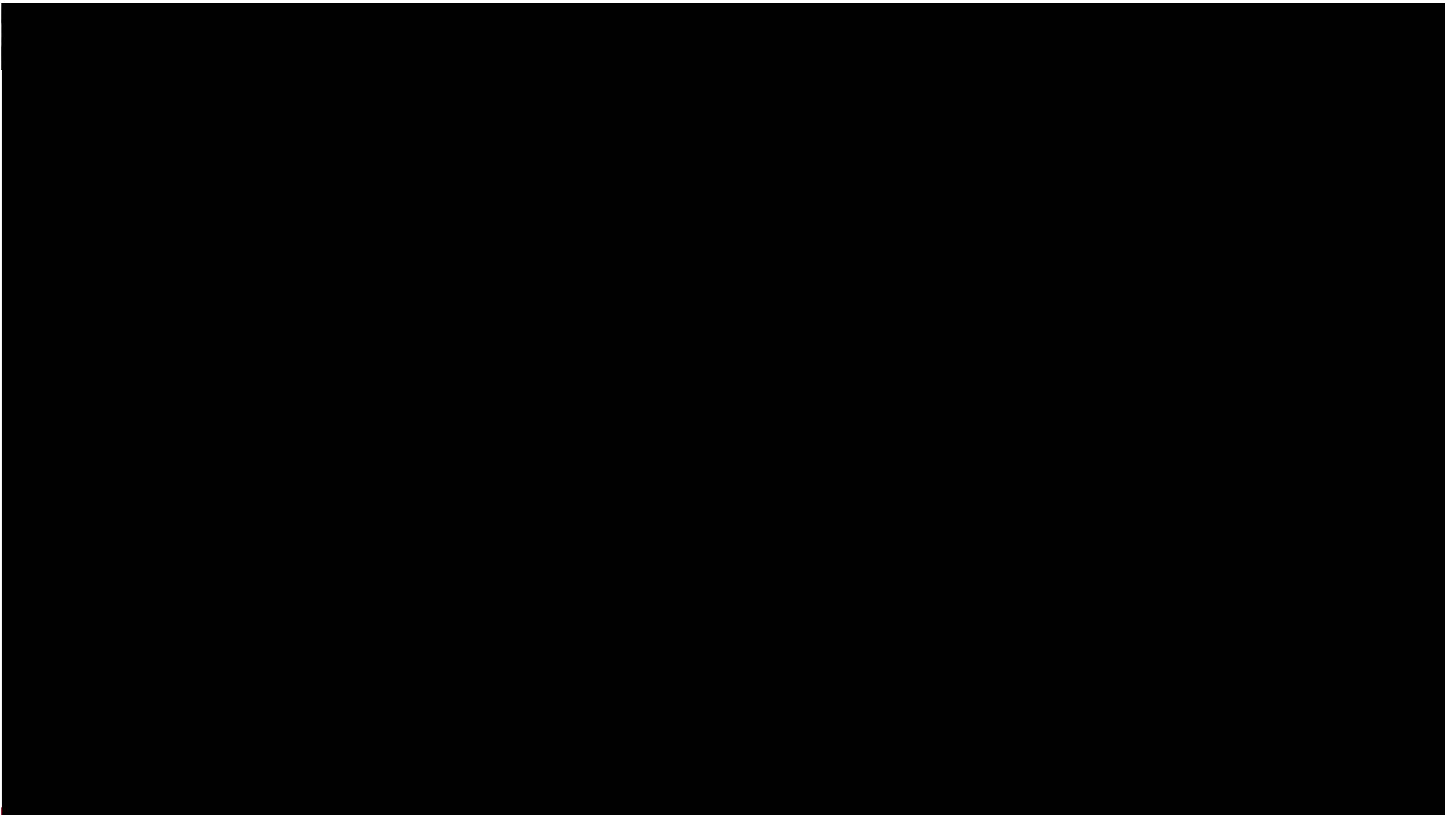
with: Jaesung Huh\* & Jacob Chalk\*  
Vangelis Kazakos Andrew Zisserman



# Motivation

with: Jaesung Huh\* & Jacob Chalk\*  
Vangelis Kazakos Andrew Zisserman





## EPIC-KITCHENS VIDEOS

100 hours

45 kitchens

Visual Action Annotations  
90K visual actions  
97 verb classes  
300 noun classes

## EPIC-Sounds

Audio-Based Annotations  
79K categorised audio events  
44 sound categories  
39K uncategorised events



spray





with: Jacob Chalk\* Jaesung Huh\*  
Vangelis Kazakos Andrew Zisserman



# TIM: A Time Interval Machine for Audio-Visual Action Recognition

Jacob Chalk\*, Jaesung Huh\*, Evangelos Kazakos, Andrew Zisserman, Dima Damen

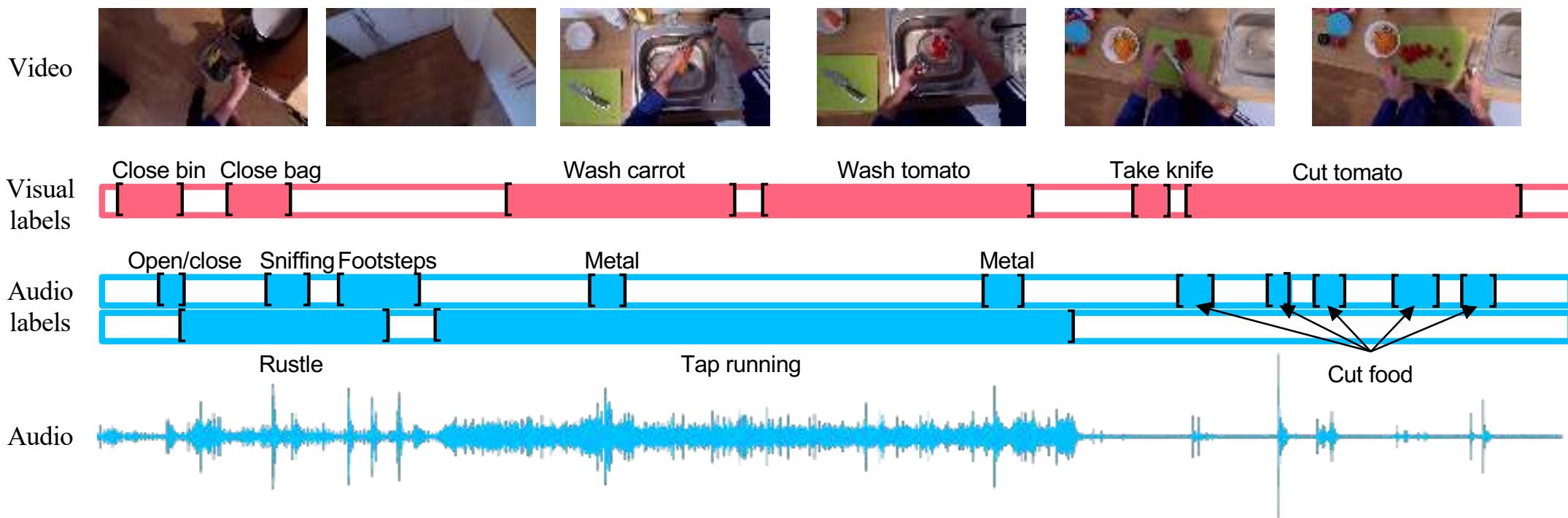
\* : Equal contribution



Dima Damen  
ICVSS2025

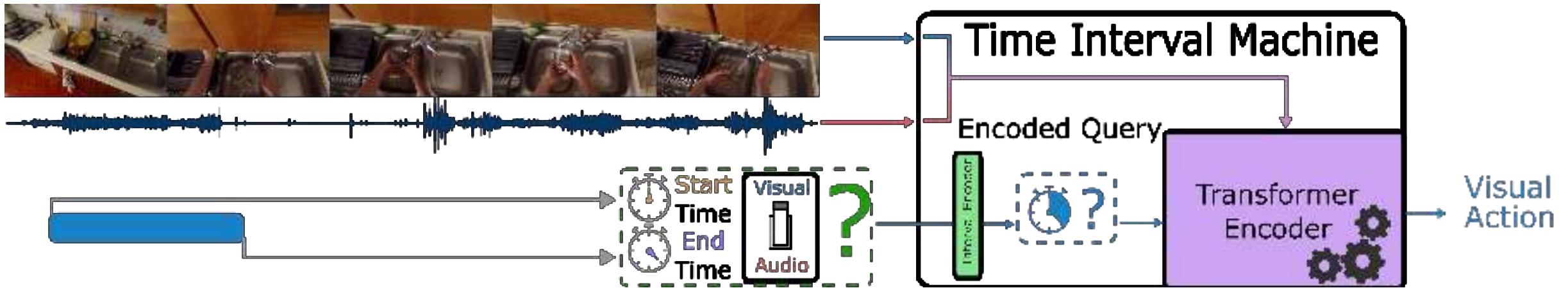
# Multi-Modal Long-Form Dataset

with: Jacob Chalk\* Jaesung Huh\*  
Vangelis Kazakos Andrew Zisserman



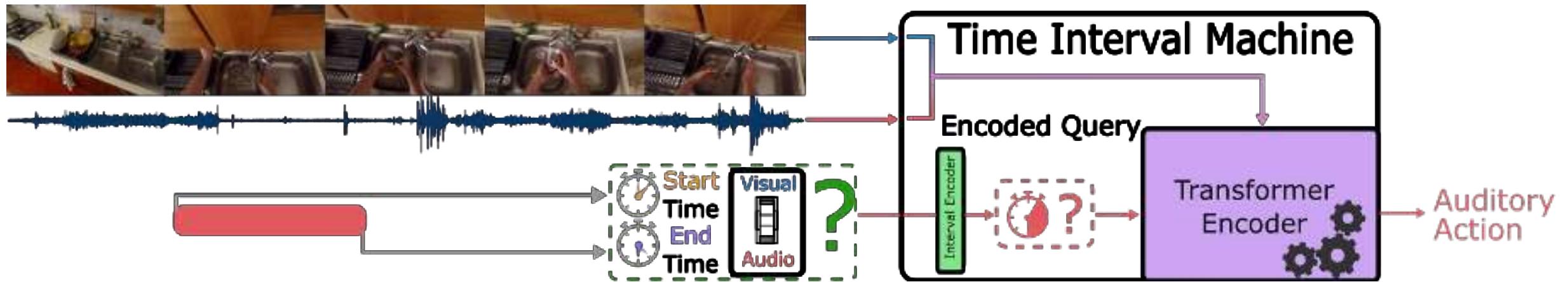
# TIM: A Time-Interval Audio-Visual Machine

with: Jacob Chalk\* Jaesung Huh\*  
Vangelis Kazakos Andrew Zisserman



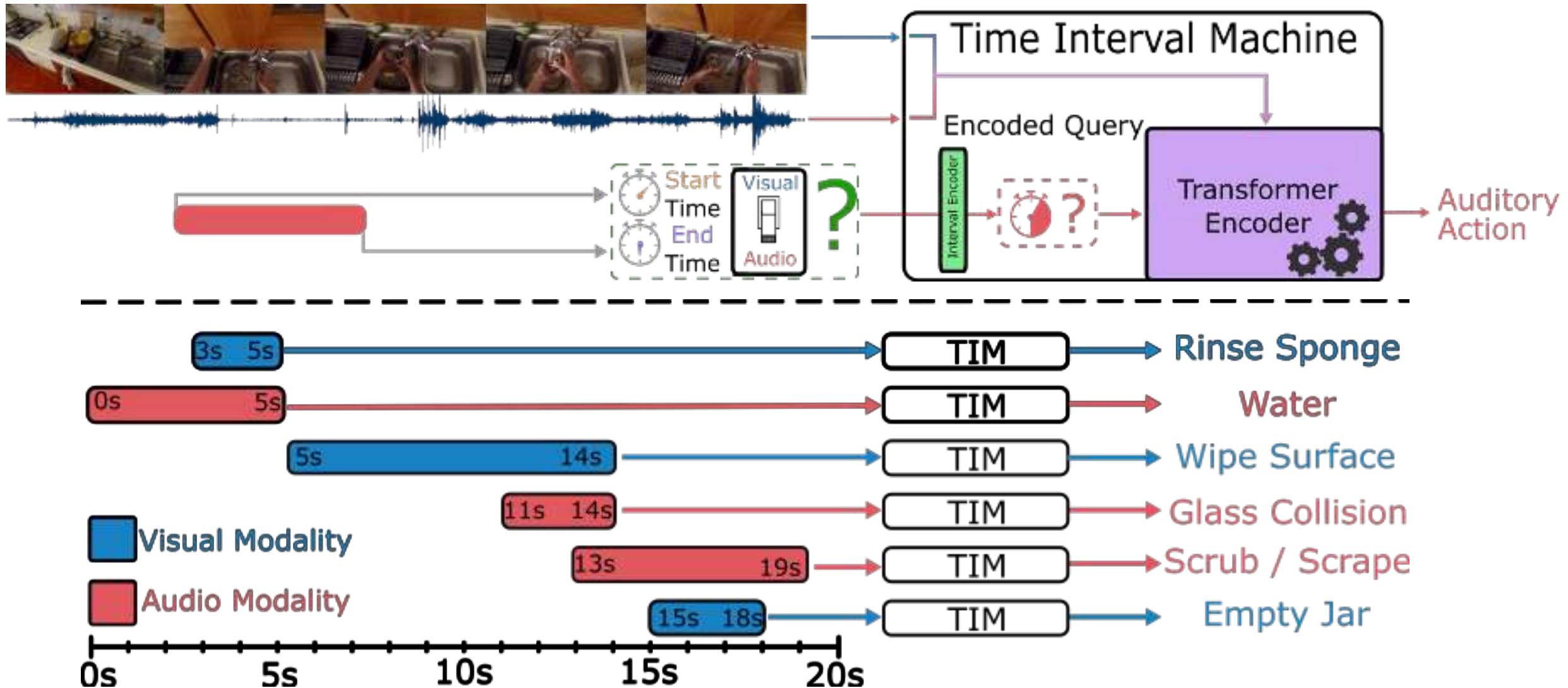
# TIM: A Time-Interval Audio-Visual Machine

with: Jacob Chalk\* Jaesung Huh\*  
Vangelis Kazakos Andrew Zisserman



# TIM: A Time-Interval Audio-Visual Machine

with: Jacob Chalk\* Jaesung Huh\*  
Vangelis Kazakos Andrew Zisserman



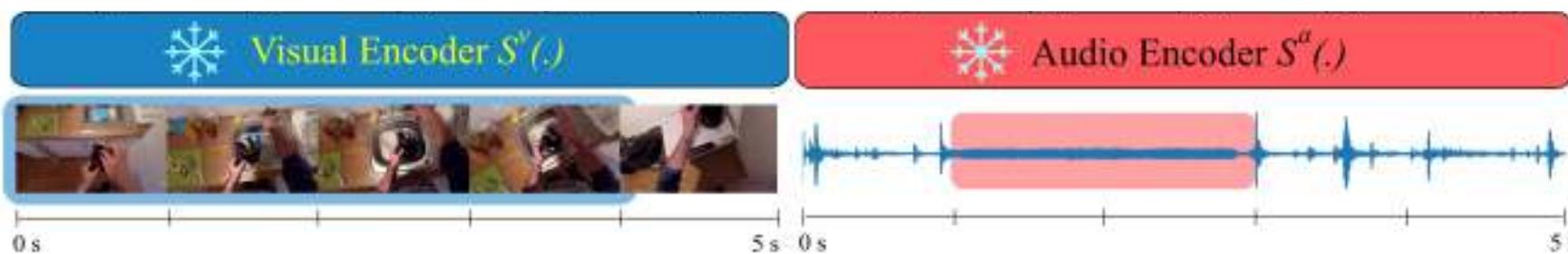
# TIM: A Time-Interval Audio-Visual Machine

with: Jacob Chalk\* Jaesung Huh\*  
Vangelis Kazakos Andrew Zisserman



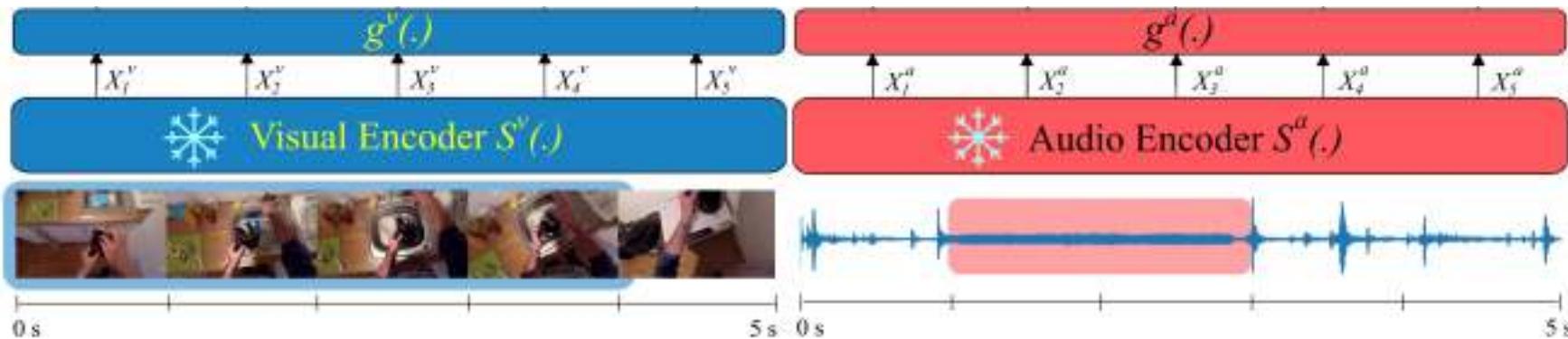
# TIM: A Time-Interval Audio-Visual Machine

with: Jacob Chalk\* Jaesung Huh\*  
Vangelis Kazakos Andrew Zisserman



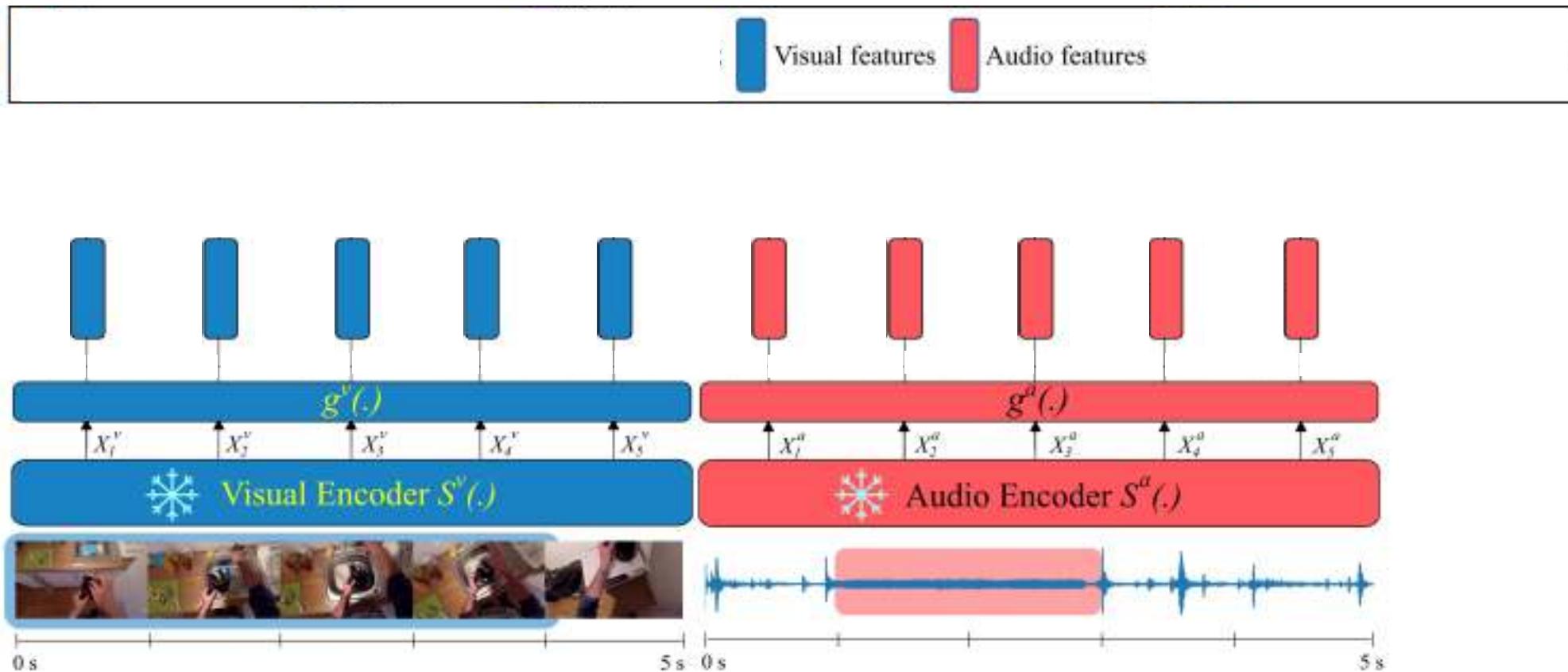
# TIM: A Time-Interval Audio-Visual Machine

with: Jacob Chalk\* Jaesung Huh\*  
Vangelis Kazakos Andrew Zisserman



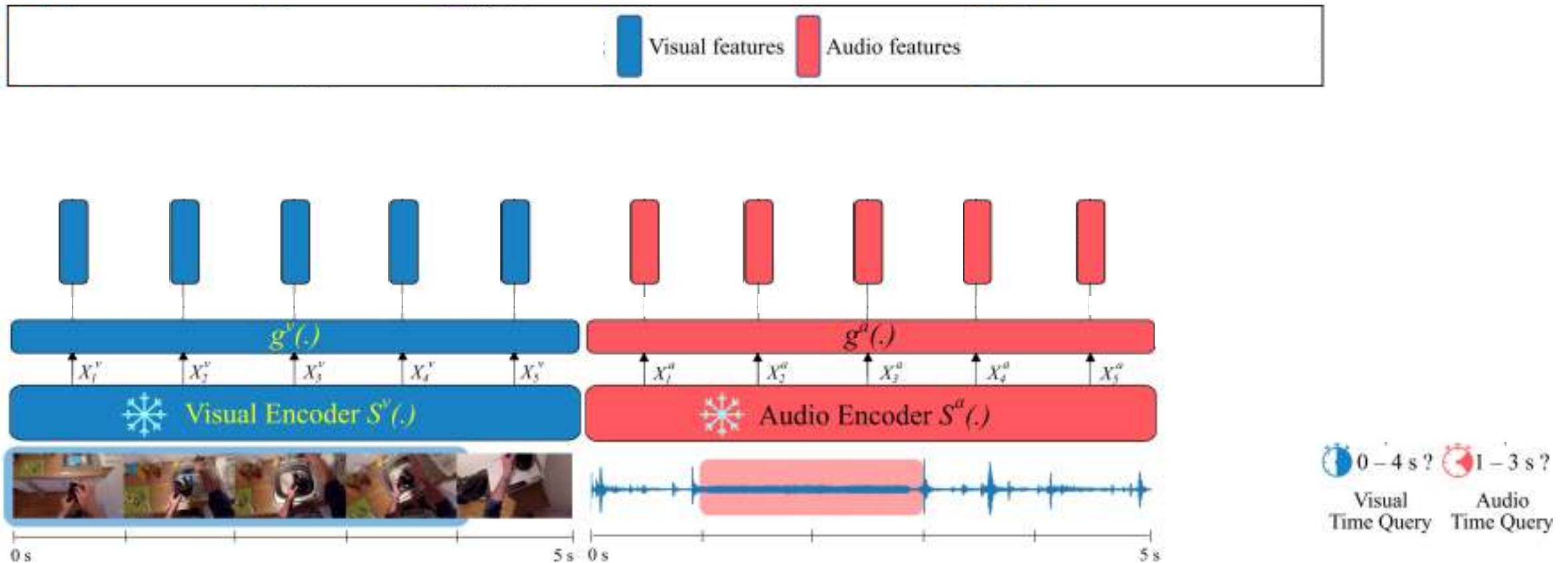
# TIM: A Time-Interval Audio-Visual Machine

with: Jacob Chalk\* Jaesung Huh\*  
Vangelis Kazakos Andrew Zisserman



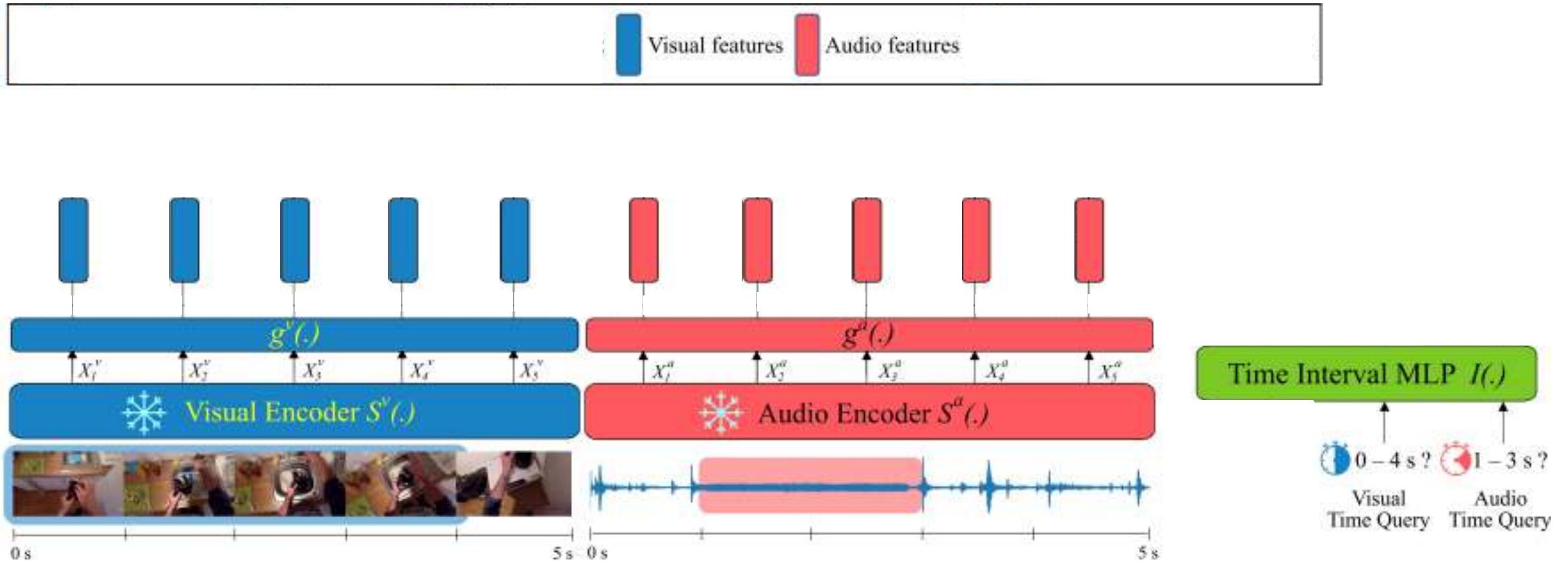
# TIM: A Time-Interval Audio-Visual Machine

with: Jacob Chalk\* Jaesung Huh\*  
Vangelis Kazakos Andrew Zisserman



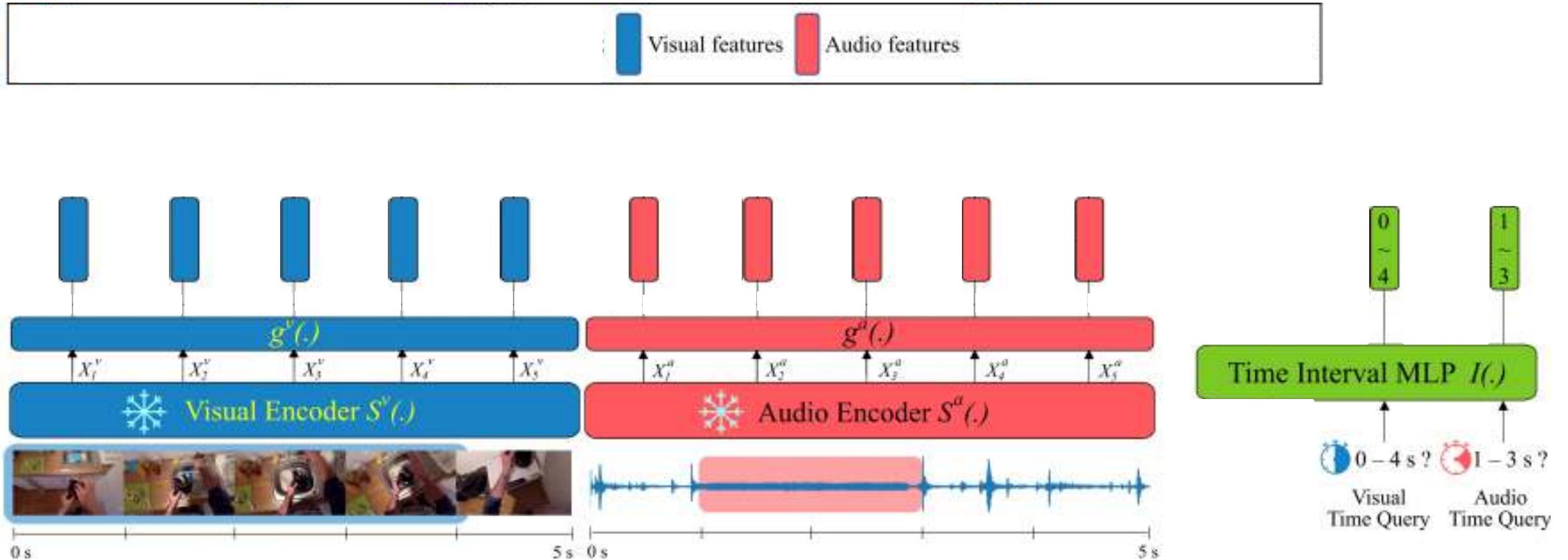
# TIM: A Time-Interval Audio-Visual Machine

with: Jacob Chalk\* Jaesung Huh\*  
Vangelis Kazakos Andrew Zisserman



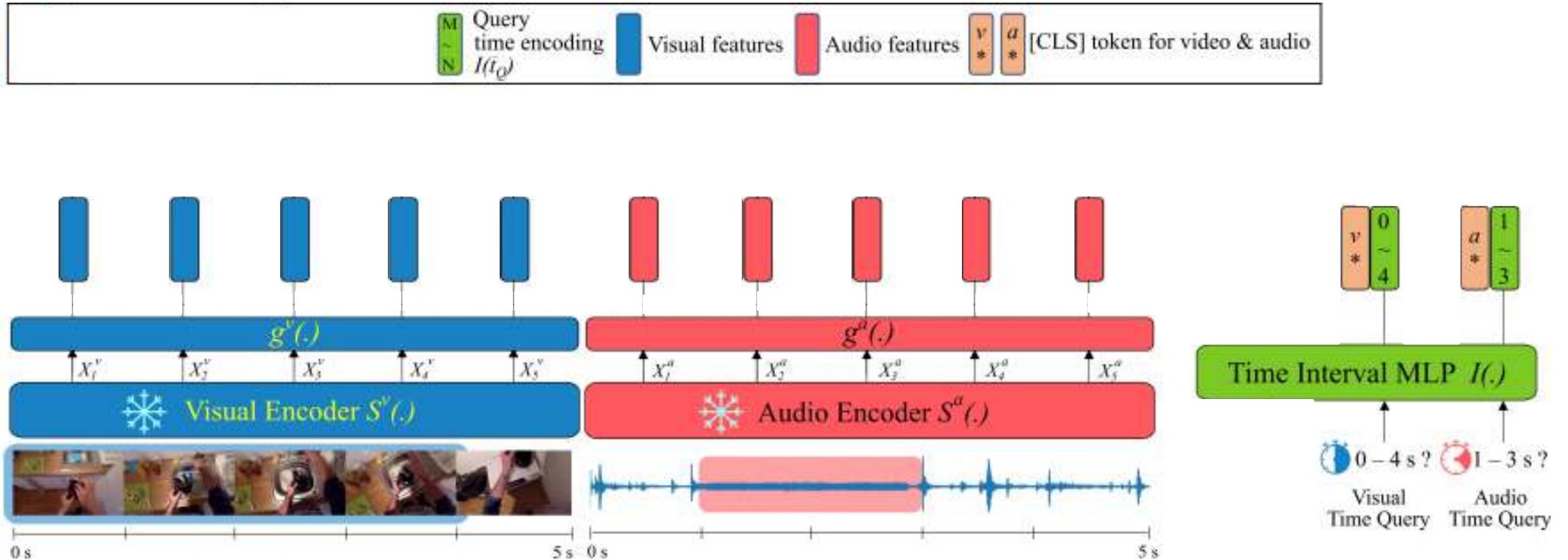
# TIM: A Time-Interval Audio-Visual Machine

with: Jacob Chalk\* Jaesung Huh\*  
Vangelis Kazakos Andrew Zisserman



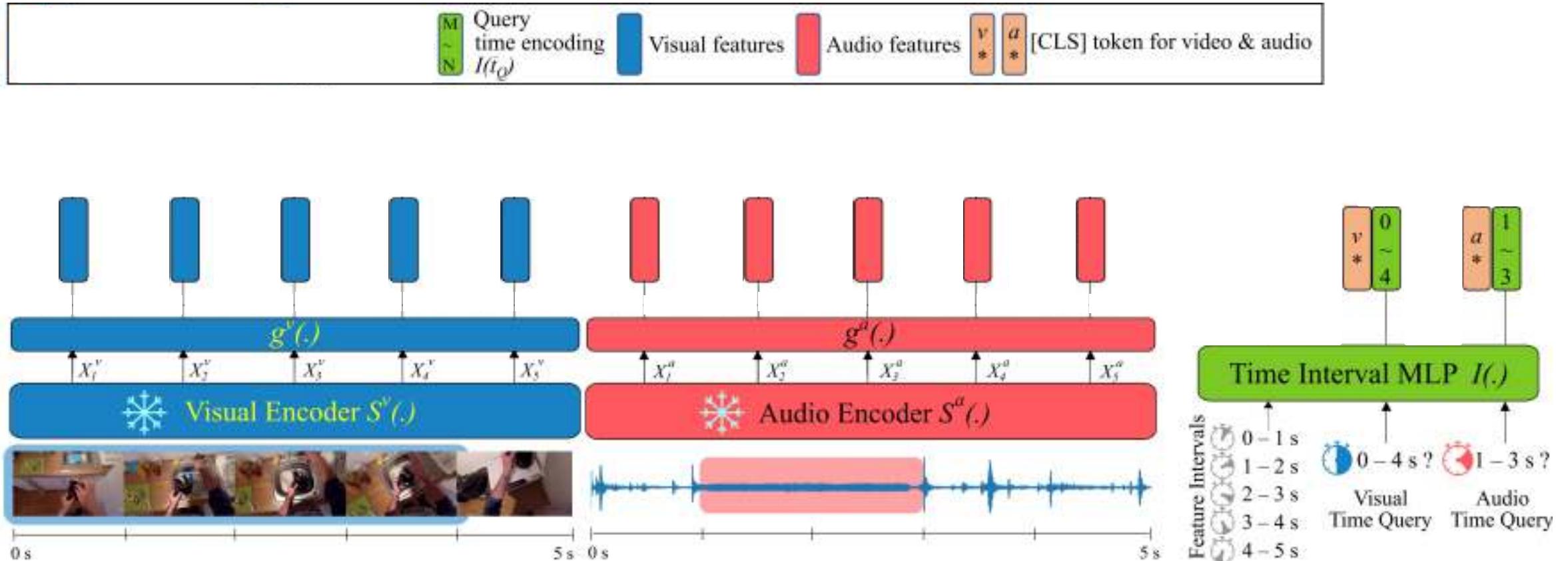
# TIM: A Time-Interval Audio-Visual Machine

with: Jacob Chalk\* Jaesung Huh\*  
Vangelis Kazakos Andrew Zisserman



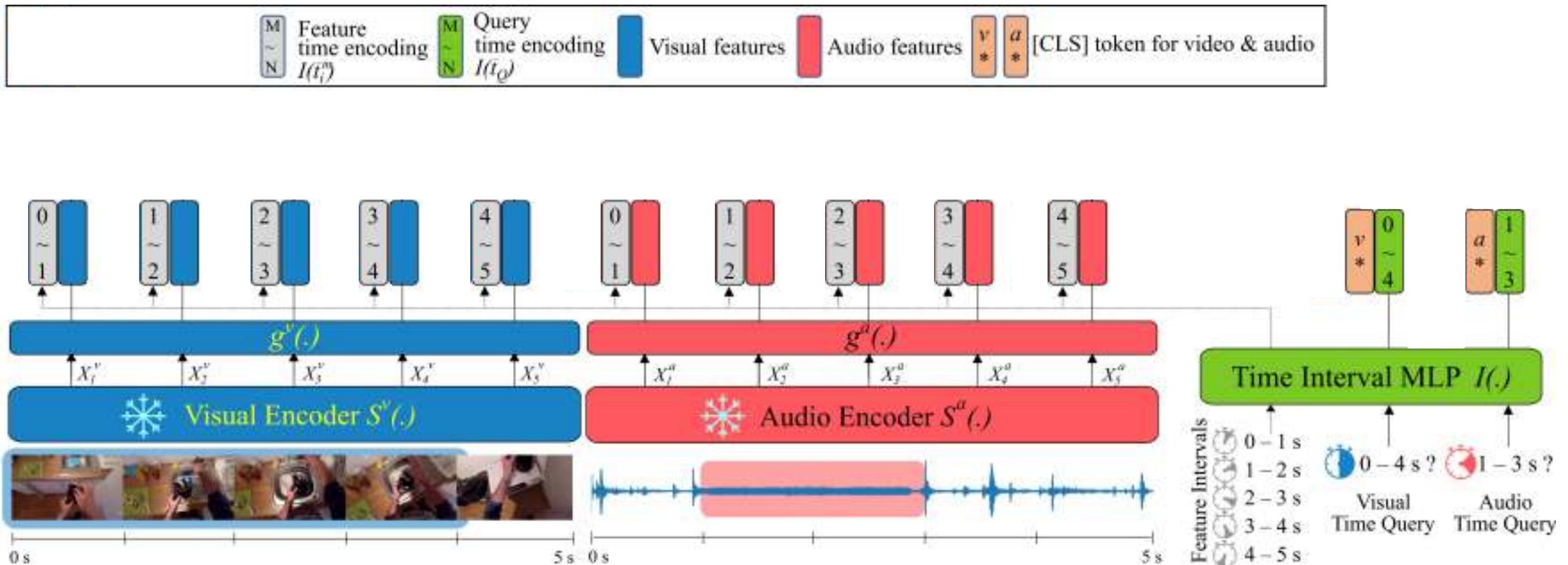
# TIM: A Time-Interval Audio-Visual Machine

with: Jacob Chalk\* Jaesung Huh\*  
Vangelis Kazakos Andrew Zisserman



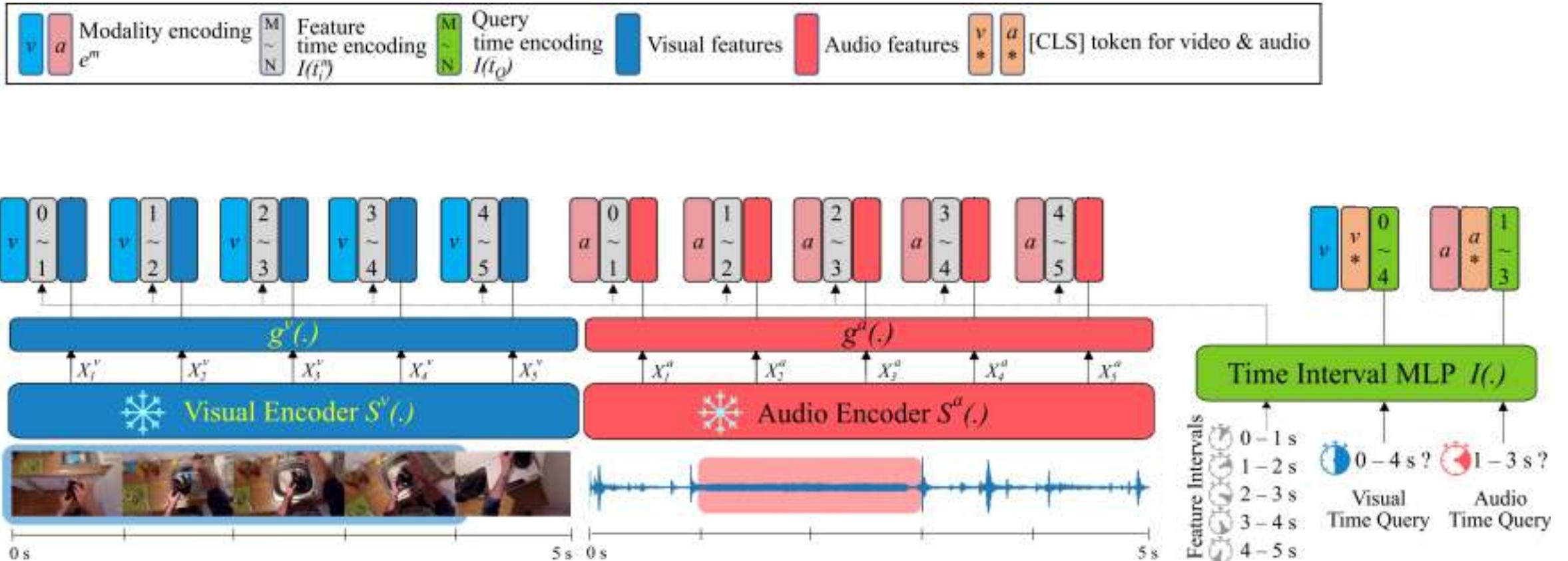
# TIM: A Time-Interval Audio-Visual Machine

with: Jacob Chalk\* Jaesung Huh\*  
Vangelis Kazakos Andrew Zisserman



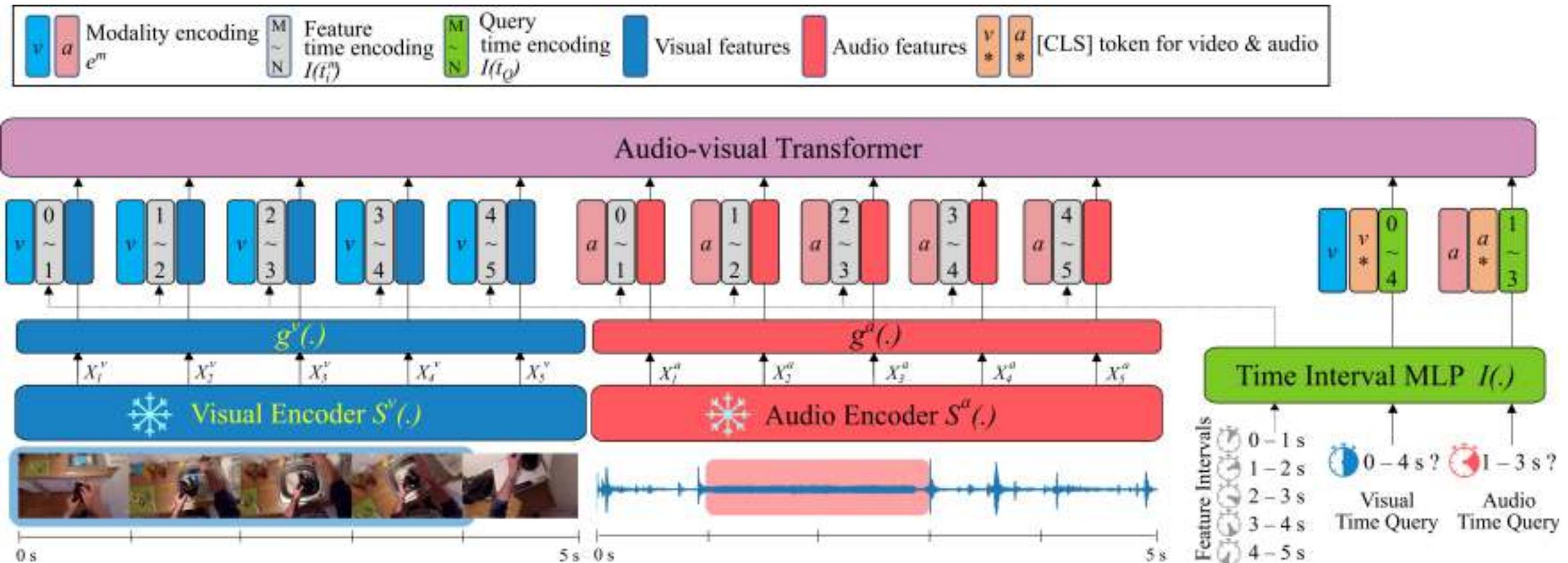
# TIM: A Time-Interval Audio-Visual Machine

with: Jacob Chalk\* Jaesung Huh\*  
Vangelis Kazakos Andrew Zisserman



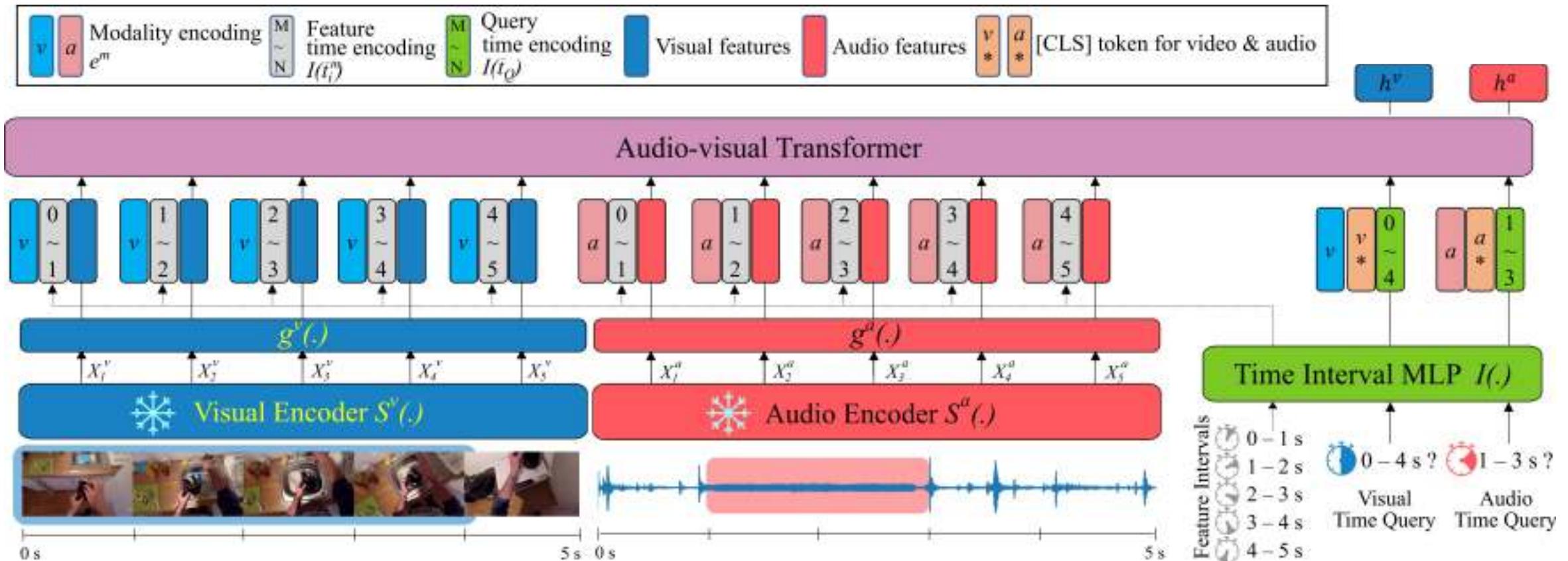
# TIM: A Time-Interval Audio-Visual Machine

with: Jacob Chalk\* Jaesung Huh\*  
Vangelis Kazakos Andrew Zisserman



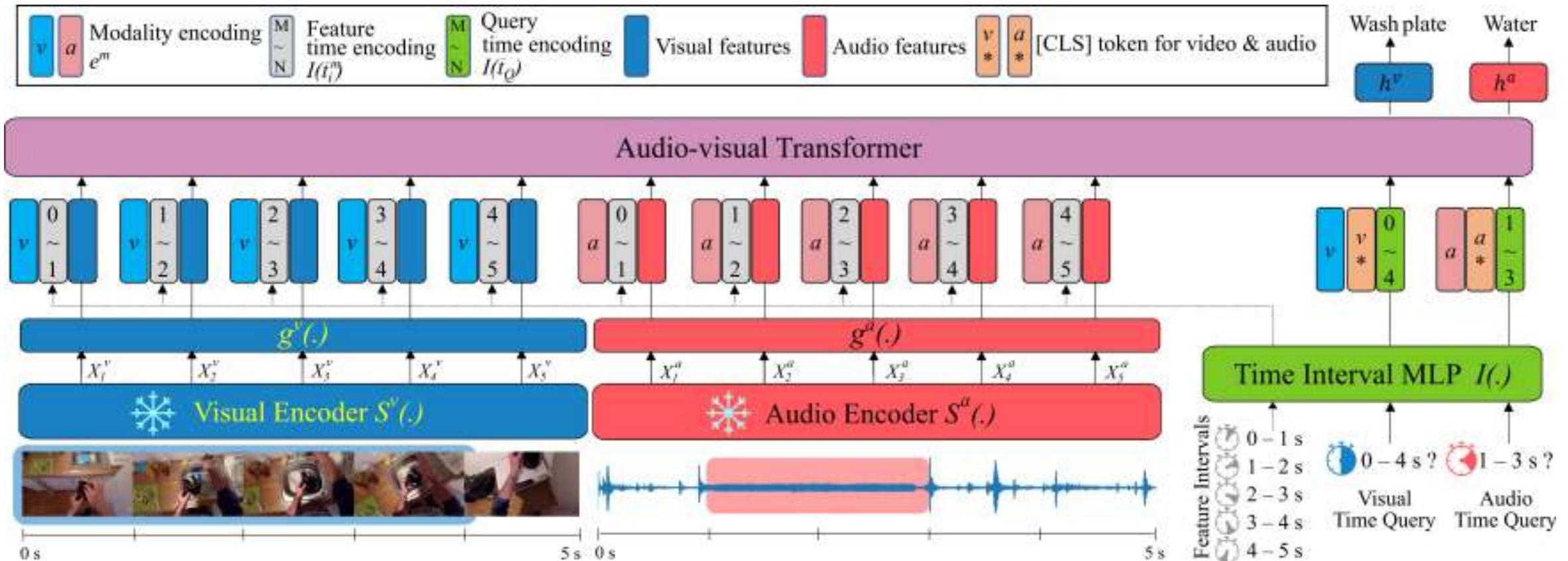
# TIM: A Time-Interval Audio-Visual Machine

with: Jacob Chalk\* Jaesung Huh\*  
Vangelis Kazakos Andrew Zisserman



# TIM: A Time-Interval Audio-Visual Machine

with: Jacob Chalk\* Jaesung Huh\*  
Vangelis Kazakos Andrew Zisserman



# TIM: A Time-Interval Audio-Visual Machine

with: Jacob Chalk\* Jaesung Huh\*  
Vangelis Kazakos Andrew Zisserman

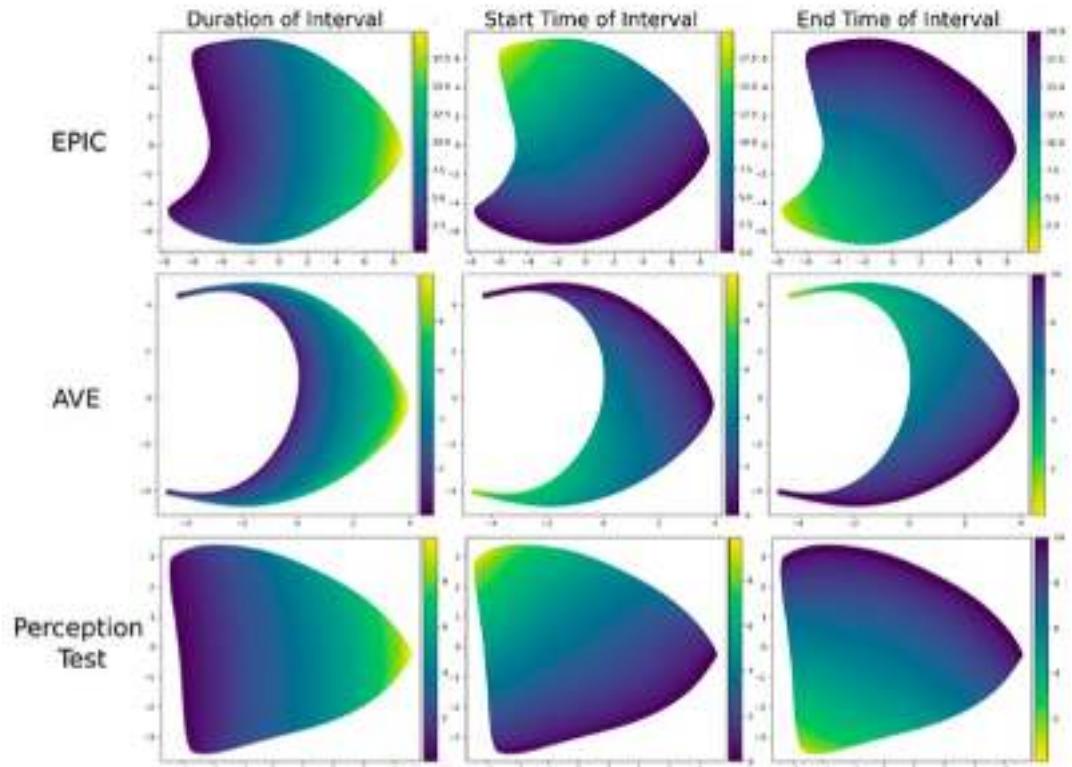
Model	xp	LLM	Verb	Noun	Action
<i>Visual-only models</i>					
MFormer-HR [37]	336p	✗	67.0	58.5	44.5
MoViNet-A6 [27]	320p	✗	72.2	57.3	47.7
MeMViT [55]	224p	✗	71.4	60.3	48.4
Omnivore [14]	224p	✗	69.5	61.7	49.9
MTV [59]	280p	✗	69.9	63.9	50.5
LaViLa (TSF-L) [63]	224p	✓	72.0	62.9	51.0
AVION (ViT-L) [62]	224p	✓	73.0	65.4	54.4
<b>TIM (ours)</b>	224p	✗	<b>76.2</b>	<b>66.4</b>	<b>56.4</b>
<i>Audio-visual models</i>					
TBN [24]	224p	✗	66.0	47.2	36.7
MBT [34]	224p	✗	64.8	58.0	43.4
MTCN [25]	336p	✗	70.7	62.1	49.6
M&M [57]	420p	✗	72.0	66.3	53.6
<b>TIM (ours)</b>	224p	✗	<b>77.5</b>	<b>67.4</b>	<b>57.9</b>

Perception Test Action				
Model	MLP (V)	MTCN [25](A+V)	<b>TIM (V)</b>	<b>TIM (A+V)</b>
<b>Top-1 acc</b>	43.7	51.2	56.1	<b>61.1</b>
Perception Test Sound				
Model	MLP (A)	MTCN [25](A+V)	<b>TIM (A)</b>	<b>TIM (A+V)</b>
<b>Top-1 acc</b>	50.6	52.9	54.8	<b>56.1</b>

Table 5. Comparisons to trained recognition baselines on the Perception Test validation split. We show both action and sound recognition and the benefit of including audio-visual in TIM for both challenges. **V** : visual and **A** : audio input features. MLP is the result by training an MLP classifier with the features directly.

# TIM: A Time-Interval Audio-Visual Machine

with: Jacob Chalk\* Jaesung Huh\*  
Vangelis Kazakos Andrew Zisserman



# In today's tutorial



Motivation and Datasets in  
Egocentric Video Understanding



Video Understanding  
Out of the Frame



Video Understanding:  
Data and Tasks



Teaser: The Wizard of Oz  
at the Sphere



Videos are Multimodal



Outlook into the Future of  
Egocentric Vision



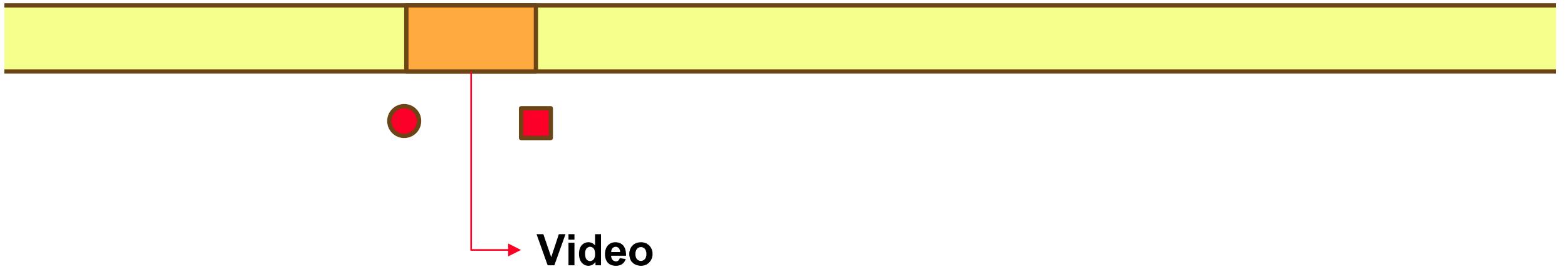
Connected Videos of One's Life



Conclusion

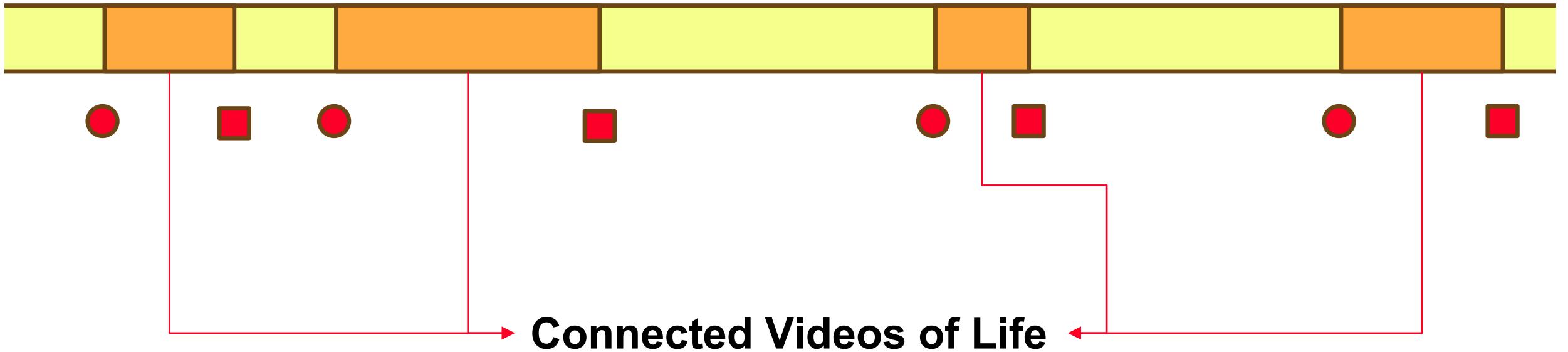


# Current Video Understanding...





# Upcoming Video Understanding...



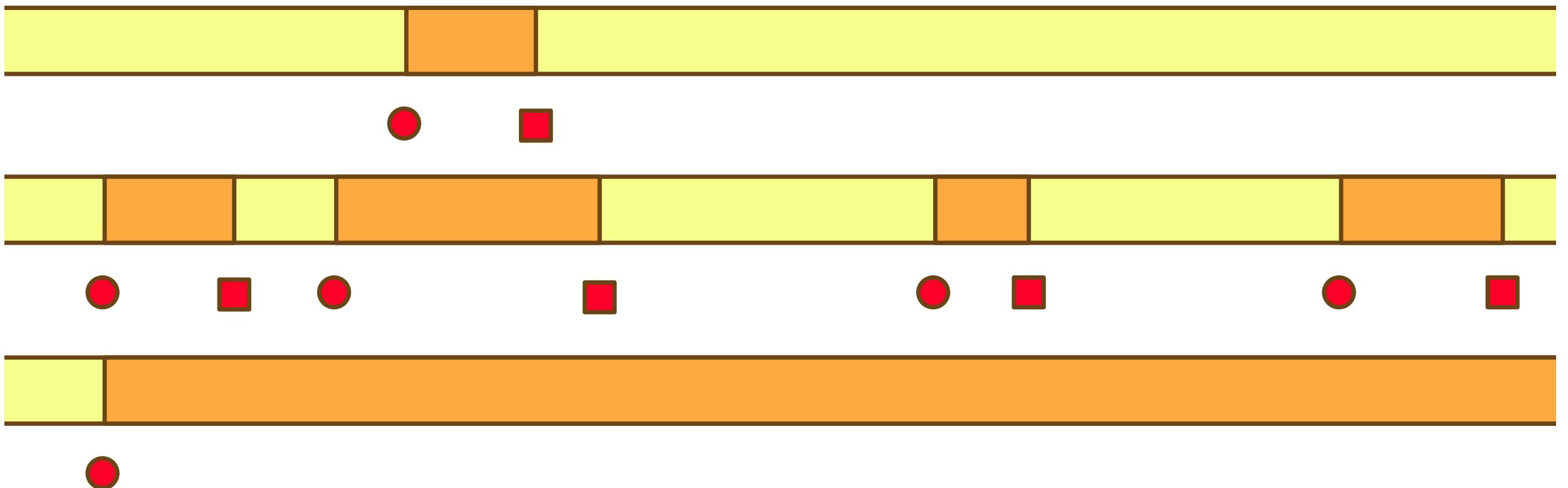


# Eventually...



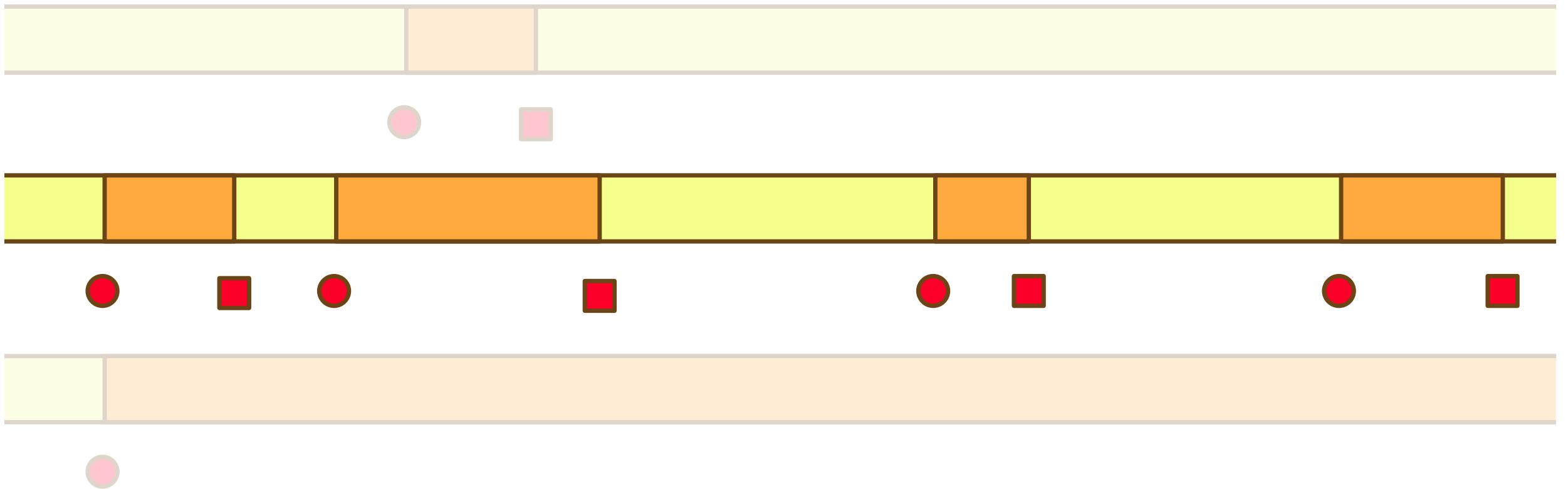


# Egocentric Video Understanding





# Egocentric Video Understanding





# It's Just Another Day: Unique Video Captioning by Discriminative Prompting

Toby Perrett, Tengda Han, Dima Damen, Andrew Zisserman



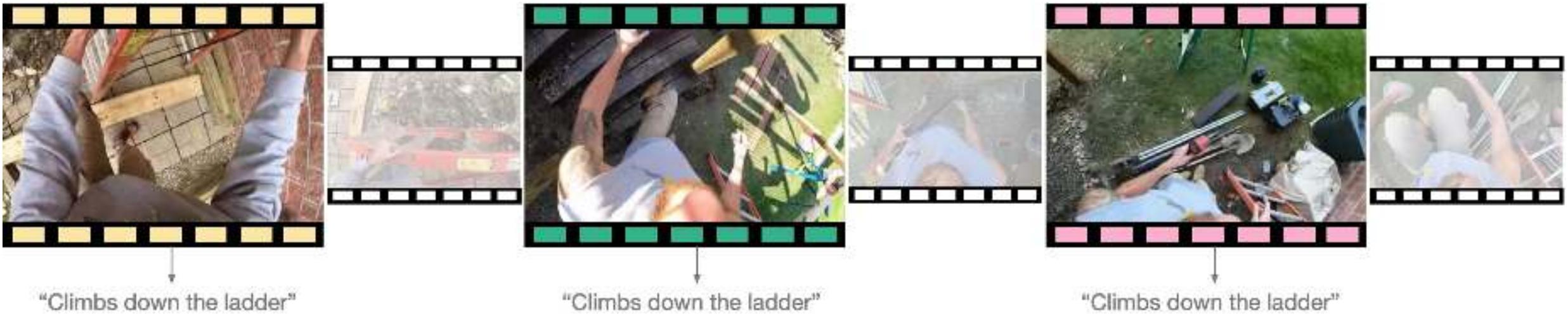
# Unique Video Captioning

with: Toby Perrett  
Tengda Han  
Andrew Zisserman

Life is repetitive...

# Unique Video Captioning

with: Toby Perrett  
Tengda Han  
Andrew Zisserman



- Current methods caption clips independently
- They generate the same caption for similar clips

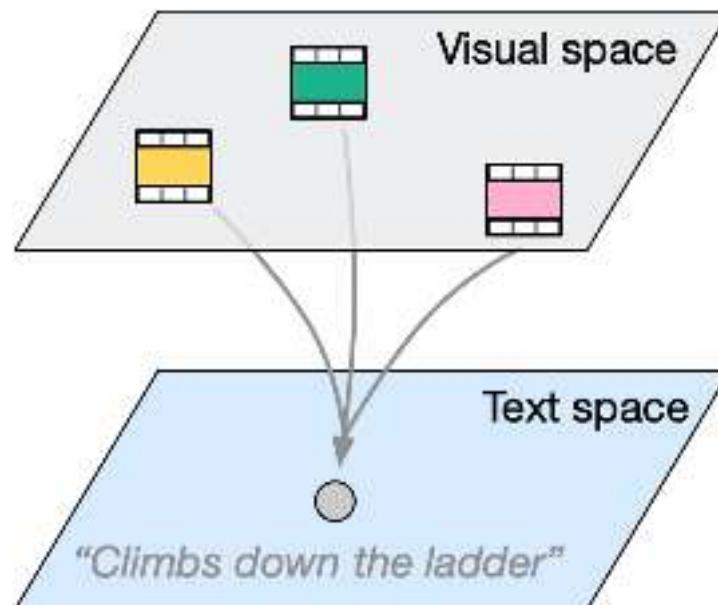
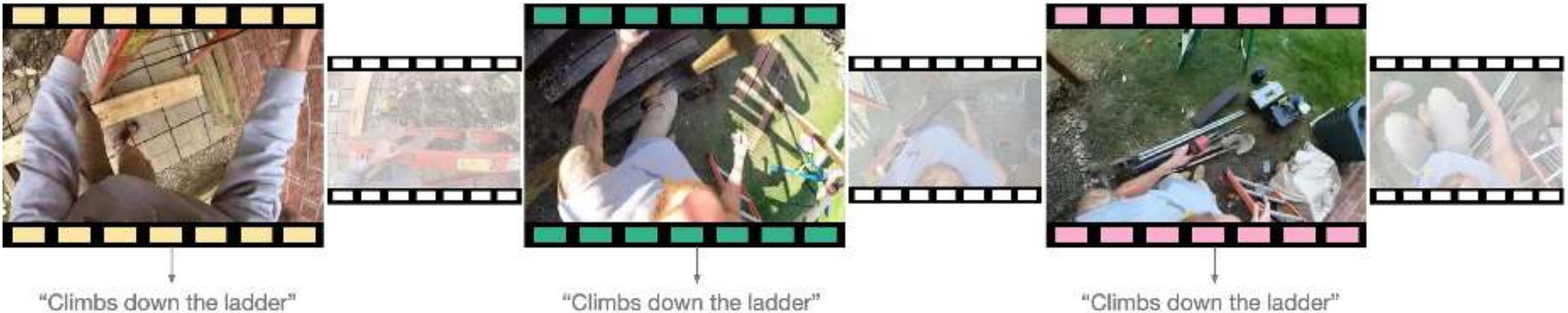
# Unique Video Captioning

with: Toby Perrett  
Tengda Han  
Andrew Zisserman

Goal:  
Generate a unique caption for every clip in a set

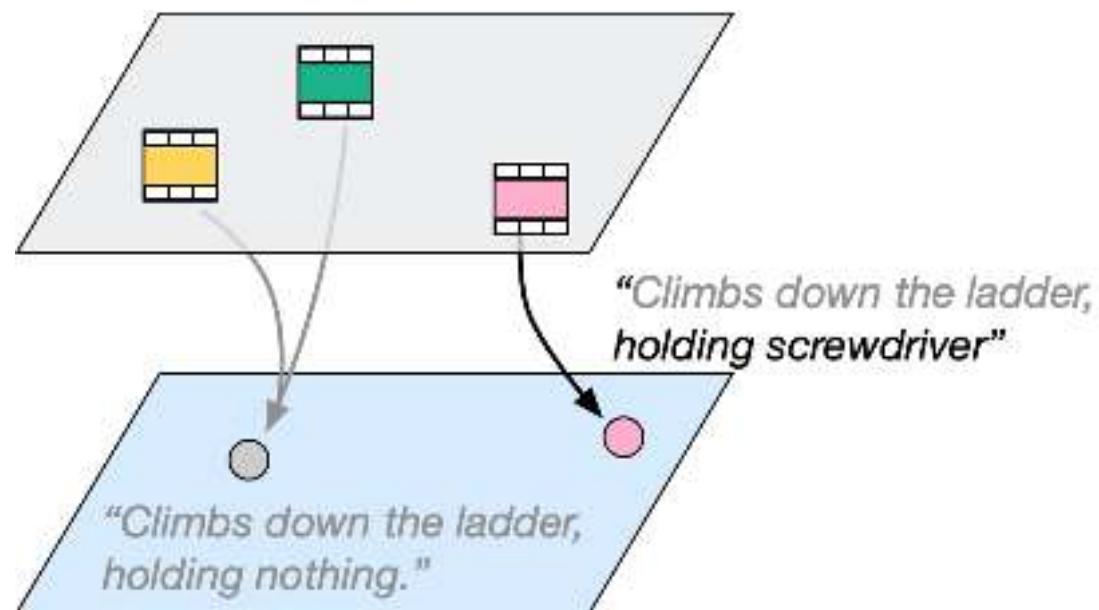
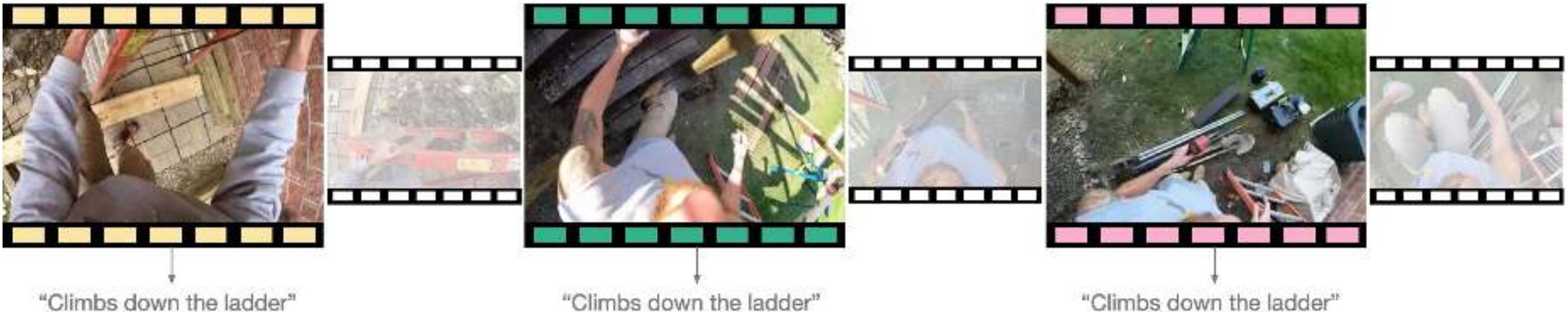
# Unique Video Captioning

with: Toby Perrett  
Tengda Han  
Andrew Zisserman



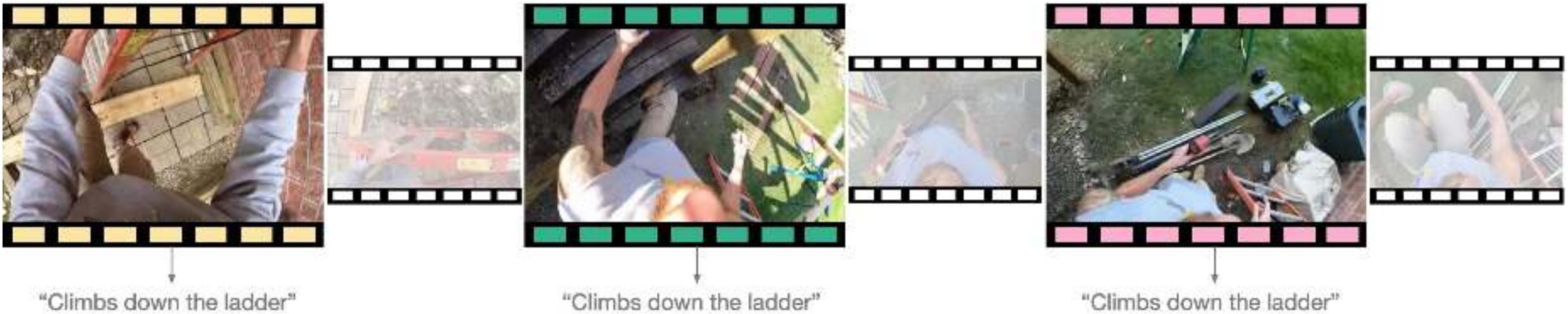
# Unique Video Captioning

with: Toby Perrett  
Tengda Han  
Andrew Zisserman



# Unique Video Captioning

with: Toby Perrett  
Tengda Han  
Andrew Zisserman



"Climbs down the ladder"

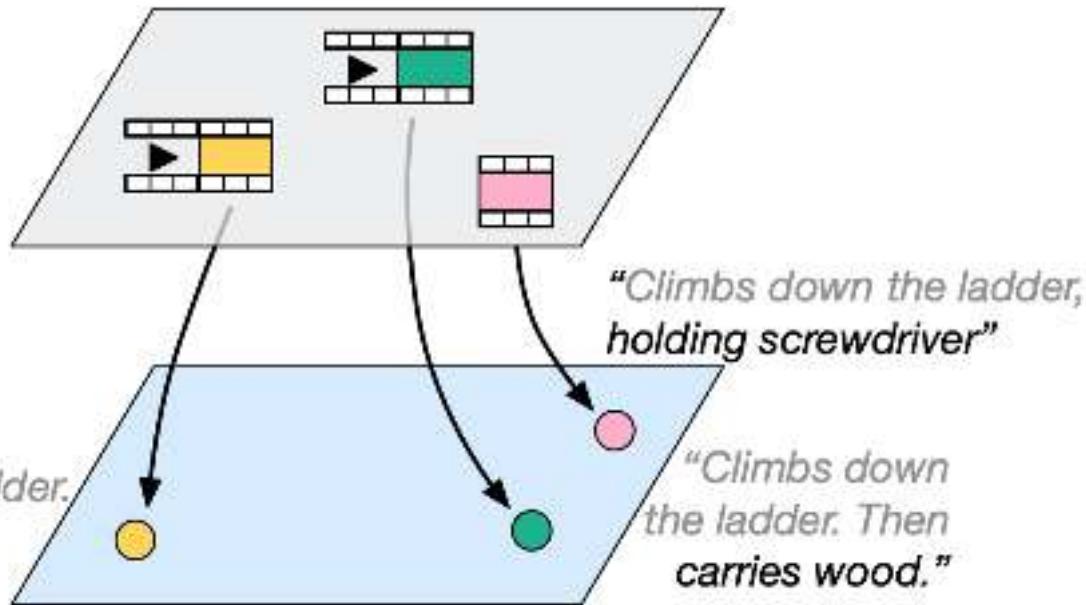
"Climbs down the ladder"

"Climbs down the ladder"

"Climbs down the ladder.  
Then **carries ladder.**"

"Climbs down the ladder,  
**holding screwdriver**"

"Climbs down  
the ladder. Then  
**carries wood.**"



# Captioning

with: Toby Perrett  
Tengda Han  
Andrew Zisserman



*Captioner*

*Response*



Climbs down a ladder

# Captioning

with: Toby Perrett  
Tengda Han  
Andrew Zisserman



*Discriminative prompt*

The person walks around

*Captioner*

*Response*

Climbs down a ladder  
and walks around  
a building site.

# Captioning by Discriminative Prompting

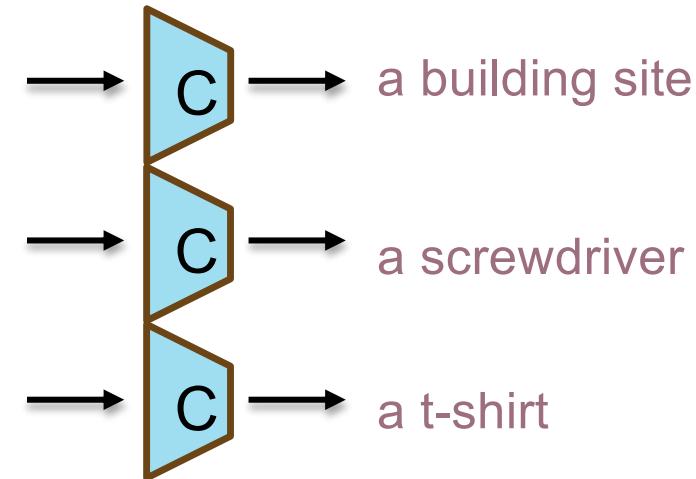
with: Toby Perrett  
Tengda Han  
Andrew Zisserman



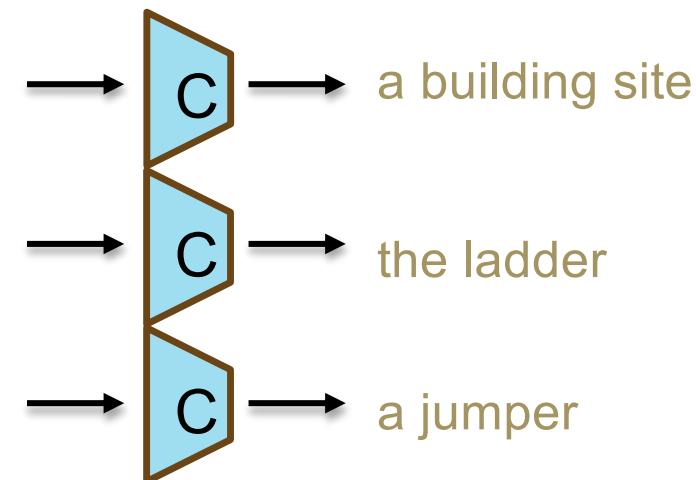
*Discriminative prompts*

The person walks around  
The person holds  
The person is wearing  
...

*Responses*

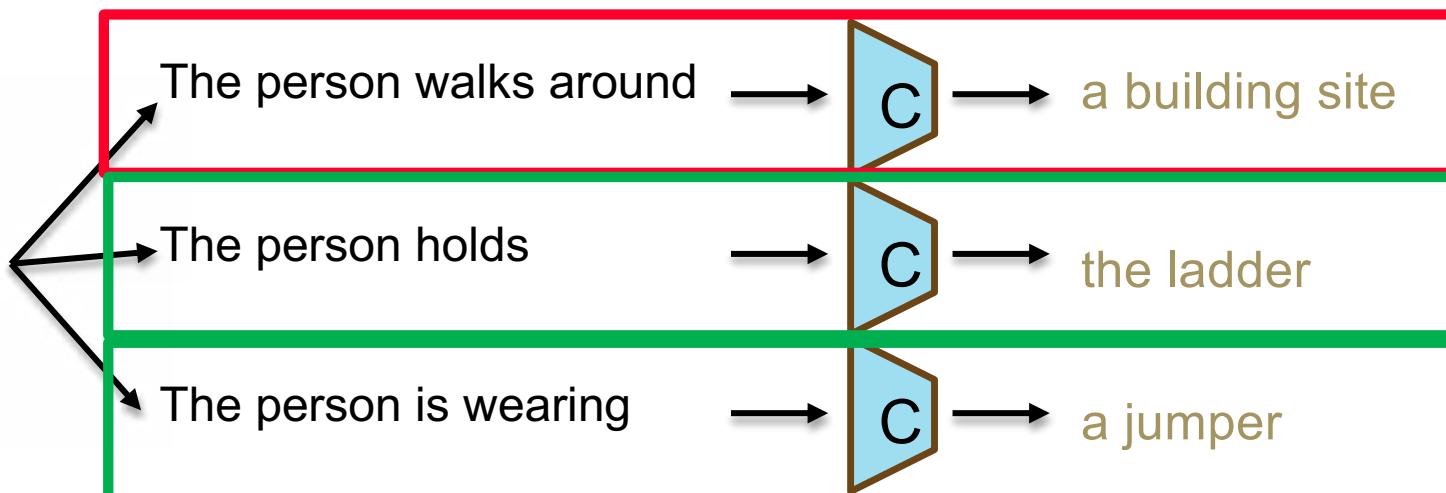
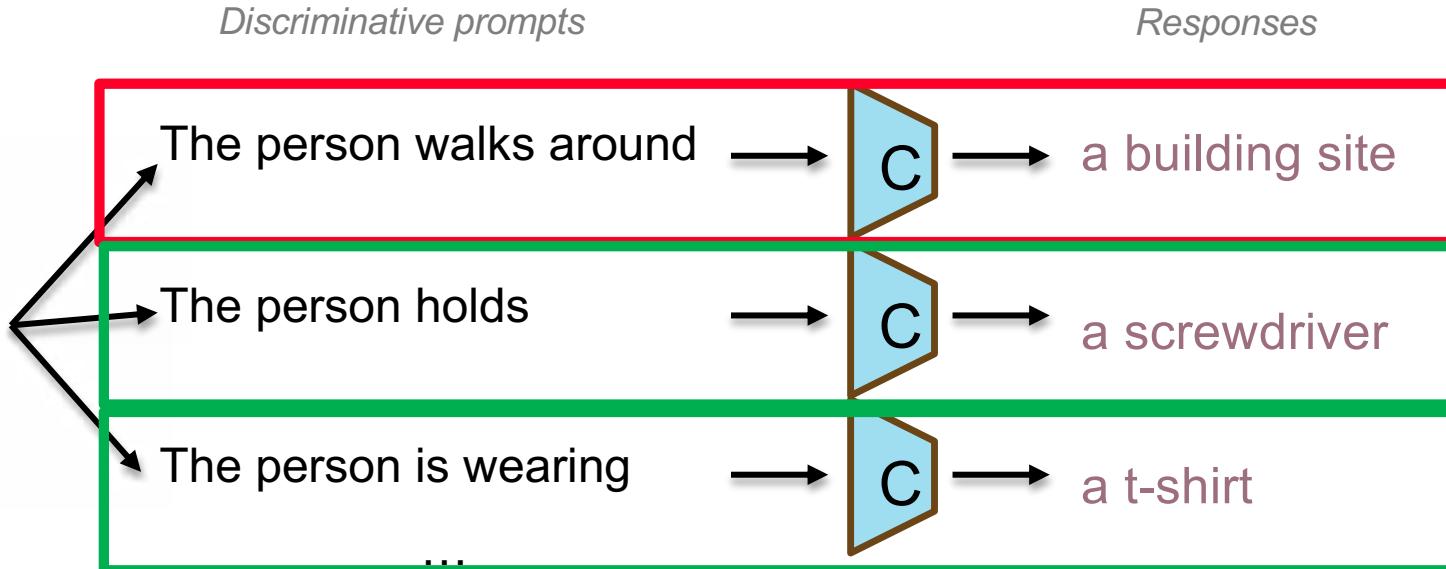


The person walks around  
The person holds  
The person is wearing  
...



# Captioning by Discriminative Prompting

with: Toby Perrett  
Tengda Han  
Andrew Zisserman



# Captioning by Discriminative Prompting

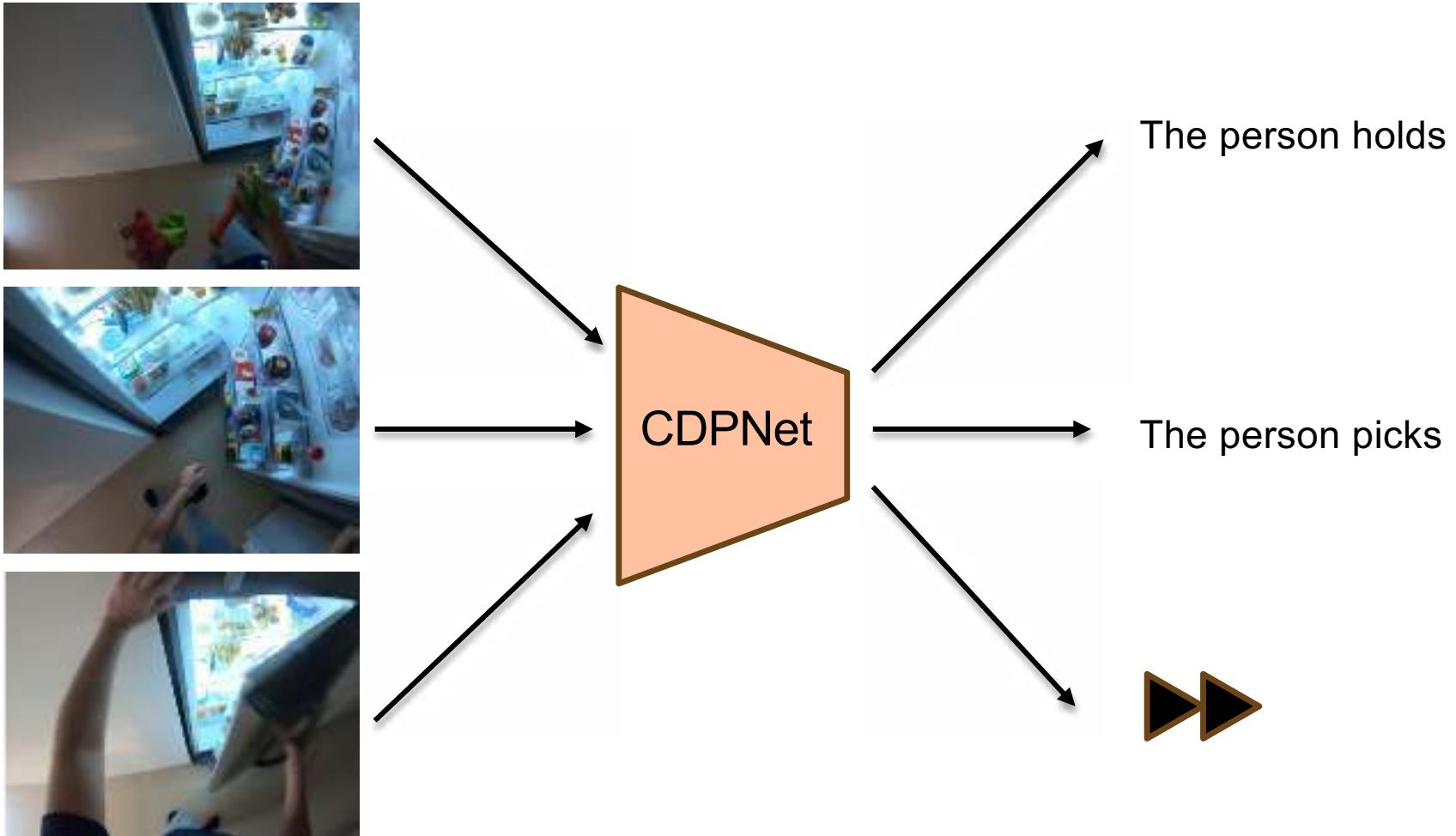
with: Toby Perrett  
Tengda Han  
Andrew Zisserman

We propose to...  
consider clips jointly  
use a bank of discriminative prompts

But...  
Expensive £££  
What if there isn't a suitable prompt?

# Captioning by Discriminative Prompting

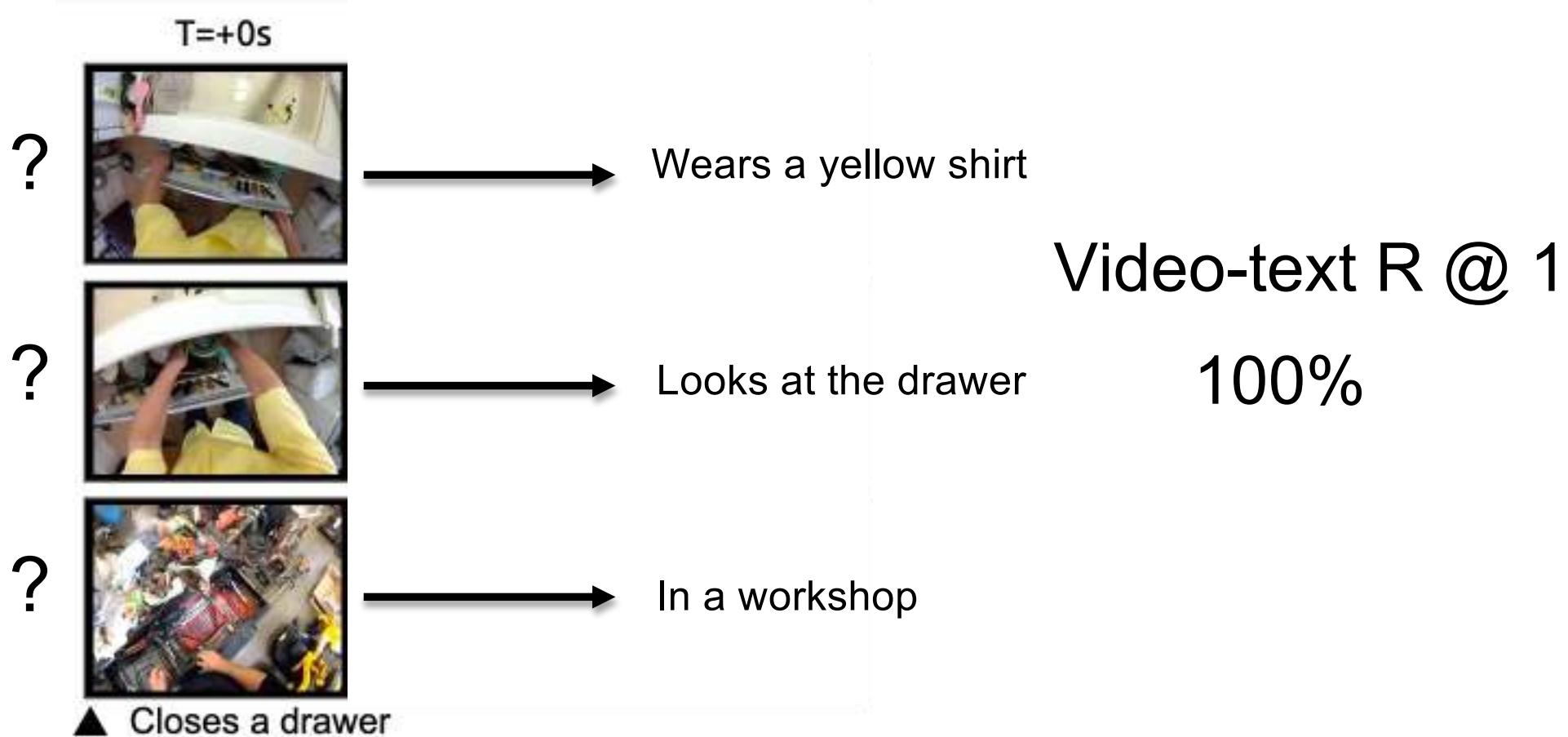
with: Toby Perrett  
Tengda Han  
Andrew Zisserman



# Benchmarks

with: Toby Perrett  
Tengda Han  
Andrew Zisserman

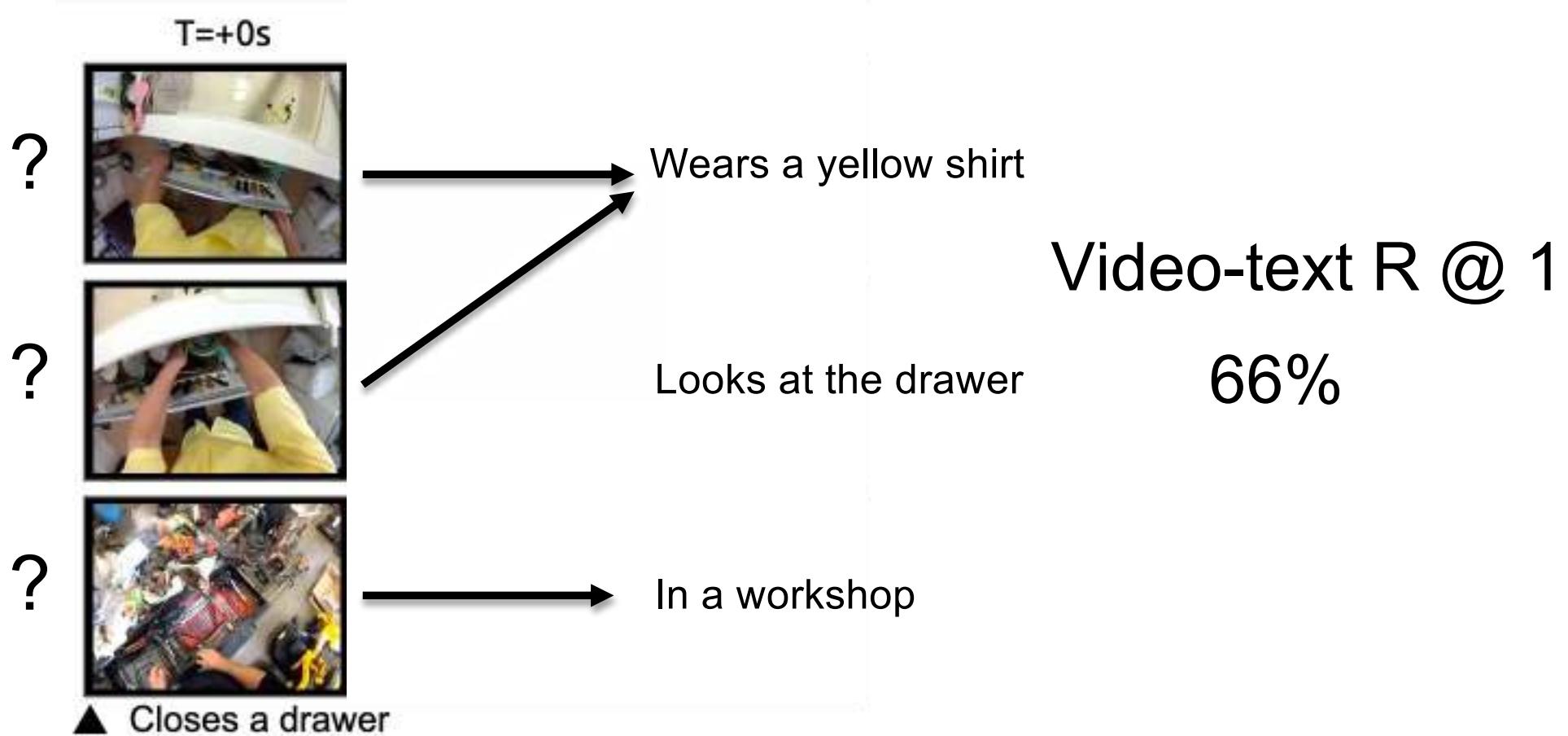
Task: Caption every clip, then evaluate retrieval.



# Benchmarks

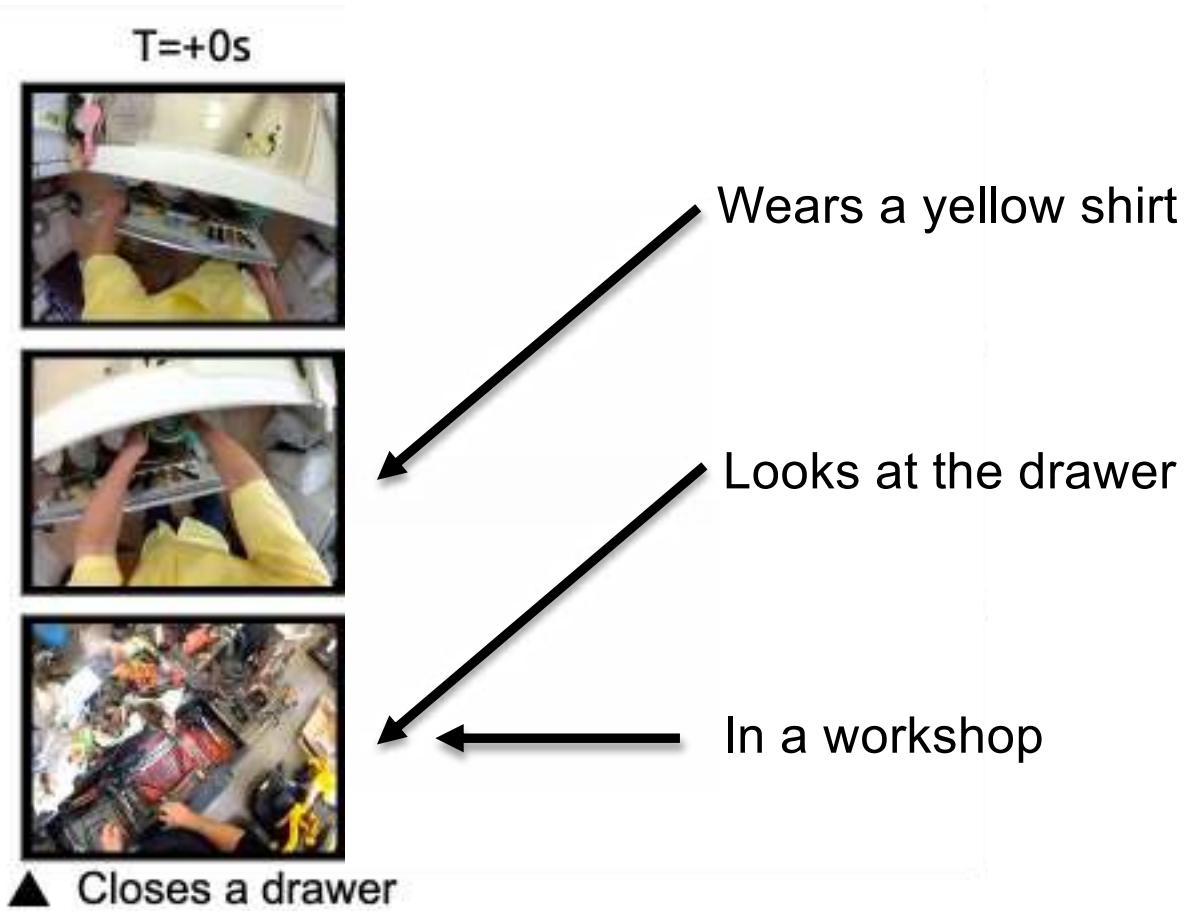
with: Toby Perrett  
Tengda Han  
Andrew Zisserman

Task: Caption every clip, then evaluate retrieval.



# Benchmarks

with: Toby Perrett  
Tengda Han  
Andrew Zisserman



$$\begin{aligned} \text{Avg Recall @ 1} \\ &= (\text{Video-text} + \text{Text-video}) / 2 \\ &= (33 + 66) / 2 \\ &= 30\% \end{aligned}$$

# Unique Video Captioning

with: Toby Perrett  
Tengda Han  
Andrew Zisserman

Average recall @ 1

<b>Egocentric</b>	<b>+0s</b>
LaViLa	37
LaViLa + CDP	45

# Unique Video Captioning

with: Toby Perrett  
Tengda Han  
Andrew Zisserman

 Climbs the stairs

# Unique Video Captioning

with: Toby Perrett  
Tengda Han  
Andrew Zisserman

🔍 Climbs the stairs



Climbs the stairs and  
holds the phone



Climbs the stairs and  
picks up the drill



Climbs the stairs and  
holds a tape measure

# Unique Video Captioning

with: Toby Perrett  
Tengda Han  
Andrew Zisserman



Looks around the shelves

# Unique Video Captioning

with: Toby Perrett  
Tengda Han  
Andrew Zisserman



Looks around the shelves



Looks around the shelves and  
the other man picks up a packet  
of biscuits from the shelf with his  
left hand

Looks around the shelves and  
looks at the list

Looks around the shelves and  
then  
picks up a packet of cough rubs

# Unique Video Captioning

with: Toby Perrett  
Tengda Han  
Andrew Zisserman

The screenshot shows a web browser window with the URL <https://tobyperrett.github.io/its-just-another-day/>. The main content is a presentation slide with the following text:

## It's Just Another Day: Unique Video Captioning by Discriminative Prompting

ACCV 2024 Oral

Toby Perrett, Tengda Han, Dima Damen, Andrew Zisserman

[arXiv](#) | [Code/Benchmark](#)

---

### Introduction

This paper investigates unique video captioning. We introduce a method, Captioning by Discriminative Prompting (CDP) and challenging unique captioning benchmarks on Egocentric video and Timeloop movies.

---

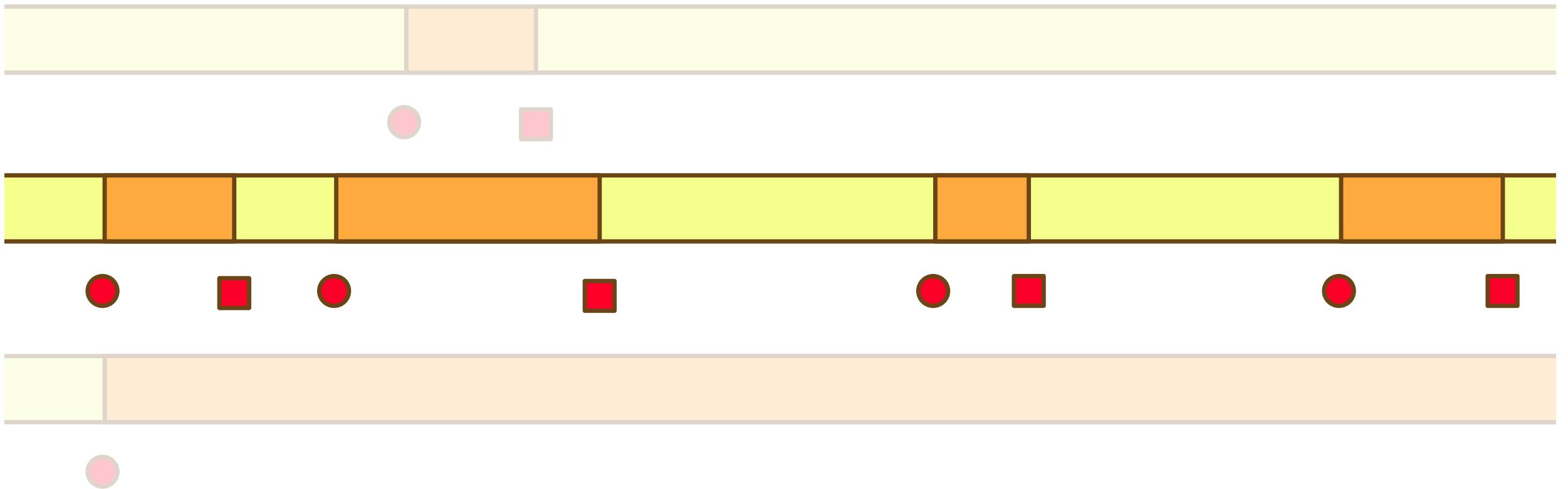
### Problem statement

The following figures highlight a shortcoming of current captioning approaches. They caption each clip independently, giving similar captions for similar clips. First, in a timeloop movie:

A film strip consisting of three frames. The first frame shows a man looking up with the caption "A man wakes up". The second frame shows the same man looking up with the caption "A man sits up". The third frame shows the same man looking up with the caption "A man wakes up".



# Egocentric Video Understanding



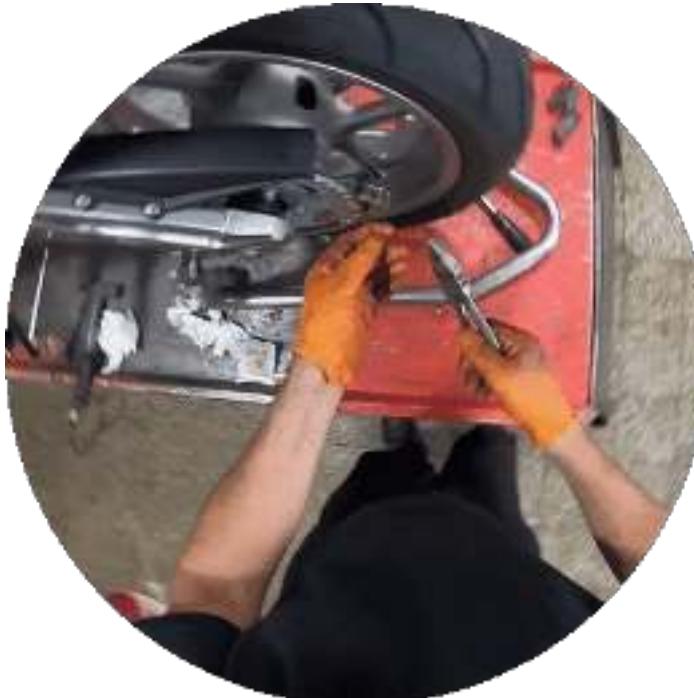
# Generalisation across Scenarios and Locations

with: Chiara Plizzari  
Toby Perrett



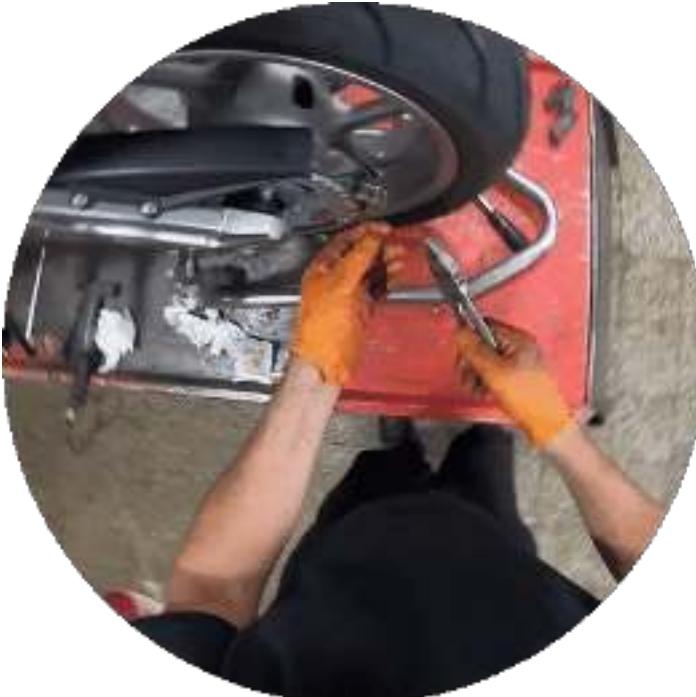
# Generalisation across Scenarios and Locations

with: Chiara Plizzari  
Toby Perrett



# Generalisation across Scenarios and Locations

with: Chiara Plizzari  
Toby Perrett



# Generalisation across Scenarios and Locations

with: Chiara Plizzari  
Toby Perrett



# Generalisation across Scenarios and Locations

with: Chiara Plizzari  
Toby Perrett



# Generalisation across Scenarios and Locations

with: Chiara Plizzari  
Toby Perrett



# Generalisation across Scenarios and Locations

with: Chiara Plizzari  
Toby Perrett



# Dataset: ARGO1M

with: Chiara Plizzari  
Toby Perrett

- We introduce **ARGO1M**, the first dataset to perform **Action Recognition Generalisation** Over Scenarios and Locations



# Dataset: ARGO1M

with: Chiara Plizzari  
Toby Perrett

- We introduce **ARGO1M**, the first dataset to perform **Action Recognition Generalisation** Over Scenarios and Locations



# Dataset: ARGO1M

with: Chiara Plizzari  
Toby Perrett

- We introduce **ARGO1M**, the first dataset to perform **Action Recognition Generalisation** Over Scenarios and Locations

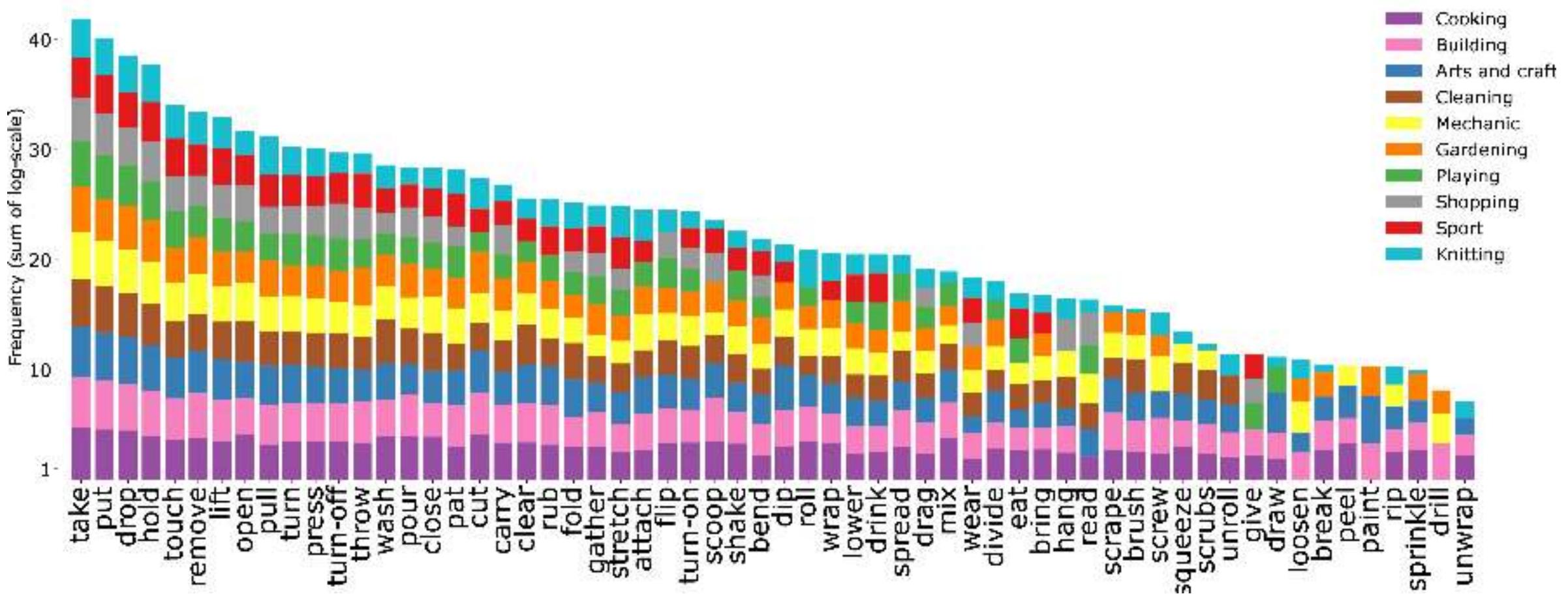
**1.1M samples**



# Generalisation across Scenarios and Locations

with: Chiara Plizzari  
Toby Perrett

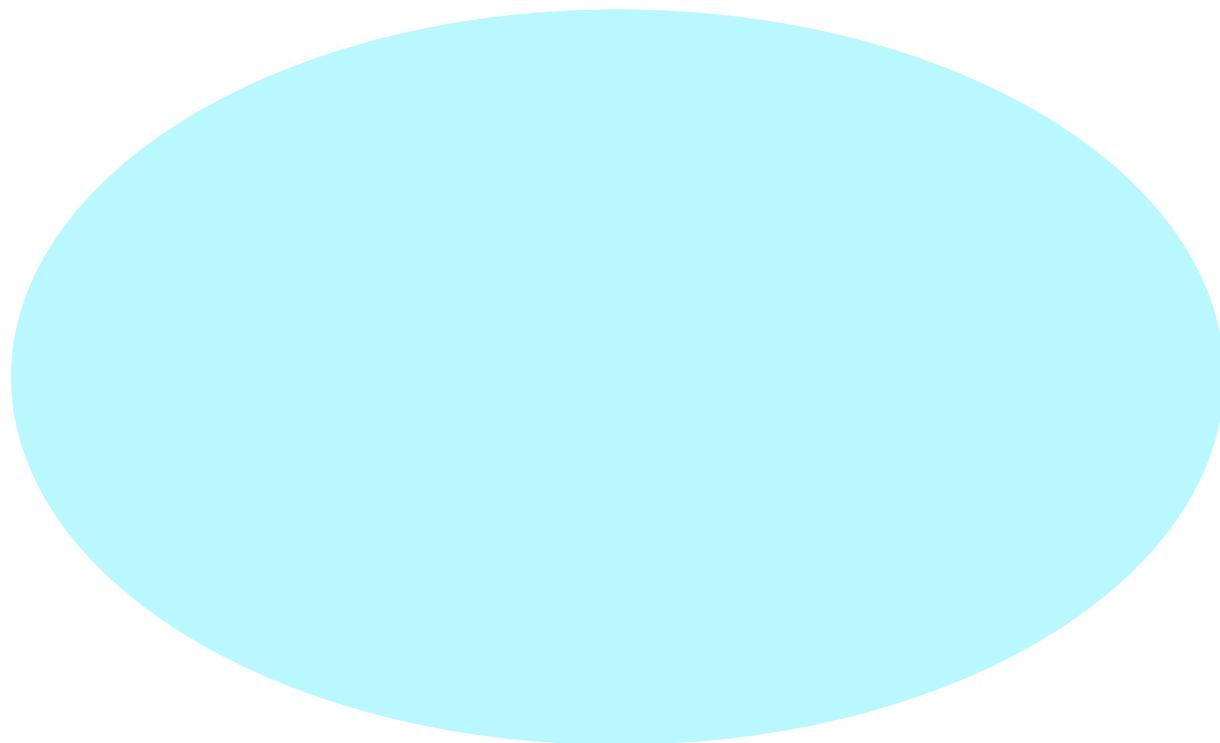
ARGO1M: 1.05M action clips from 60 action classes recorded in 13 locations within 10 scenarios



# ARGO1M Splits

with: Chiara Plizzari  
Toby Perrett

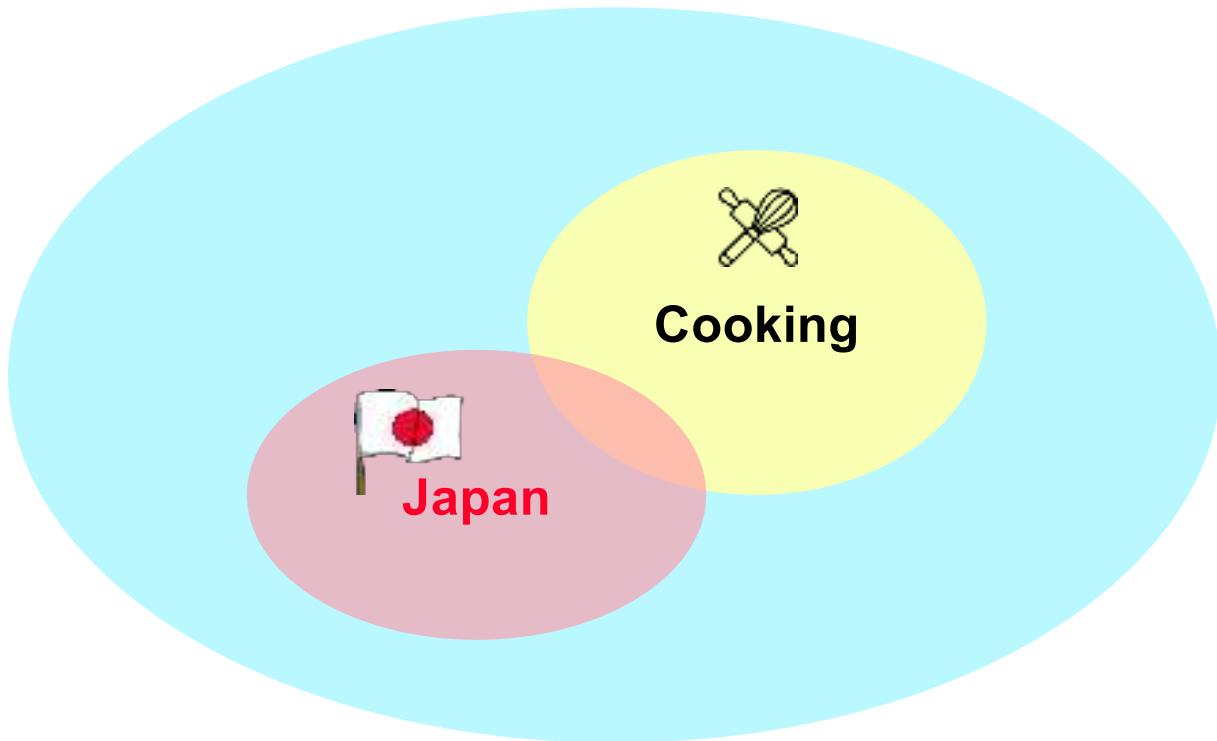
## ARGO1M



# ARGO1M Splits

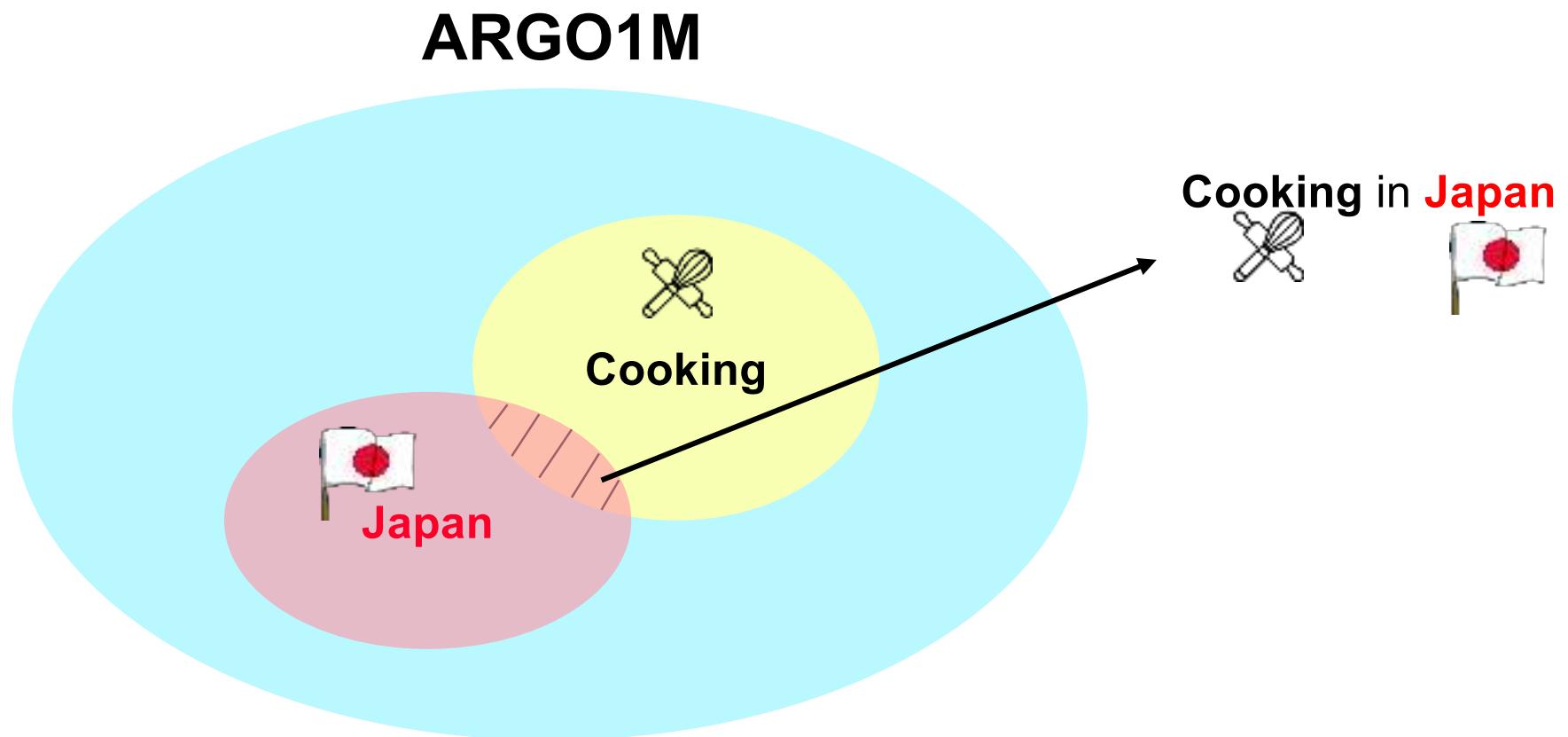
with: Chiara Plizzari  
Toby Perrett

## ARGO1M



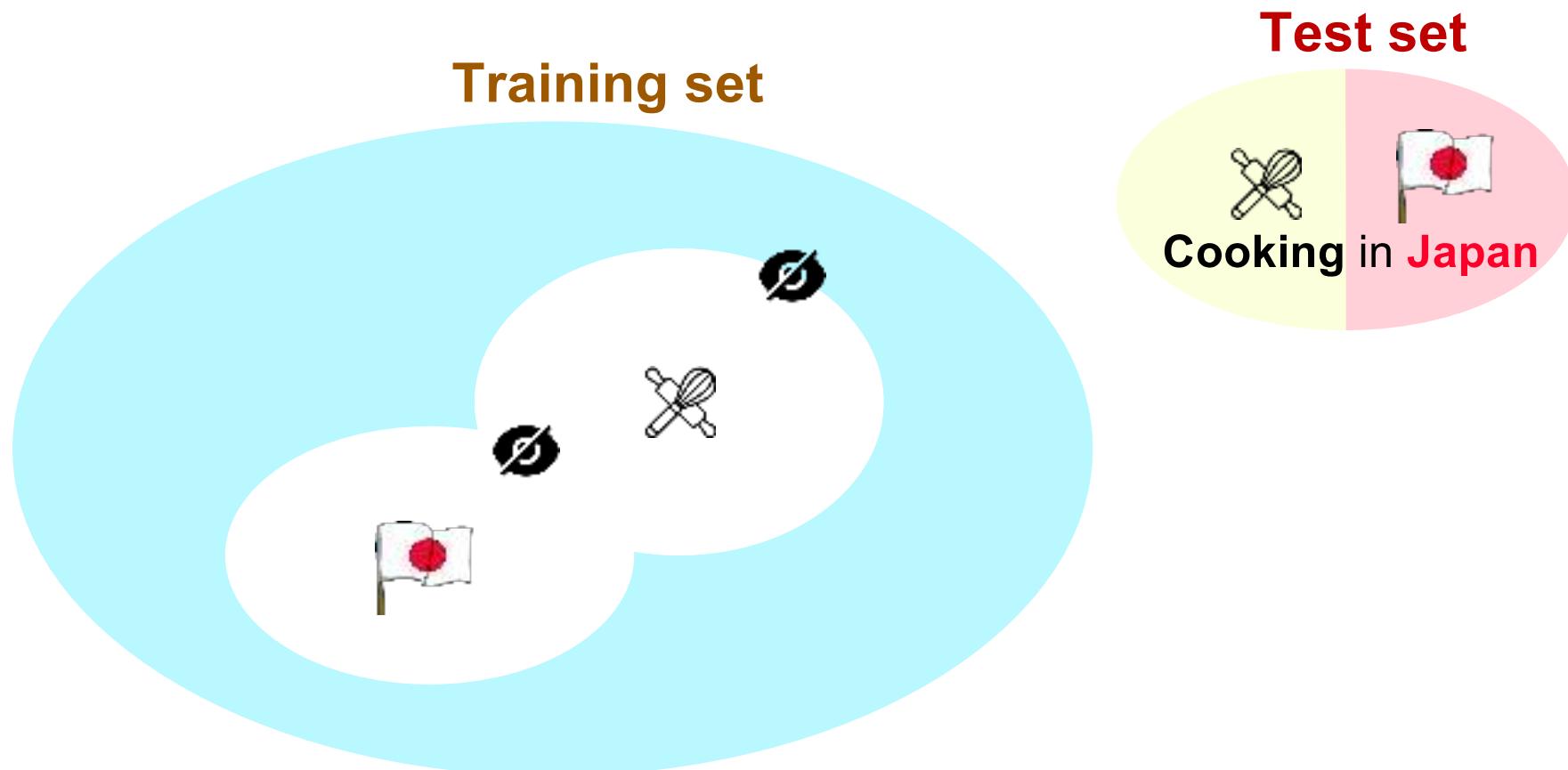
# ARGO1M Splits

with: Chiara Plizzari  
Toby Perrett



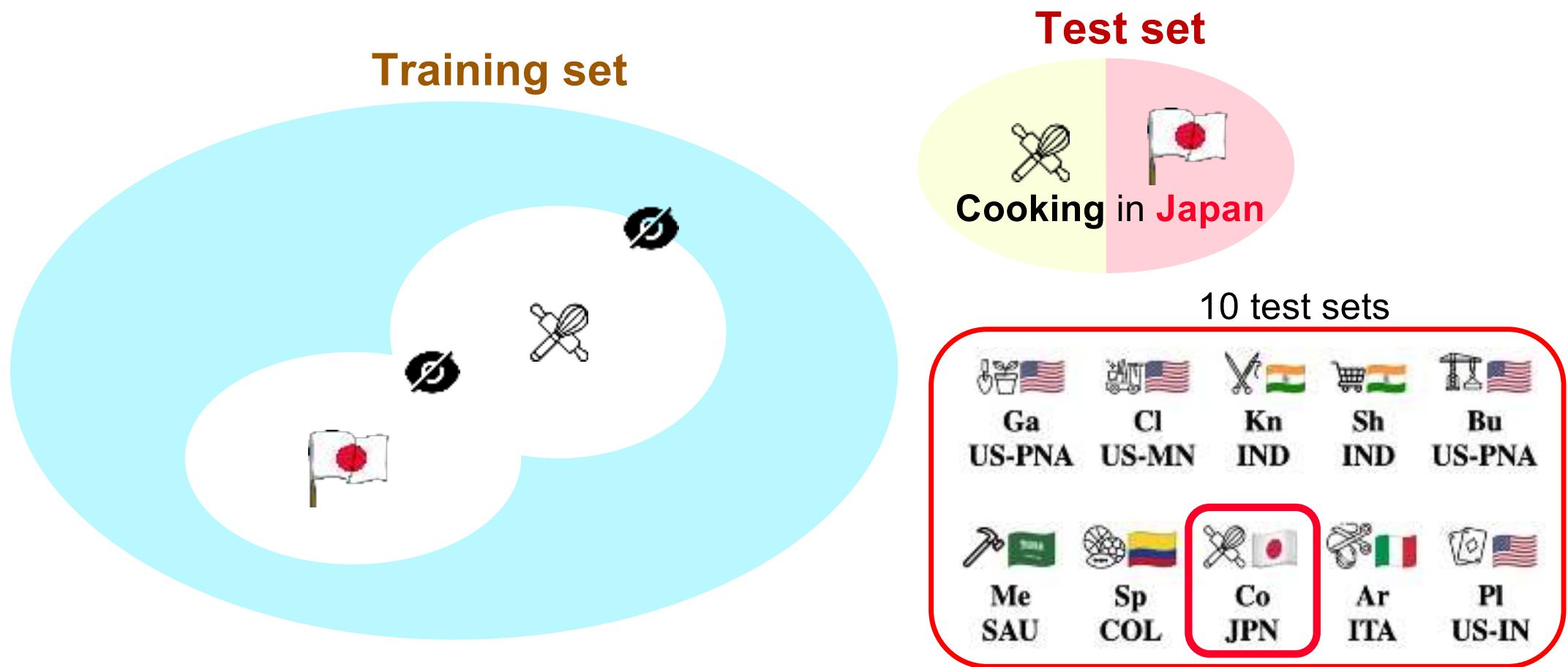
# ARGO1M Splits

with: Chiara Plizzari  
Toby Perrett



# ARGO1M Splits

with: Chiara Plizzari  
Toby Perrett



# Generalisation across Scenarios and Locations

with: Chiara Plizzari  
Toby Perrett

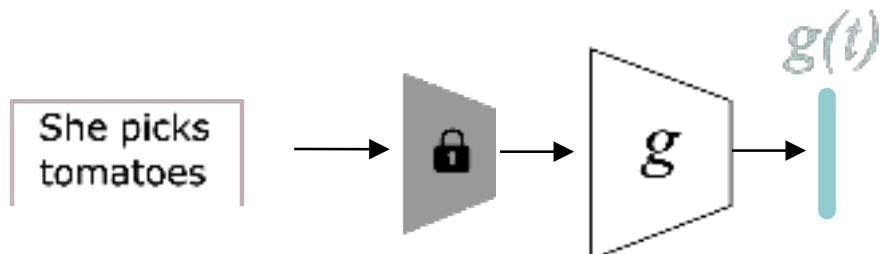
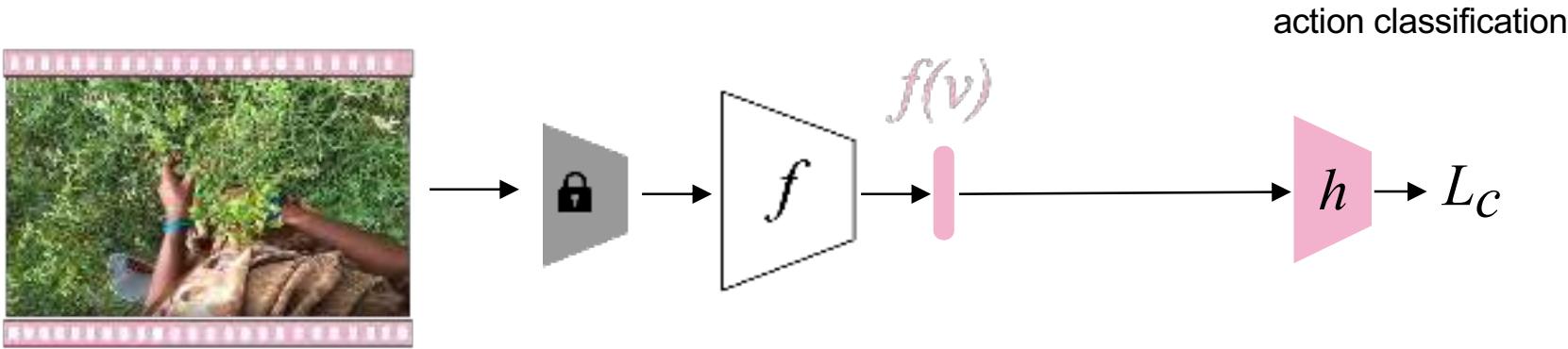


*He cuts the lemon strand*



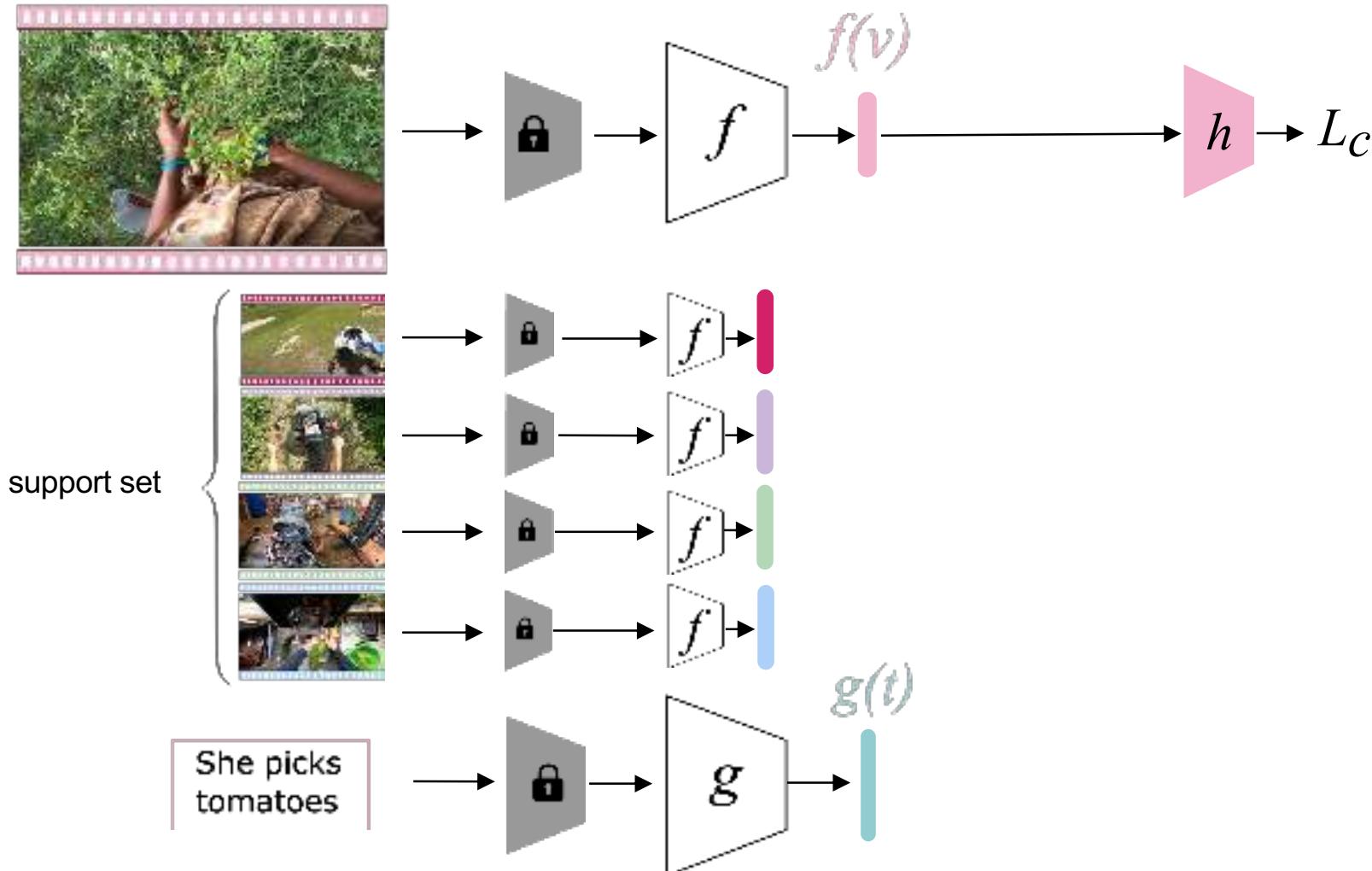
# Proposed method: CIR

with: Chiara Plizzari  
Toby Perrett



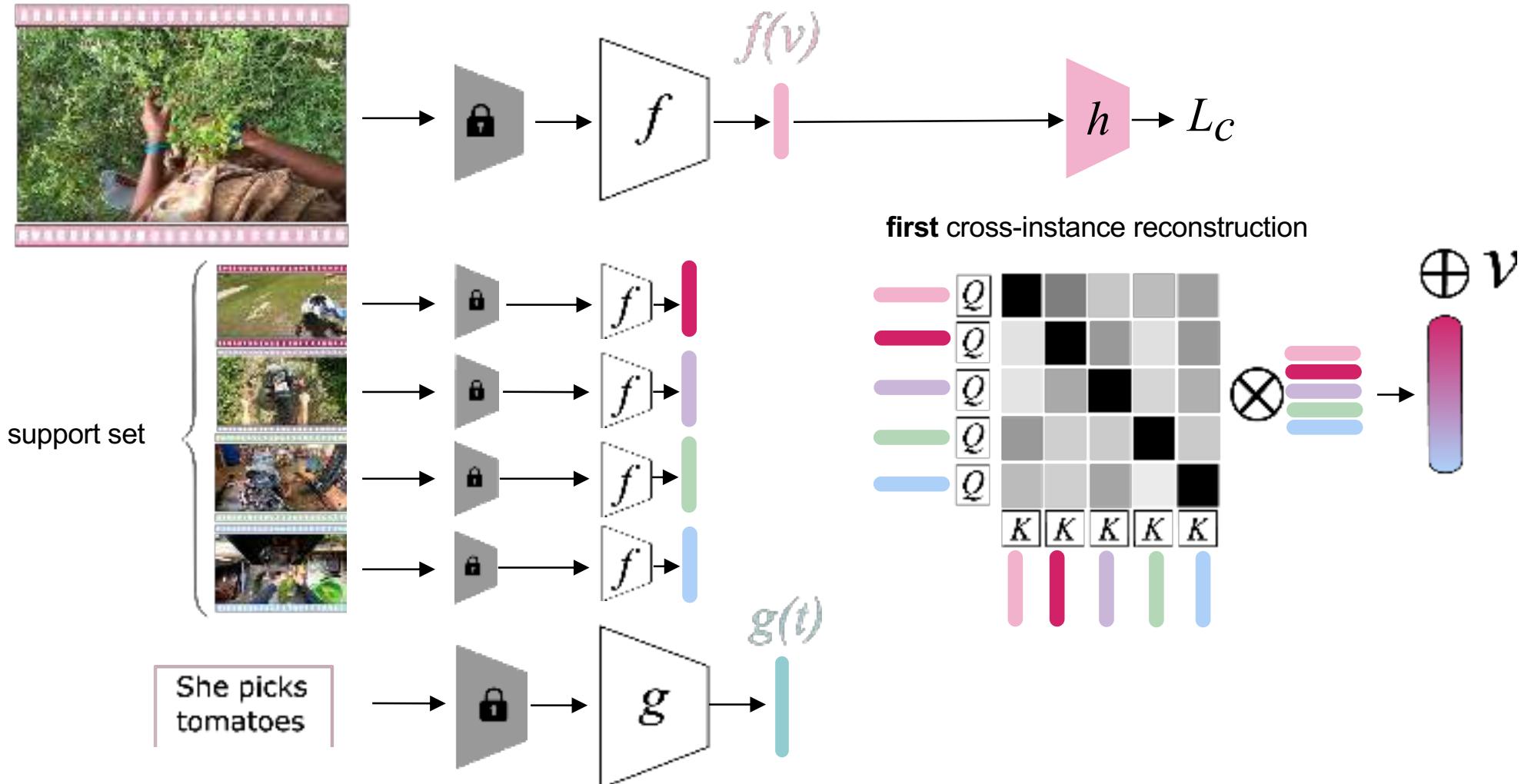
# Proposed method: CIR

with: Chiara Plizzari  
Toby Perrett



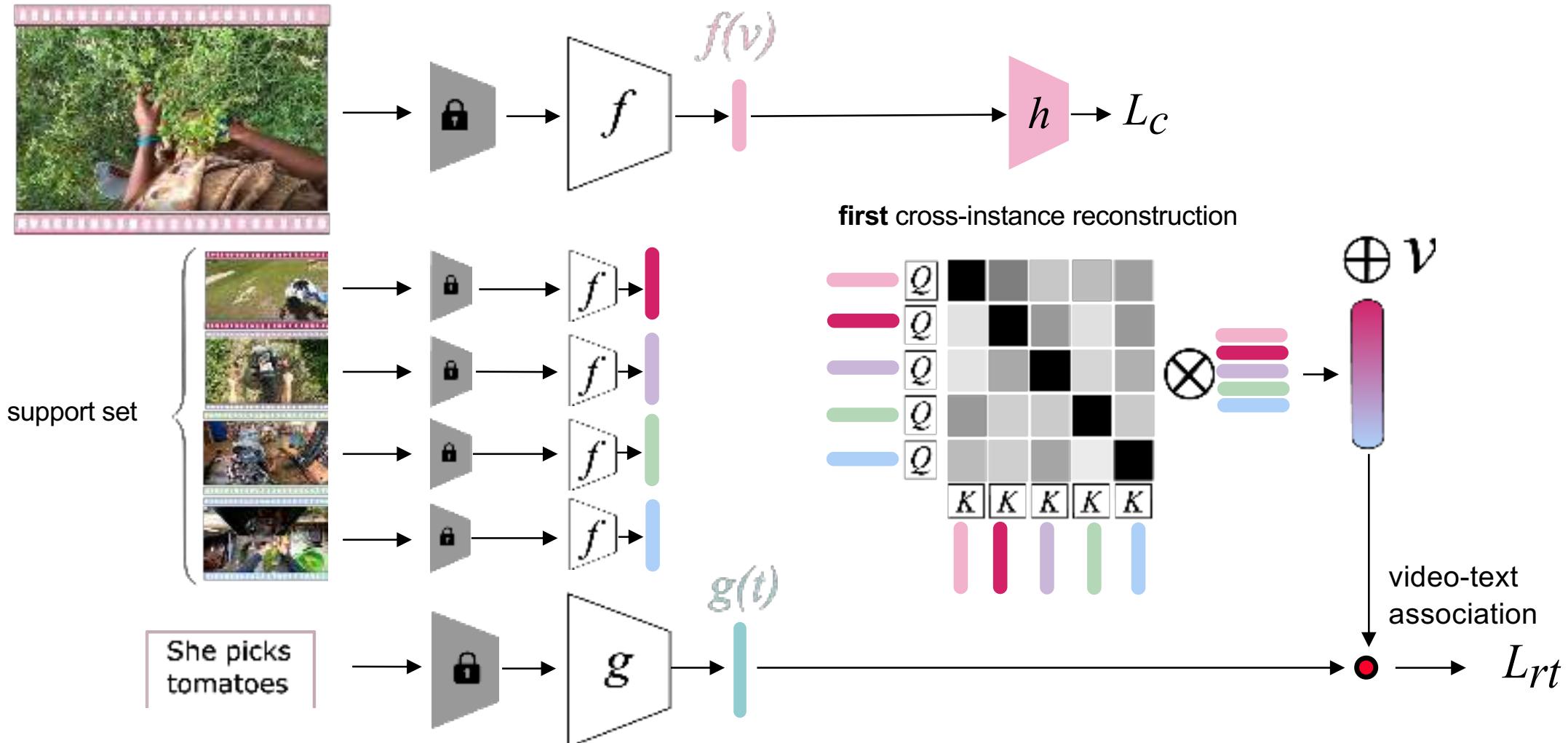
# Proposed method: CIR

with: Chiara Plizzari  
Toby Perrett



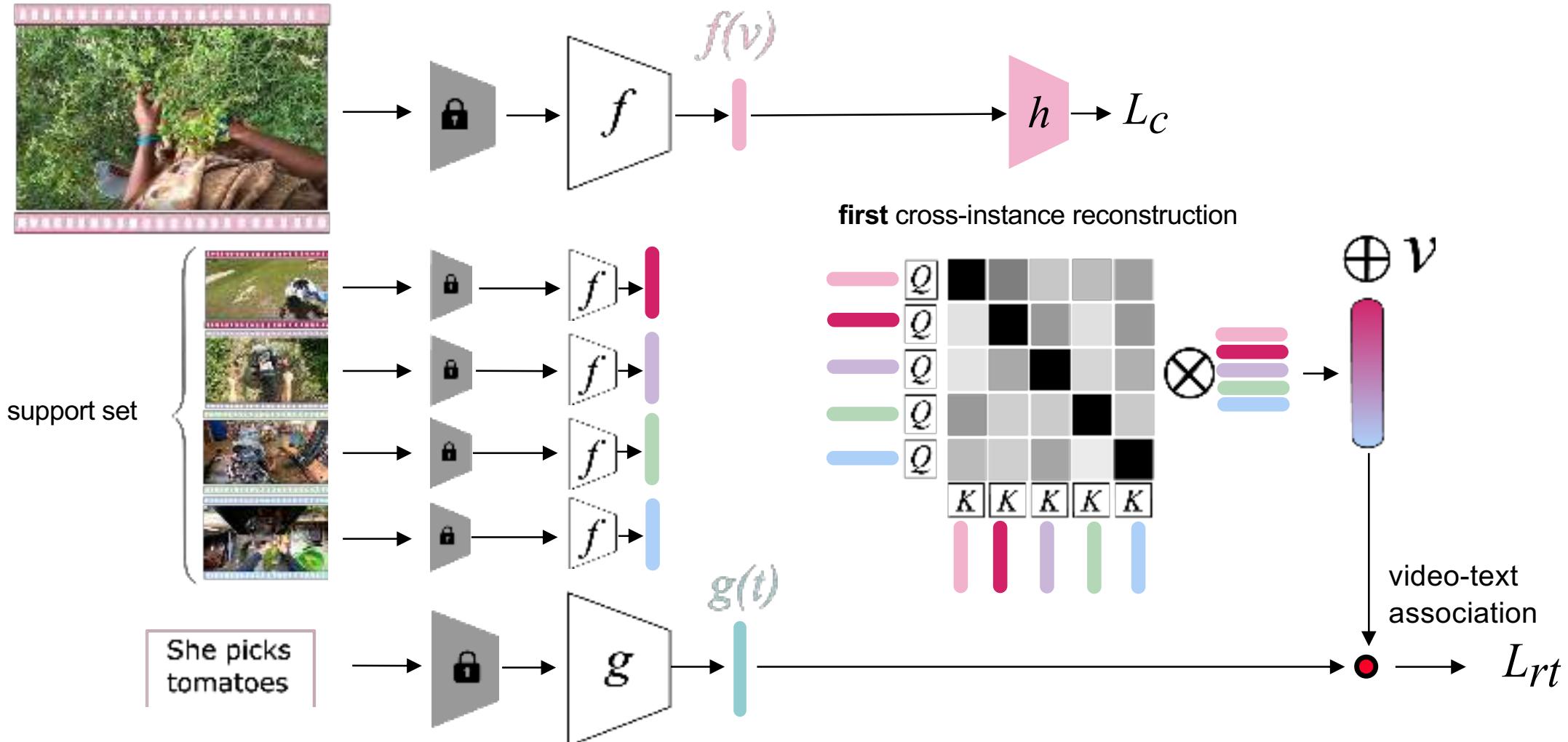
# Proposed method: CIR

with: Chiara Plizzari  
Toby Perrett



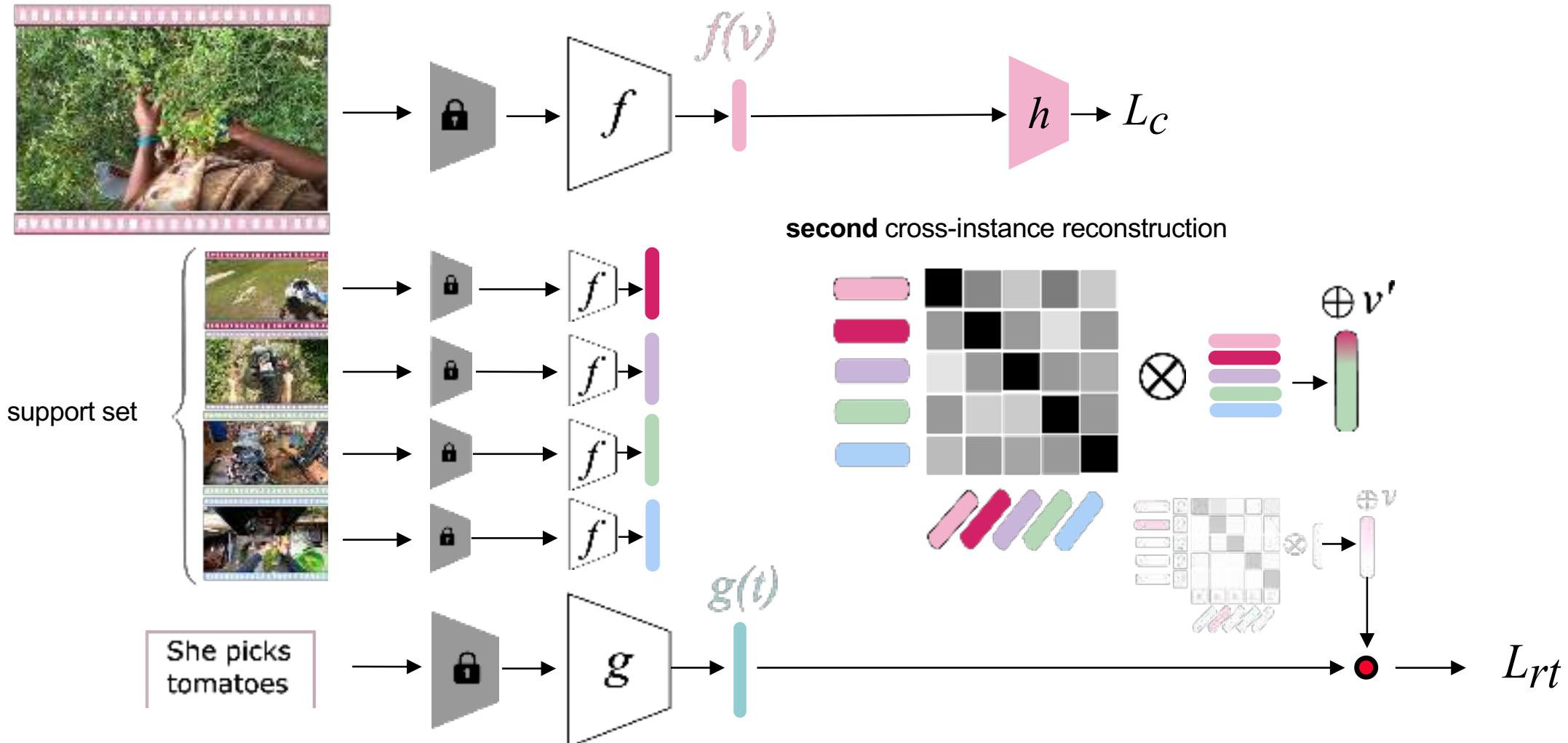
# Proposed method: CIR

with: Chiara Plizzari  
Toby Perrett



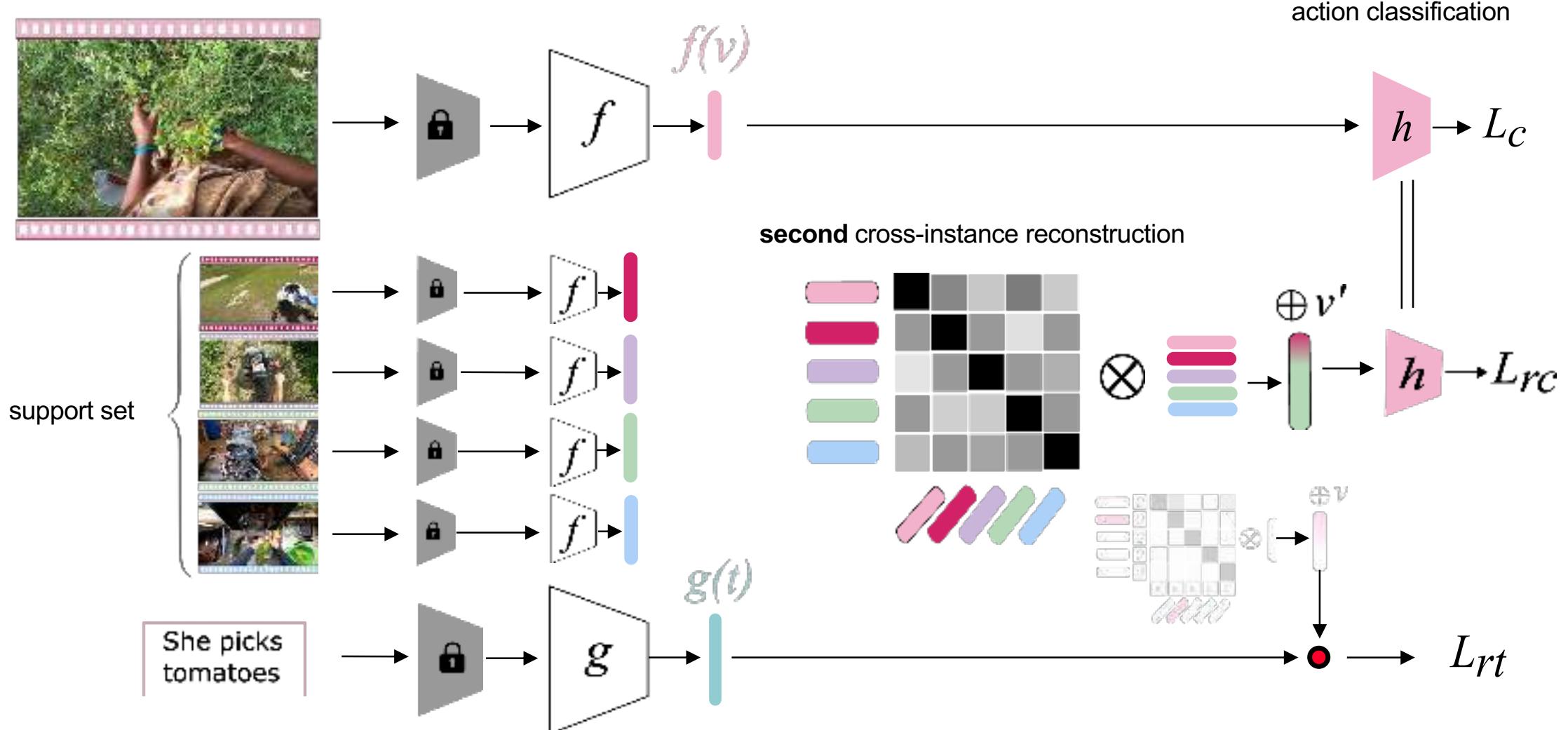
# Proposed method: CIR

with: Chiara Plizzari  
Toby Perrett



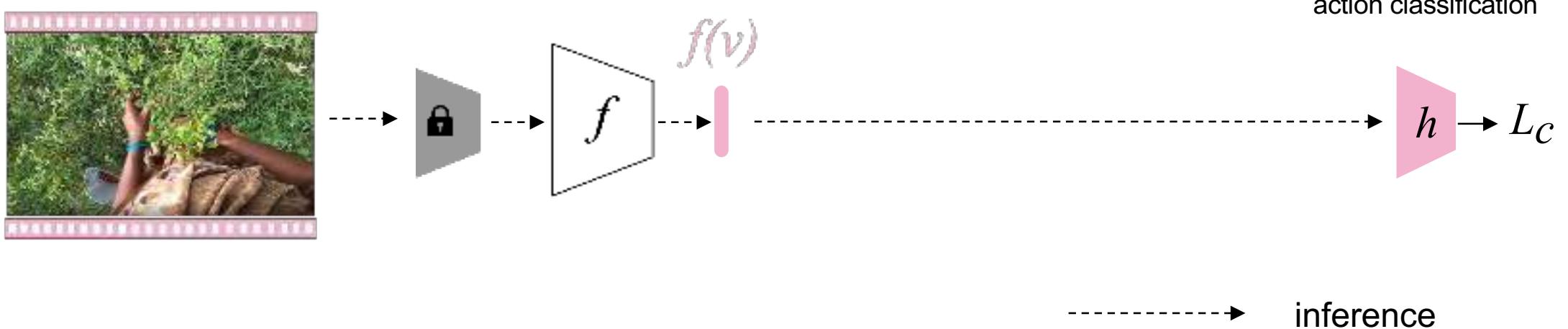
# Proposed method: CIR

with: Chiara Plizzari  
Toby Perrett



# Proposed method: CIR

with: Chiara Plizzari  
Toby Perrett



# Examples

Chiara Plizzari  
Toby Perrett  
Dima Damen

#C C drops the cut vegetables



query



support 1

support 2

support 3

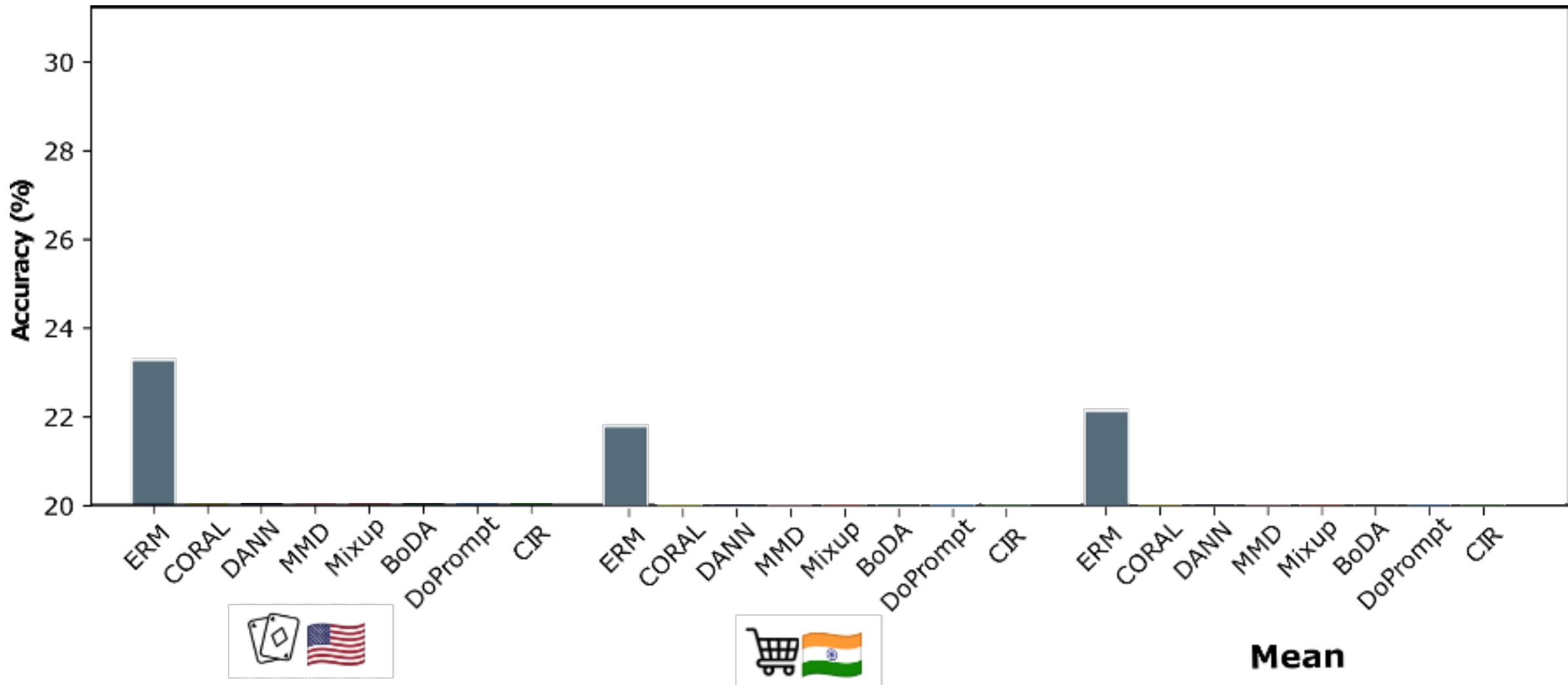
support 4

support 5



# Proposed method: CIR

with: Chiara Plizzari  
Toby Perrett



# What can a cook in Italy teach a mechanic in India?

with: Chiara Plizzari  
Toby Perrett

## What can a cook in Italy teach a mechanic in India? Action Recognition Generalisation Over Scenarios and Locations

Chiara Plizzari\*

Toby Perrett\*

Barbara Caputo\*

Dima Damen\*

\* Politecnico di Torino, Italy

\* University of Bristol, United Kingdom

### Abstract

We propose and address a new generalisation problem: can a model trained for action recognition successfully classify actions when they are performed within a previously unseen scenario and in a previously unseen location? To answer this question, we introduce the Action Recognition Generalisation Over scenarios and locations dataset (ARGO1M), which contains 1.1M video clips from the large-scale Egoid dataset, across 10 scenarios and 10 locations. We demonstrate recognition models struggle to generalise over 10 prepared test splits, each of an unseen scenario in an unseen location. We then propose CIR, a method to represent each video as a Cross-Domain Representation of videos from other domains. Reconstructions are paired with test scenarios to guide the learning of a domain generalisable representation. We provide extensive analysis and ablations on ARGO1M that show CIR outperforms prior domain generalisation methods all test splits. *Code and data:* <https://github.com/plizzari/argo1m/>, <https://tinyurl.com/argo1m>.



Figure 1: Problem statement and samples from the ARGO1M dataset. The same action e.g. “cut”, is performed differently based on the scenario and the location in which it is carried out. We aim to generalise so as to recognise the same action within a new scenario, unseen during training, and in an unseen location, e.g., Mechanic (), Arts ().

### 1. Introduction

A notable distinction between human and machine intelligence is the ability of humans to generalise. We can use an example of the action “cut” performed by a cook in Italy, and recognise the same action performed in a different geographical location, e.g. India, despite having never visited. We can also recognise actions within new scenarios, such as a mechanic cutting metal, even if we are unfamiliar with the tools they use.

This problem is known as domain generalisation [50], where a model trained on a set of labelled data fails to generalise to a different distribution in inference. The gap between distributions is known as domain shift. To date, work have focused on generalising over visual-domain shifts [22, 49, 14, 43, 31]. In this paper, we introduce the scenario shift, where the same action is performed as part

\*Work carried during Chiara's research visit to the University of Bristol



## ARGO1M Dataset CIR Method Code and Models

RELEASED



# HD-EPIC: A Highly-Detailed Egocentric Video Dataset



Toby Perrett



Ahmad Darkhalil



Saptarshi Sinha



Omar Emara



Sam Pollard



Kranti Parida



Kaiting Liu



Prajwal Gatti



Siddhant Bansal



Kevin Flanagan



Jacob Chalk



Zhifan Zhu



Rhodri Guerrier



Fahd Abdelazim



Bin Zhu



Davide Moltisanti



Michael Wray



Hazel Doughty



Dima Damen

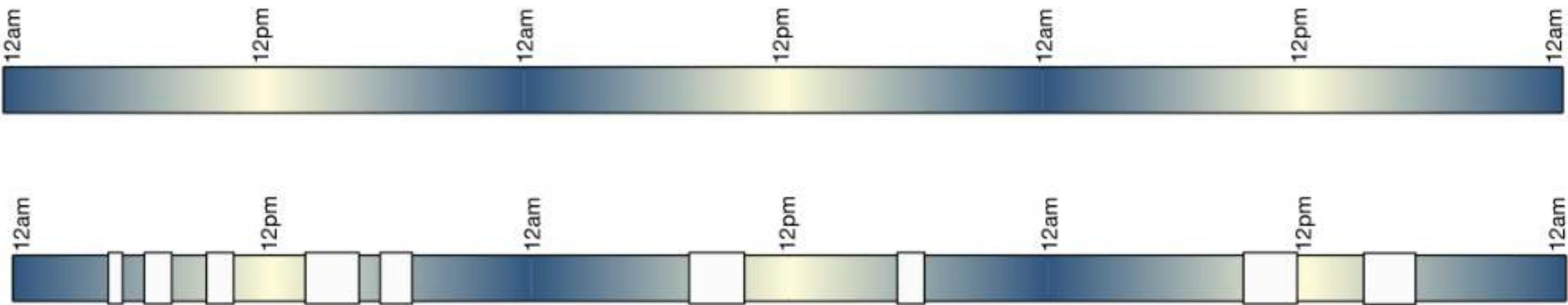


# HD-EPIC





# HD-EPIC





# HD-EPIC



Recorded over 3 days



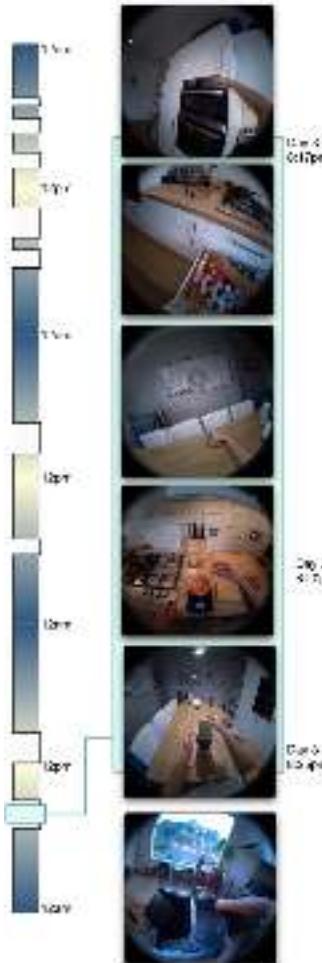


# HD-EPIC





# HD-EPIC





# HD-EPIC



## Recipe: Southwestern Salad

May 3  
1:17pm

1: Preheat the oven to 400F

2: Wash and peel the sweet potatoes and chop into bite-sized pieces. Put the sweet potatoes in a bowl and add the olive oil, cumin, and chili powder. Pour onto tray and roast for 10 mins.

3: Pulse all the dressing ingredients in a food processor until mostly smooth.

**Recipe  
and nutrition**



# HD-EPIC



## Cacio e Pepe (modified)

Ingredients:

200 g penne

400g of pasta of your choice  
(we recommend bucatini)

2 tablespoon of black peppercorn

30 g parmesano

200g of freshly grated pecorino cheese

+25g of slightly salted butter



Steps:

1. Toast the peppercorns until fragrant in a dry frying pan over medium heat, about 2 minutes. Keep them moving to prevent them from burning.

~~Once toasted, roughly crush.~~

→ step 2



2. Cook your choice of pasta in a large pot of generously salted boiling water ~~for around 4-6 minutes~~, or until al dente.

→ step 1



3. While the pasta cooks, add freshly grated cheese and crushed black peppercorns to a large serving bowl. Gradually add a cup of the boiling cooking water constantly mixing to obtain a silky, smooth sauce that's able to completely coat the pasta.

→ step 3





# HD-EPIK



- The **prep** of a corresponding **step** is defined as all essential actions the participant takes to get ready to execute a given step.
- For example, the **step** ‘chop tomato’:
  - **Prep:** retrieve tomato from storage, wash tomato, retrieve a knife and chopping board.
- the **step** ‘add chopped onions and stir’:
  - **Prep:** retrieve onion from storage, retrieve a knife and chopping board, **and chop the onions.**



# HD-EPIC



- Prep



- Step



Cook the pasta in a pan of boiling salted water according to the packet instructions.

Slice the bacon and place in a non-stick frying pan on a medium heat with half a tablespoon of olive oil and ...

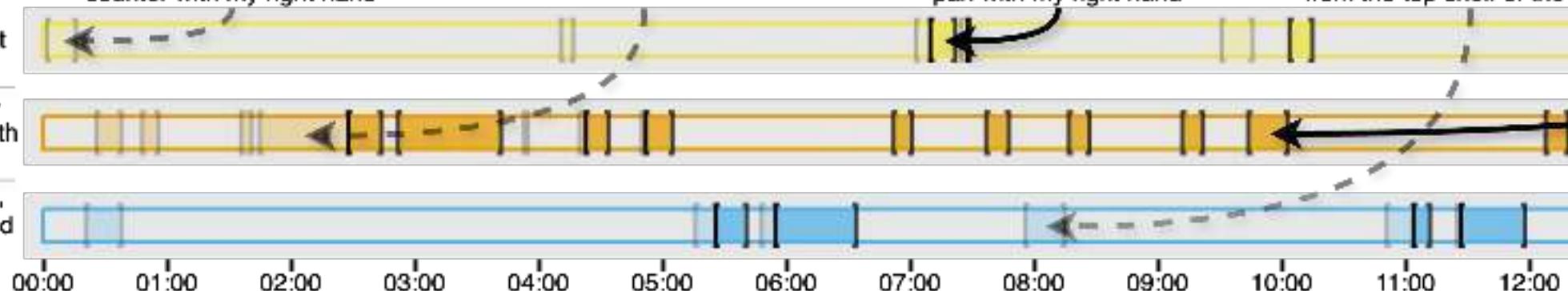
Meanwhile, beat the eggs in a bowl, then finely grate in the Parmesan and mix well.

pick up kettle from its base on the counter with my right hand

pick up packet of bacon

pour water from kettle into the pan with my right hand

pick up block of cheese that from the top shelf of the

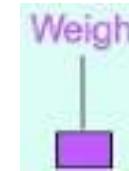




# HD-EPIC



```
"P01_R03_I01": {  
    "name": "penne pasta",  
    "amount": 125,  
    "amount_unit": "g",  
    "calories": 445,  
    "fat": 1.9,  
    "carbs": 90,  
    "protein": 15,
```





# HD-EPIC

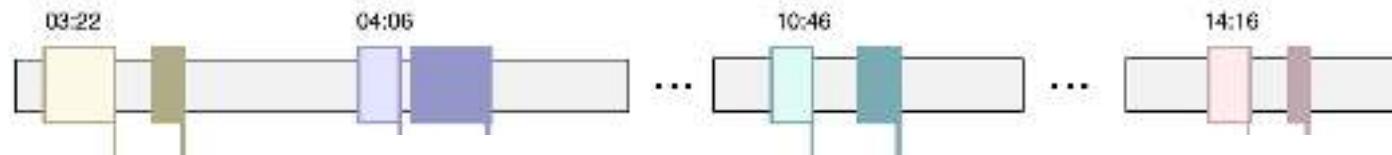


```
"P01_R03_I01": {  
    "name": "penne pasta",  
    "amount": 125,  
    "amount_unit": "g",  
    "calories": 445,  
    "fat": 1.9,  
    "carbs": 90,  
    "protein": 15,
```





# HD-EPIC



Weigh



Ingredient  
nutrition



Sugar  
Qty: 200g  
Cal: 760  
Fat: 0g  
Carbs: 196g  
Protein: 0g

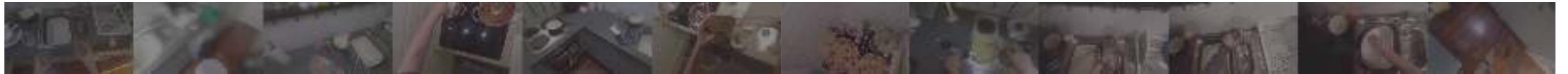


# HD-EPIC

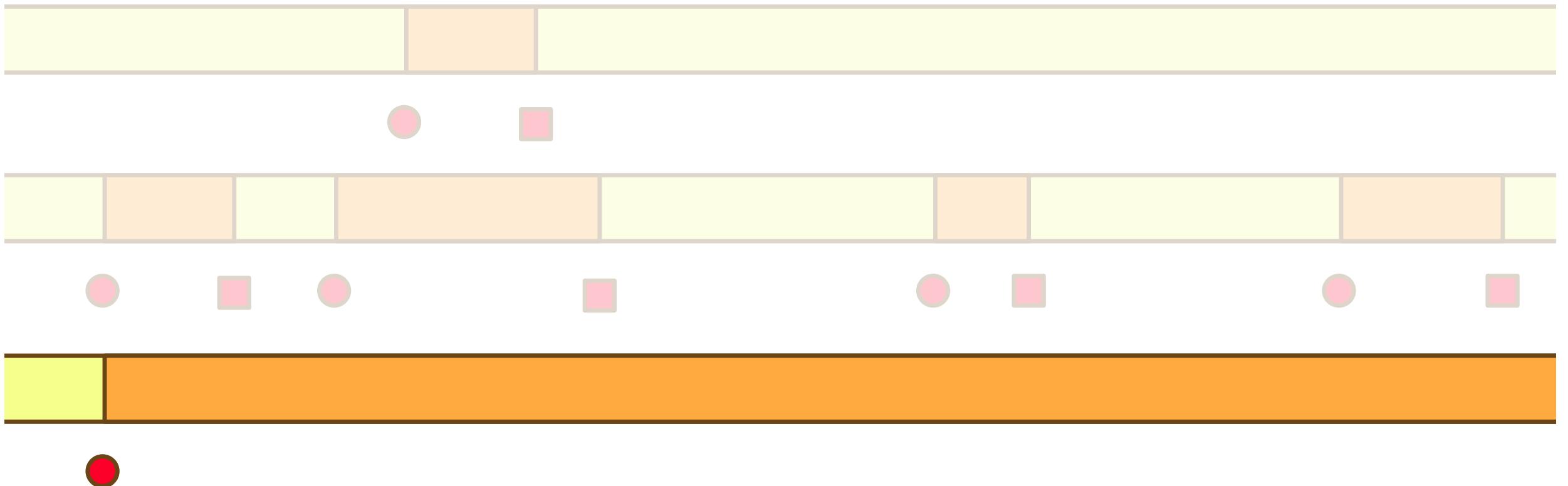


- 69 recipes. Avg: 6.6 steps, 8.1 ingredients, 4.8 hours, 2.1 videos.
- Our longest recipe took 2 days and 6 hours to complete.

HD-EPIC Videos	Recipe Videos
Multiple recipes / video, Multiple videos / recipe	One video == one recipe
All the action including preparatory and cleaning/clearing	Only the steps (prep edited out)
Faster (no explanations)	Slower (with descriptions)
Ingredients weighed on camera	Ingredients pre-weighed



# Egocentric Video Understanding





# Eventually...



- No current model has the context required for this ...
- Impossible to store and process this influx of data ...

But....

- Immense potential ...

# **Learning from Streaming Video with Orthogonal Gradients**

Tengda Han<sup>◦</sup>, Dilara Gokay<sup>◦</sup>, Joseph Heyward<sup>◦</sup>, Chuhan Zhang<sup>◦</sup>  
Daniel Zoran<sup>◦</sup>, Viorica Pătrăucean<sup>◦</sup>, João Carreira<sup>◦</sup>, Dima Damen<sup>◦†</sup>, Andrew Zisserman<sup>◦‡</sup>  
<sup>◦</sup>Google DeepMind, <sup>†</sup>University of Bristol, <sup>‡</sup>University of Oxford



Han et al (2025). Learning from Streaming Video with Orthogonal Gradients. IEEE/CVF Computer Vision and Pattern Recognition (CVPR)

Dima Damen  
ICVSS2025

# Learning from Streaming Videos with Orthogonal Gradients

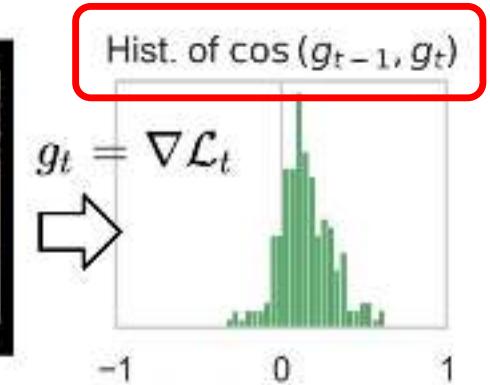


*standard IID samples*



# Learning from Streaming Videos with Orthogonal Gradients

⤒  
shuffled  
loading



gradients are almost **not correlated** over training steps

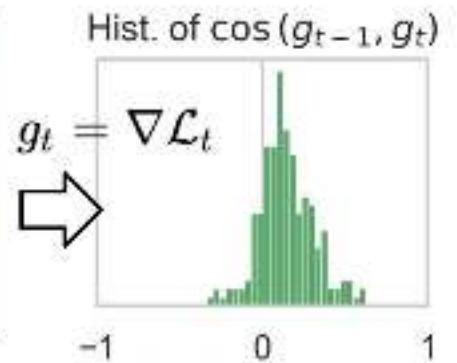


Han et al (2025). Learning from Streaming Video with Orthogonal Gradients. IEEE/CVF Computer Vision and Pattern Recognition (CVPR)

Dima Damen  
ICVSS2025

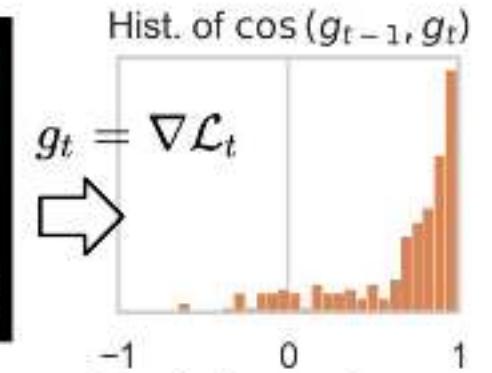
# Learning from Streaming Videos with Orthogonal Gradients

⤚  
shuffled  
loading



gradients are almost **not correlated** over training steps

⌚  
sequential  
loading



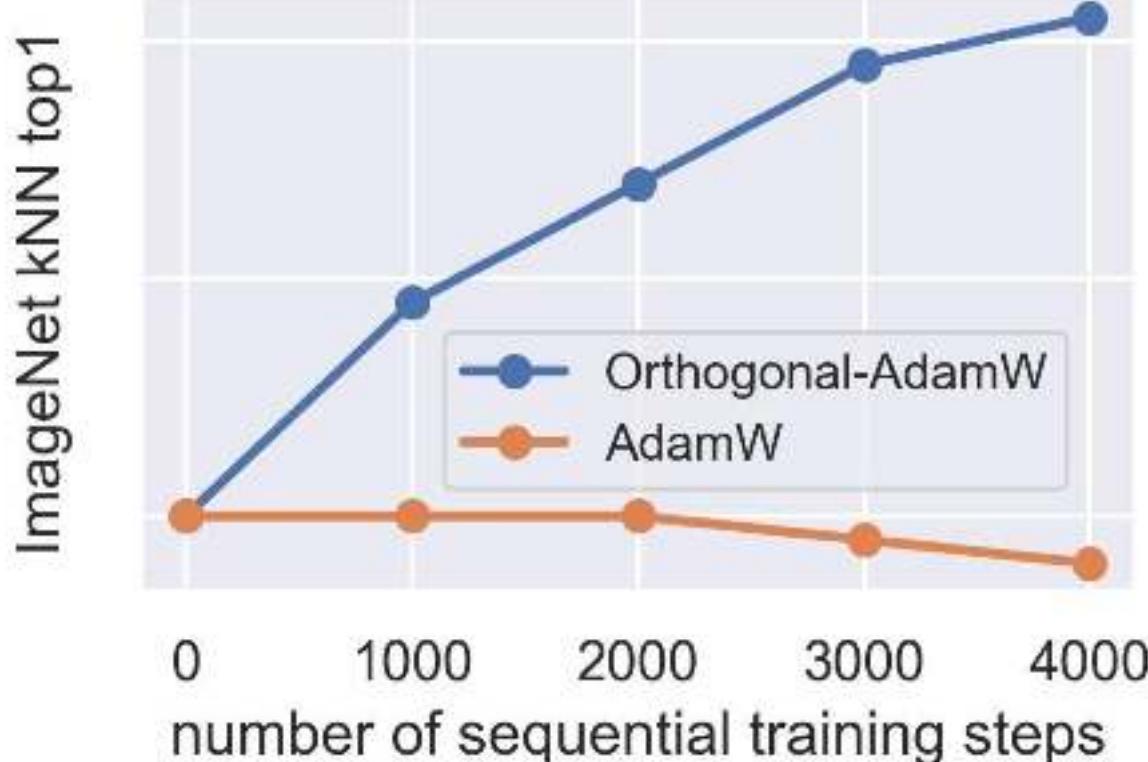
gradients are **highly correlated** over training steps



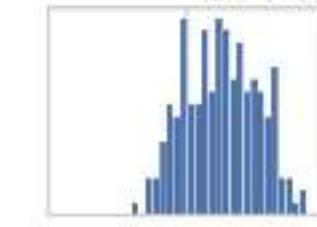
Han et al (2025). Learning from Streaming Video with Orthogonal Gradients. IEEE/CVF Computer Vision and Pattern Recognition (CVPR)

Dima Damen  
ICVSS2025

# Learning from Streaming Videos with Orthogonal Gradients

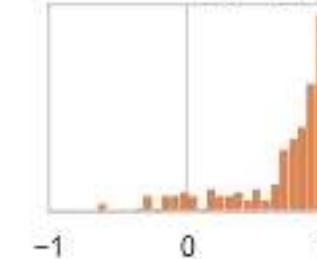


Hist. of  $\cos(g_{t-1}, g_t)$



Orthogonal Optimizer

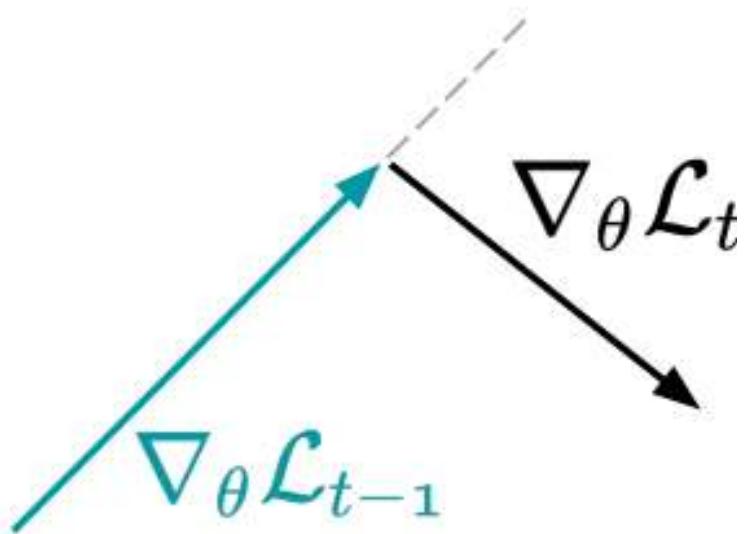
Hist. of  $\cos(g_{t-1}, g_t)$



Han et al (2025). Learning from Streaming Video with Orthogonal Gradients. IEEE/CVF Computer Vision and Pattern Recognition (CVPR)

Dima Damen  
ICVSS2025

# Learning from Streaming Videos with Orthogonal Gradients



(a)



Han et al (2025). Learning from Streaming Video with Orthogonal Gradients. IEEE/CVF Computer Vision and Pattern Recognition (CVPR)

Dima Damen  
ICVSS2025

# Learning from Streaming Videos with Orthogonal Gradients

**Algorithm 2**

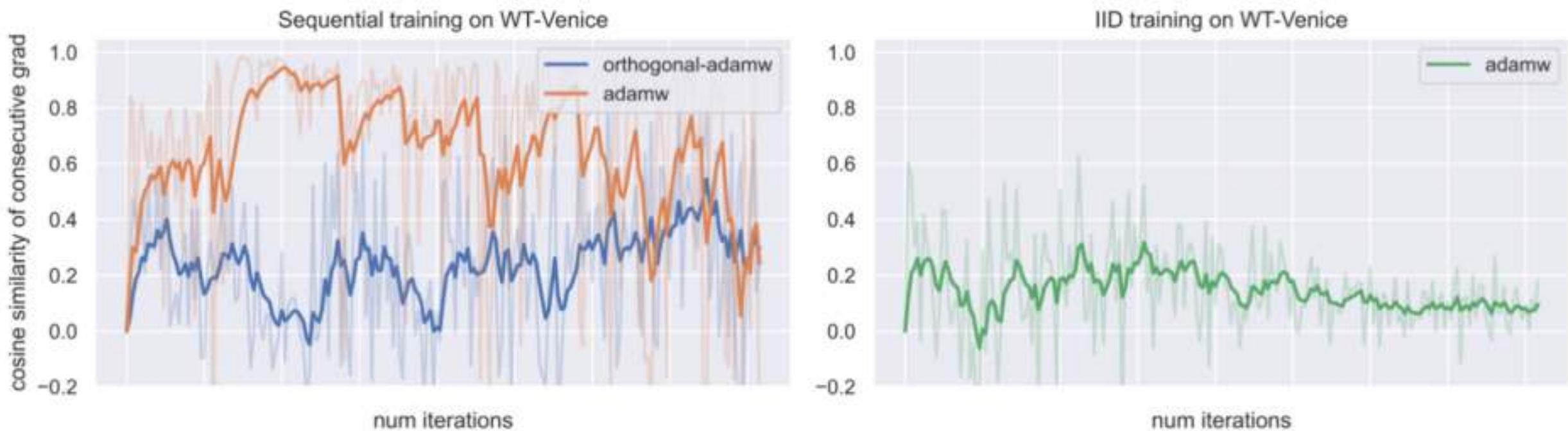
AdamW

**Require:** Learning rate  $\eta > 0$ , weight decay coefficient  $\lambda > 0$ , decay rates  $\beta_1, \beta_2 \in [0, 1]$ , small constant  $\epsilon > 0$ , initial parameters  $\theta_0$ , number of iterations  $T$

- 1: Initialize first moment vector  $m_0 = 0$ , and second moment vector  $v_0 = 0$
- 2: **for**  $t = 1$  to  $T$  **do**
- 3:     Sample a mini-batch of data  $\mathcal{B}_t$  from the training set
- 4:     Compute the gradient:  $g_t = \nabla_{\theta} \mathcal{L}(\theta_{t-1}; \mathcal{B}_t)$
- 5:     Update biased first moment estimate:  $m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$
- 6:     Update biased second moment estimate:  $v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$
- 7:     Compute bias-corrected first moment:  $\hat{m}_t = \frac{m_t}{1 - \beta_1^t}$
- 8:     Compute bias-corrected second moment:  $\hat{v}_t = \frac{v_t}{1 - \beta_2^t}$
- 9:     Apply weight decay:  $\theta_{t-1} = \theta_{t-1} - \eta \lambda \theta_{t-1}$
- 10:     Update parameters:  $\theta_t = \theta_{t-1} - \eta \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon}$
- 11: **end for**



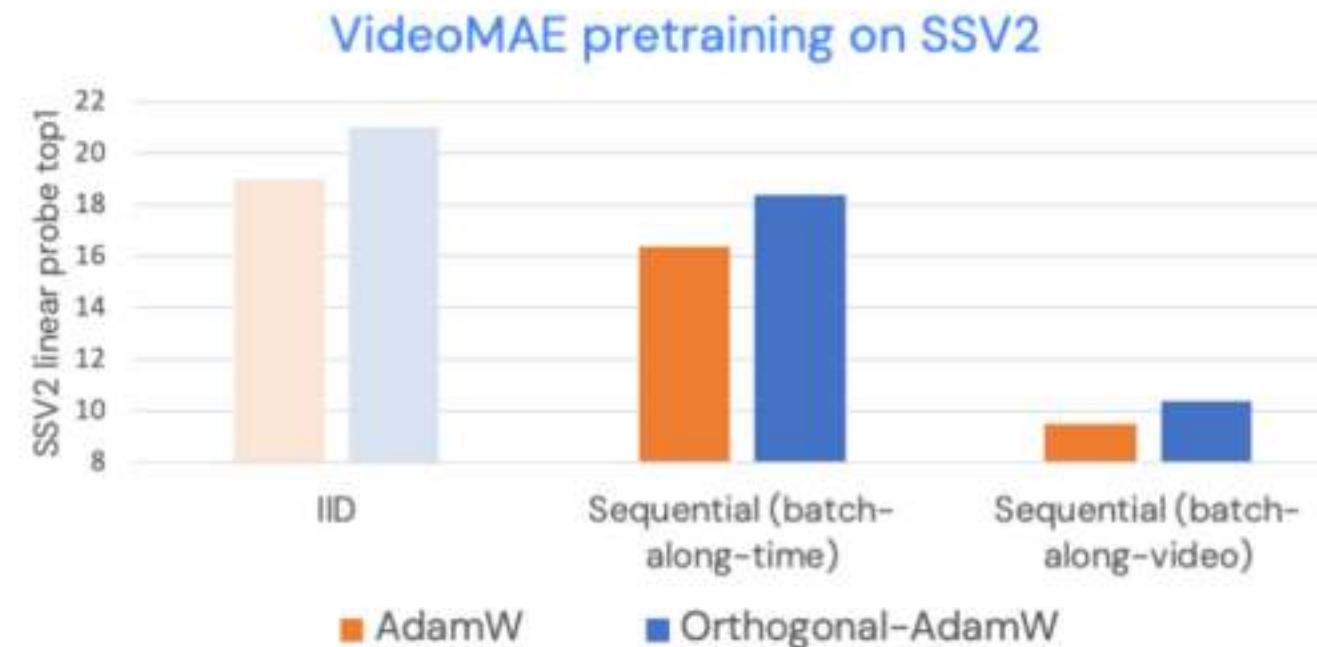
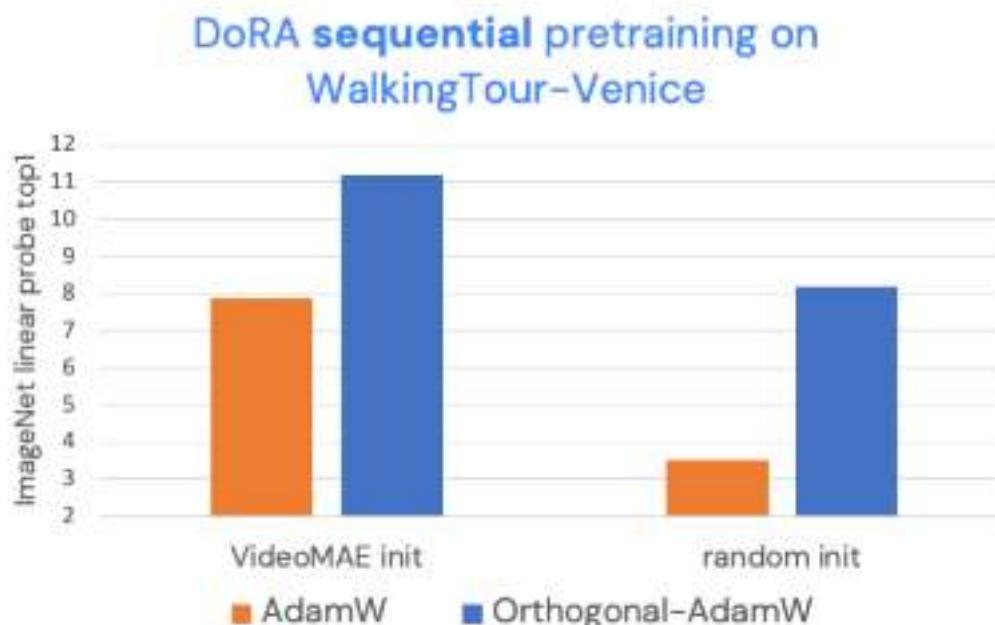
# Learning from Streaming Videos with Orthogonal Gradients



Han et al (2025). Learning from Streaming Video with Orthogonal Gradients. IEEE/CVF Computer Vision and Pattern Recognition (CVPR)

Dima Damen  
ICVSS2025

# Learning from Streaming Videos with Orthogonal Gradients



Han et al (2025). Learning from Streaming Video with Orthogonal Gradients. IEEE/CVF Computer Vision and Pattern Recognition (CVPR)

Dima Damen  
ICVSS2025

# In today's tutorial



Motivation and Datasets in  
Egocentric Video Understanding



Video Understanding  
Out of the Frame



Video Understanding:  
Data and Tasks



Teaser: The Wizard of Oz  
at the Sphere



Videos are Multimodal



Outlook into the Future of  
Egocentric Vision



Connected Videos of One's Life



Conclusion



# First-person Hyperlapse Videos

Johannes Kopf   Michel F. Cohen   Richard Szeliski  
Microsoft Research

[research.microsoft.com/hyperlapse](http://research.microsoft.com/hyperlapse)

SIGGRAPH 2014



Rendered from **single** input frame



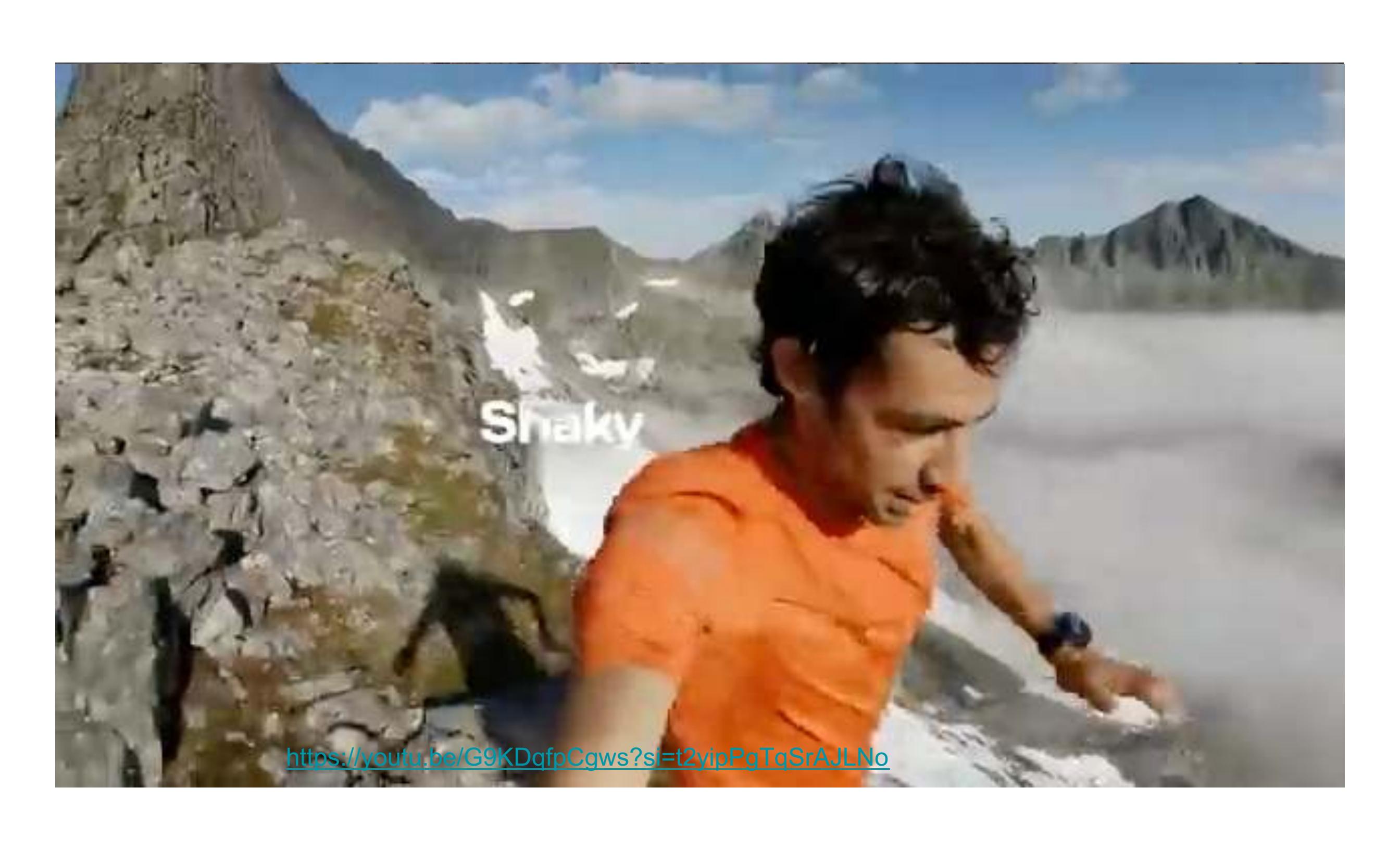
## 4 years later (2018)

The screenshot shows the GoPro website with a dark header bar containing links for Home, About, Support, and My Account. The main navigation menu includes Camera, Apps, Accessories, Lifestyle Gear, GoPro Subscription, and Shop by Activity. Below the menu, there's a search bar and a link to 'GoPro Academy'. The main content area features a large banner with the headline 'SHAKY VIDEO IS DEAD: HERO7 BLACK IS HERE' and a video player showing a person riding a roller coaster. Below the video are two smaller thumbnail images: 'HERO7 16 NEW THINGS' and 'HyperSmooth'.

Today, GoPro announced its new product lineup including the \$399 flagship, HERO7 Black, which sets a new bar for video stabilization with its standout feature, HyperSmooth.

HyperSmooth is the best in-camera video stabilization ever featured in a camera. It makes it easy to capture professional-looking, gimbal-like stabilized video without the expense or hassle of a motorized gimbal. And HyperSmooth works underwater and in high-shock and wind situations where gimbals fail. HERO7 Black with HyperSmooth video stabilization – you've got to see it to believe it.



A man with dark hair and a beard, wearing an orange t-shirt, is riding a motorcycle on a rocky, dirt road. He is looking back over his shoulder. The background features a range of mountains with some snow and a clear blue sky.

Shaky

<https://youtu.be/G9KDqfpCgws?si=t2yipPgTqSrAJLNo>



# Video Understanding Out of the Frame

While neighbouring frames have been used in on-board video stabilisation,  
approaches focused on video understanding within the frame



**EPIC-KITCHENS**

# EPIC Fields

with: V Tschernezki\*, A Darkhalil\*, Z Zhu\*,  
D Fouhey, I Laina, D Larlus, A Vedaldi

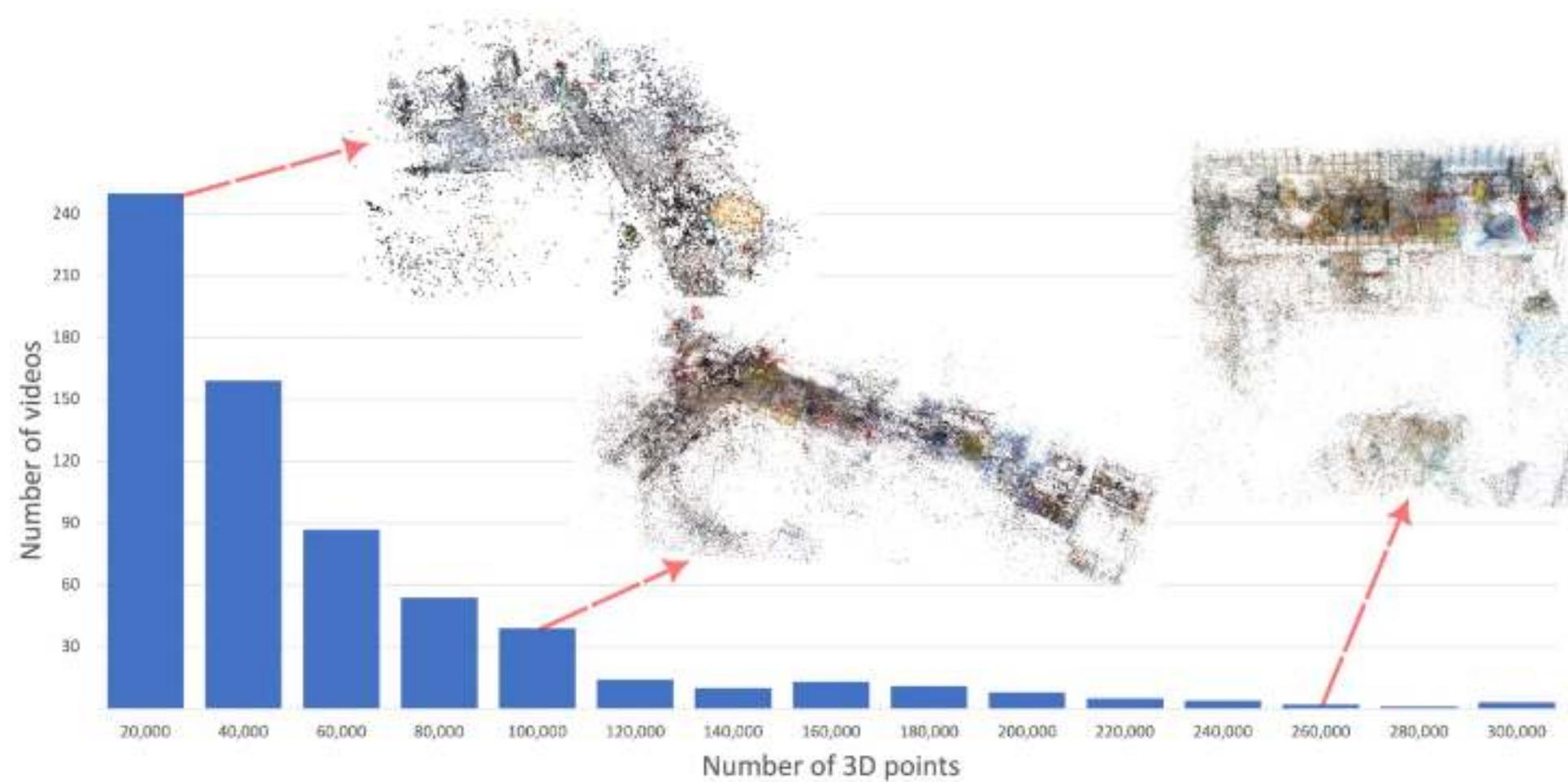


Figure 4: **Number of 3D points histogram.** The majority of our reconstructions generate less than 40,000 points that are enough to represent the kitchen. However, some reconstructions have more than 100,000, we include the point clouds for each points range showing the fine details covered by having more points

# EPIC Fields

with: V Tschernezki\*, A Darkhalil\*, Z Zhu\*,  
D Fouhey, I Laina, D Larlus, A Vedaldi

Table 1: Comparison of datasets commonly used in dynamic new-view synthesis.

Dataset	#Scenes	Seq. Length	Monocular	Semantics
Nerfies [37]	4	8–15 sec	-	-
D-NeRF [41]	8	1–3 sec	-	-
Plenoptic Video [22]	6	10–60 sec	-	-
NVIDIA Dynamic Scene Dataset [65]	12	1–5 sec	4 / 12	-
HyperNeRF [38]	16	8–15 sec	13 / 16	-
iPhone [13]	14	8–15 sec	7 / 14	-
SAFF [25]	8	1–5sec	-	✓
<b>EPIC Fields (ours)</b>	50	6–37 min (Avg 22)	50 / 50	✓



# Video Understanding Out of the Frame

What can we now do with these reconstructions:

- Point Tracking
- Object Tracking
- Gaze Estimation



# EgoPoints: Advancing Point Tracking for Egocentric Videos

Ahmad Darkhalil<sup>1</sup> Rhodri Guerrier<sup>1</sup> Adam W. Harley<sup>2</sup> Dima Damen<sup>1</sup>

<sup>1</sup>University of Bristol      <sup>2</sup>Stanford University

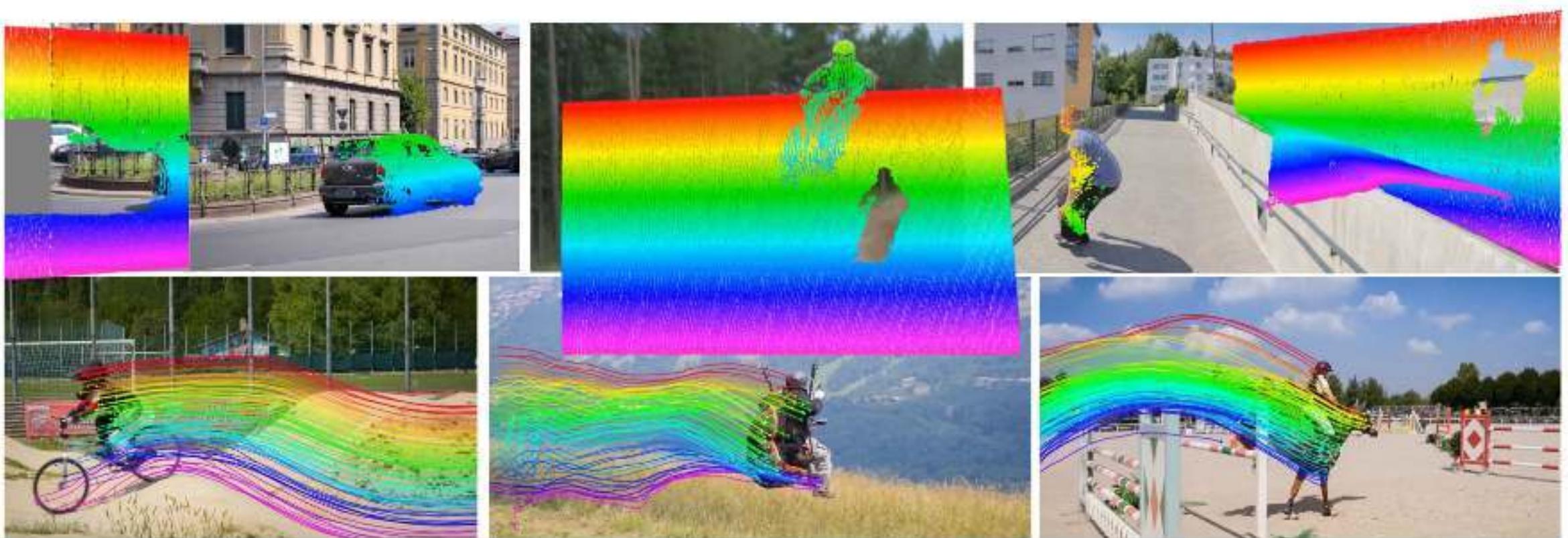


Dima Damen  
ICVSS2025

# What is point tracking?

with: Ahmad Darkhalil  
Rhodri Guerrier  
Adam W Harley

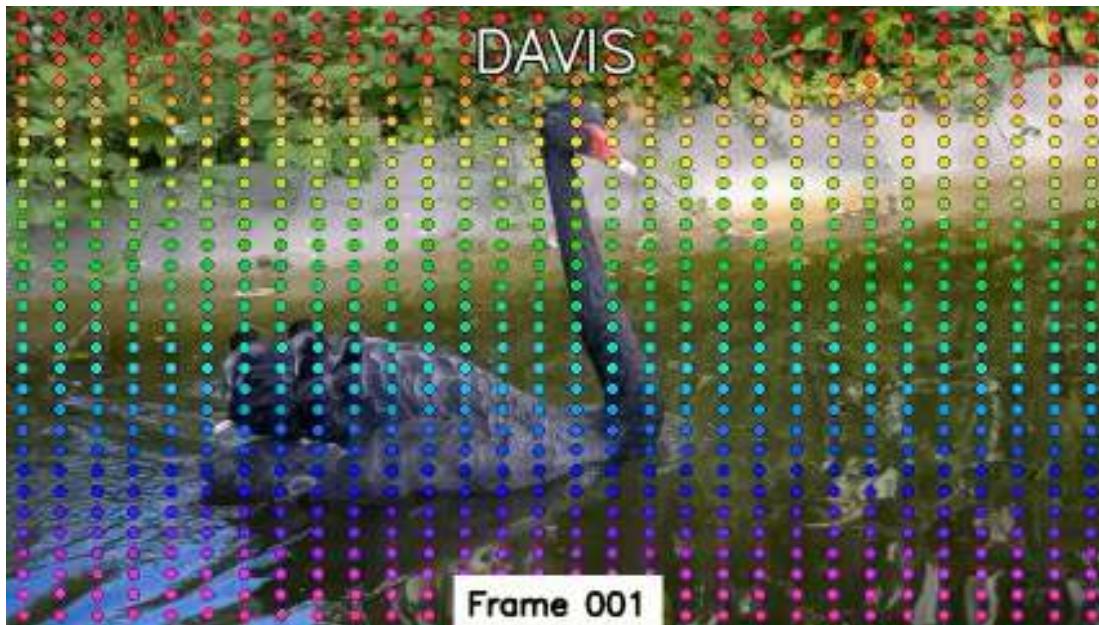
- Given: Query points in one frame
- Track these points throughout the video



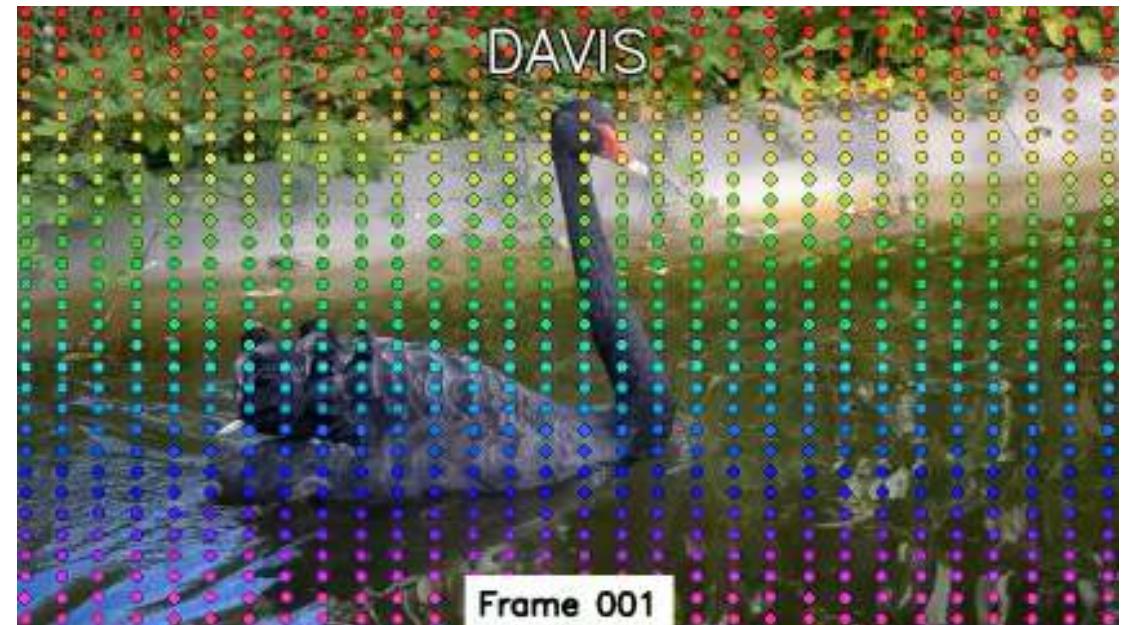
# SOTA on current benchmarks

with: Ahmad Darkhalil  
Rhodri Guerrier  
Adam W Harley

LocoTrack



CoTracker3

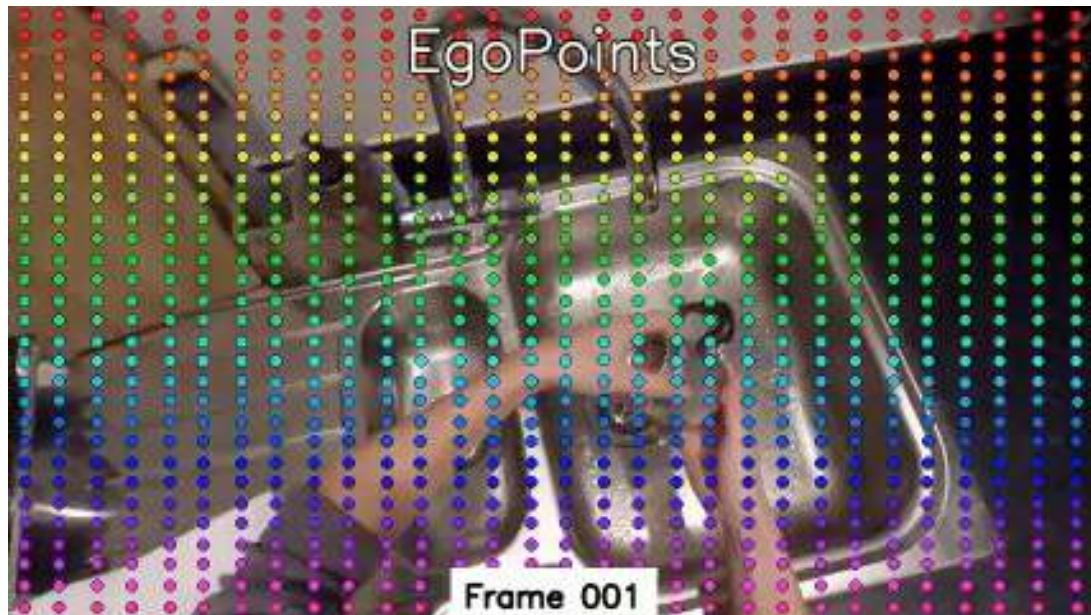


# Current Models Struggle with Egocentric Videos

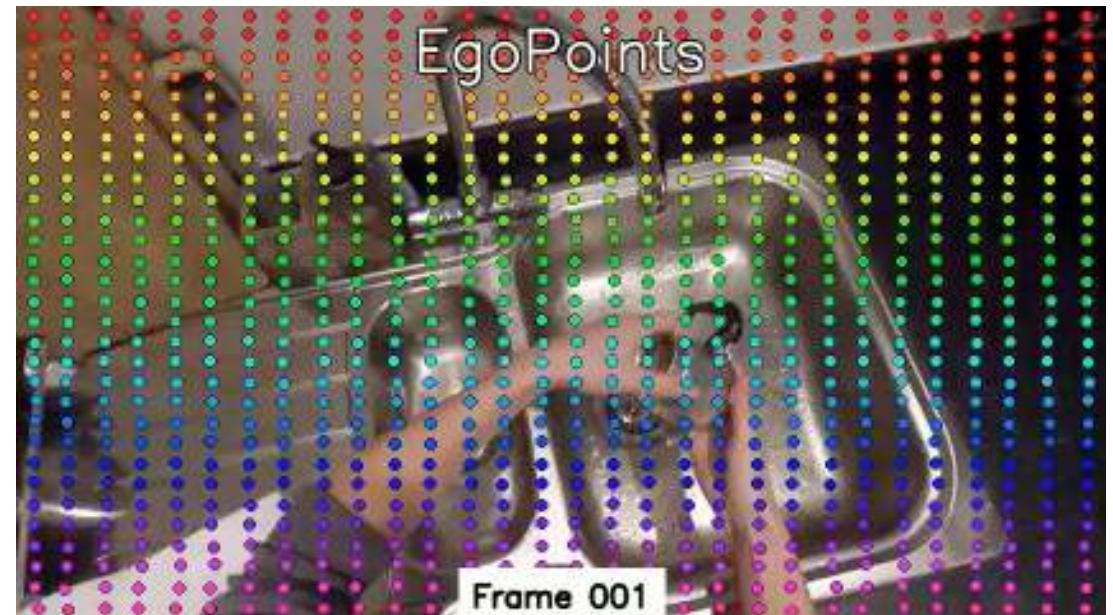
with: Ahmad Darkhalil  
Rhodri Guerrier  
Adam W Harley

- Head motion and motion blur
- Frequent re-identification

LocoTrack



CoTracker3



# Main Contributions

with: Ahmad Darkhalil  
Rhodri Guerrier  
Adam W Harley

- Identify challenges that point trackers face in egocentric videos.
- Propose a new benchmark (EgoPoints) and new metrics to showcase these challenges
- Propose K-EPIC, a pipeline to generate semi-real training data

# EgoPoints Annotation interface

with: Ahmad Darkhalil  
Rhodri Guerrier  
Adam W Harley



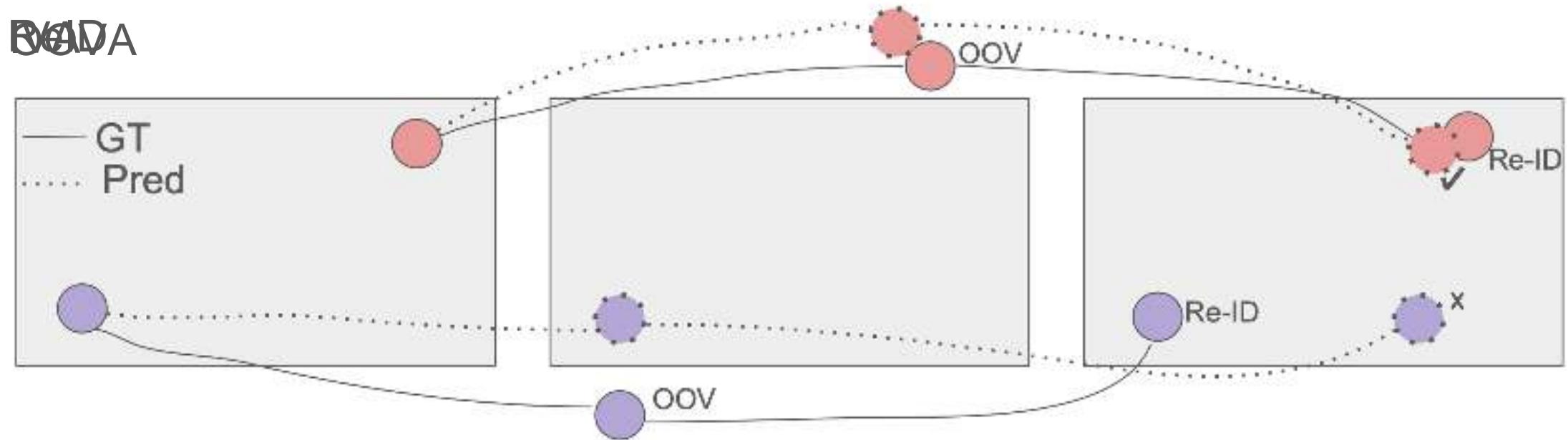
# EgoPoints Benchmark

with: Ahmad Darkhalil  
Rhodri Guerrier  
Adam W Harley



# Proposed Metrics

with: Ahmad Darkhalil  
Rhodri Guerrier  
Adam W Harley



# EgoPoints Benchmark

with: Ahmad Darkhalil  
Rhodri Guerrier  
Adam W Harley

Dataset	Total Tracks	OOV Tracks	ReID Tracks	Avg. Video Length	Avg. Points/Frame
TAP-Vid-DAVIS	650	94	10	66.6	<b>21.7</b>
EgoPoints	<b>4703</b>	<b>875</b>	<b>593</b>	<b>511.0</b>	8.5

Comparisons of our annotated sequences, EgoPoints, and the commonly used TAP-Vid-DAVIS [7] point tracking benchmarks

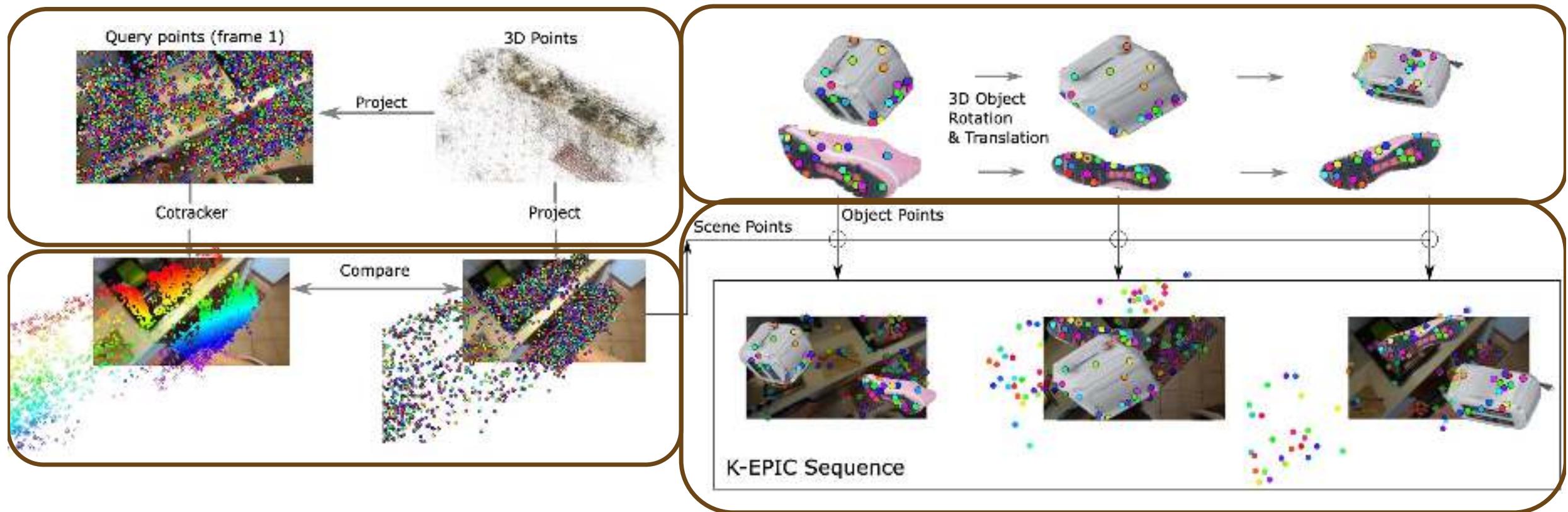
# SOTA Models Struggle on EgoPoints

with: Ahmad Darkhalil  
Rhodri Guerrier  
Adam W Harley

Model	TAP-Vid-DAVIS		EgoPoints			
	$\delta_{\text{avg}} \uparrow$	$\delta_{\text{avg}} \uparrow$	ReID	$\delta_{\text{avg}} \uparrow$	OOVA↑	IVA↑
PIPs++ [42]	64.0	36.9		14.6	50.4	89.2
CoTracker [22]	74.7	38.5		4.8	<b>81.4</b>	73.4
BootsTAPIR Online [8]	65.2	39.6		0.0	0.0	<b>100.0</b>
LocoTrack [4]	75.3	<b>59.4</b>		0.1	0.2	99.9
CoTracker v3 [21]	<b>77.2</b>	50.0		<b>15.0</b>	31.8	99.3

# Pipeline of K-EPIC

with: Ahmad Darkhalil  
Rhodri Guerrier  
Adam W Harley



# Examples from K-EPIC

with: Ahmad Darkhalil  
Rhodri Guerrier  
Adam W Harley



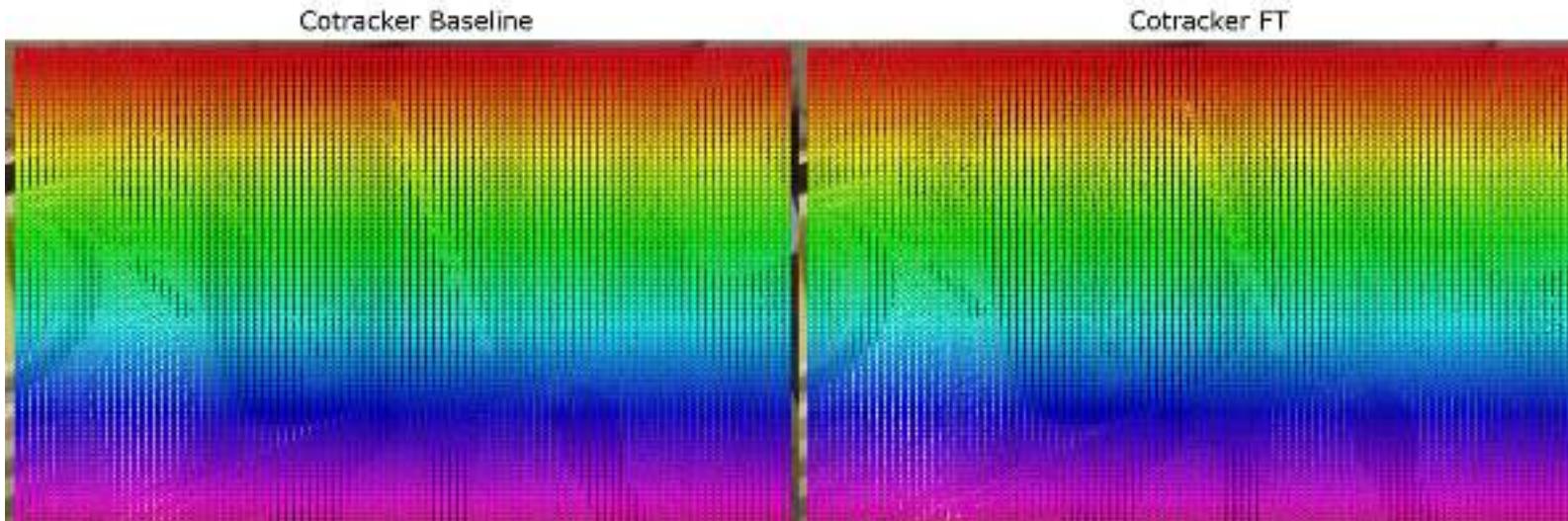
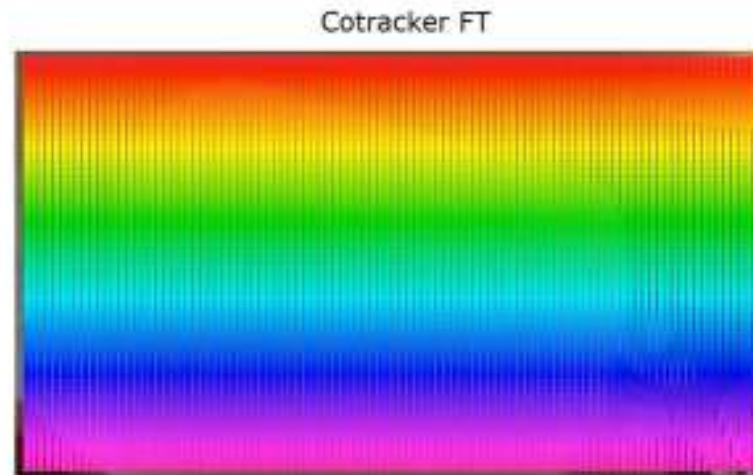
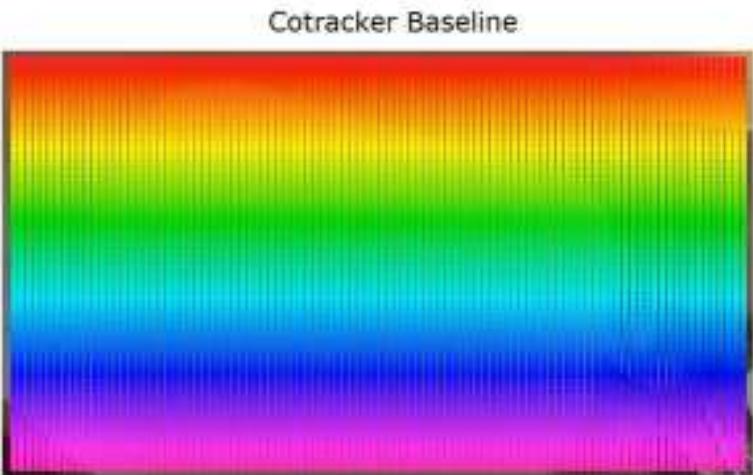
# Improvements After Fine-Tuning on K-EPIC

with: Ahmad Darkhalil  
Rhodri Guerrier  
Adam W Harley

Model	$\delta$ Metrics			Accuracy Metrics			Error
	$\delta_{\text{avg}} \uparrow$	$\delta_{\text{avg}}^* \uparrow$	ReID $\delta_{\text{avg}} \uparrow$	IVA $\uparrow$	OOVA $\uparrow$	OA $\uparrow$	MTE $\downarrow$
PIPs++ [42]	<b>36.9</b>	57.8	14.0	89.2	50.4	–	22.9
PIPs++ w. K-EPIC FT (scene points only)	36.3	57.8	13.0	<b>90.1</b>	<b>53.0</b>	–	22.9
PIPs++ w. K-EPIC FT (scene and object points)	36.6	<b>58.1</b>	<b>16.8</b>	89.9	52.0	–	<b>22.2</b>
CoTracker [22]	38.5	54.8	4.8	73.4	81.4	81.0	52.1
CoTracker w. K-EPIC FT (scene points only)	38.9	56.0	6.3	74.8	<b>85.4</b>	80.7	51.3
CoTracker w. K-EPIC FT (scene and object points)	<b>39.6</b>	<b>57.5</b>	<b>7.2</b>	<b>78.1</b>	82.0	<b>81.8</b>	<b>40.5</b>

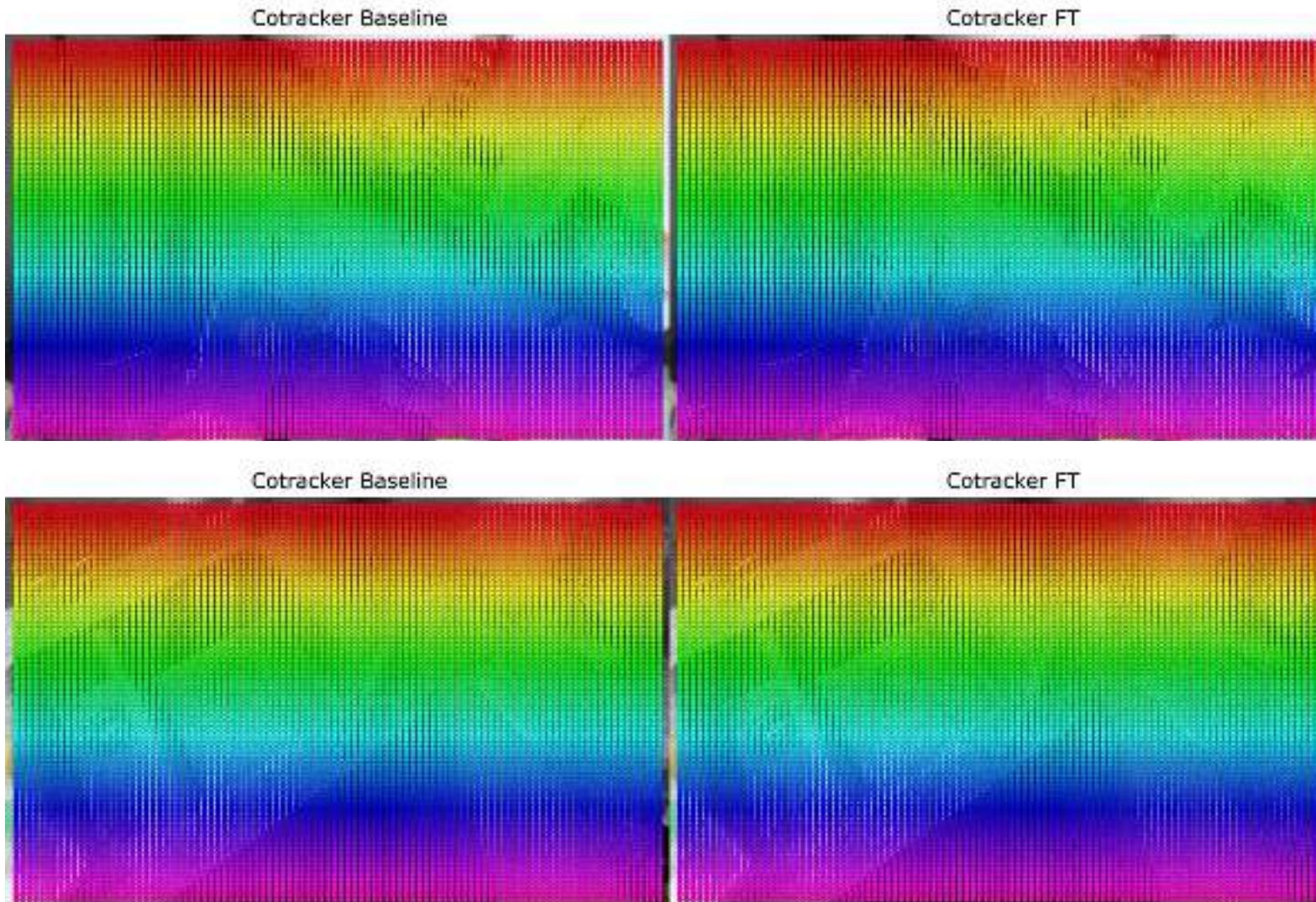
# Qualitative Examples of CoTracker

with: Ahmad Darkhalil  
Rhodri Guerrier  
Adam W Harley



# Qualitative Examples of CoTracker

with: Ahmad Darkhalil  
Rhodri Guerrier  
Adam W Harley



# AllTracker – newest work

Table 1. Comparison against recent point trackers and optical flow models, across nine datasets. We evaluate  $\delta_{\text{avg}}$  (higher is better), using an input resolution of  $384 \times 512$ . The benchmarks are BADJA [3], CroHD [46], TAPVid-DAVIS [11], DriveTrack [1], EgoPoints [10], Horse10 [33], TAPVid-Kinetics [11], RGB-Stacking [28], and RoboTAP [52].

Method	Params.	Training	Bad.	Cro.	Dav.	Dri.	Ego.	Hor.	Kin.	Rgb.	Rob.	Avg.
RAFT [47]	5.26	Flow mix	23.7	29.3	48.5	44.8	41.0	27.8	64.3	82.8	72.2	48.3
SEA-RAFT [54]	19.66	Flow mix	23.9	21.9	48.7	49.4	44.0	33.1	64.3	85.7	67.6	48.7
AccFlow [55]	11.76	Flow mix	10.3	22.2	23.5	26.4	4.0	12.1	38.8	63.2	57.9	28.7
PIPs++ [59]	17.57	PointOdyssey	34.1	27.5	62.5	51.3	38.5	21.4	64.2	70.4	73.4	49.3
LocoTrack [6]	11.52	Kubric	41.4	43.1	68.0	66.5	58.4	48.9	70.0	80.3	76.9	61.5
BootsTAPIR [13]	54.70	Kubric+15M	42.7	34.9	67.9	66.9	56.8	48.8	70.6	81.0	78.2	60.9
DELTA [37]	59.17	Kubric	44.6	42.9	75.3	67.8	40.3	41.8	66.5	83.0	74.8	59.7
CoTracker2 [23]	45.43	Kubric	40.0	31.7	70.9	67.8	43.2	33.9	65.8	73.4	73	55.5
CoTracker3-Kub [25]	25.39	Kubric	47.5	<b>48.9</b>	<b>77.4</b>	<b>69.8</b>	58.0	47.5	70.6	83.4	77.2	64.5
CoTracker3 [25]	25.39	Kubric+15k	48.3	44.5	77.1	<b>69.8</b>	60.4	47.1	71.8	84.2	81.6	65.0
AllTracker-Tiny-Kub	6.29	Kubric	45.4	39.6	73.7	65.1	55.9	45.2	70.6	86.1	79.3	62.3
AllTracker-Tiny	6.29	Kubric+mix	<b>47.5</b>	39.8	74.3	63.9	58.3	45.5	71.5	88.1	80.7	63.3
AllTracker-Kub	16.48	Kubric	46.4	42.3	75.2	66.1	60.3	<b>49.0</b>	71.3	<b>90.1</b>	82.2	64.8
AllTracker	16.48	Kubric+mix	<b>51.5</b>	44.0	76.3	65.8	<b>62.5</b>	<b>49.0</b>	<b>72.3</b>	90.0	<b>83.4</b>	<b>66.1</b>

# Out of Sight, not Out of Mind

with: Chiara Plizzari  
Toby Perrett

Shubham Goel  
Angjoo Kanazawa

## Spatial Cognition from Egocentric Video: Out of Sight, Not Out of Mind

Chiara Plizzari

Shubham Goel

Toby Perrett

Jacob Chalk

Angjoo Kanazawa

Dima Damen

<http://dimadamen.github.io/OSNOM>



Politecnico  
di Torino

Berkeley  
UNIVERSITY OF CALIFORNIA



University of  
BRISTOL



Plizzari et al (2025). Spatial Cognition from Egocentric Video: Out of Sight, Not Out of Mind. 3DV

...na Damen  
ICVSS2025

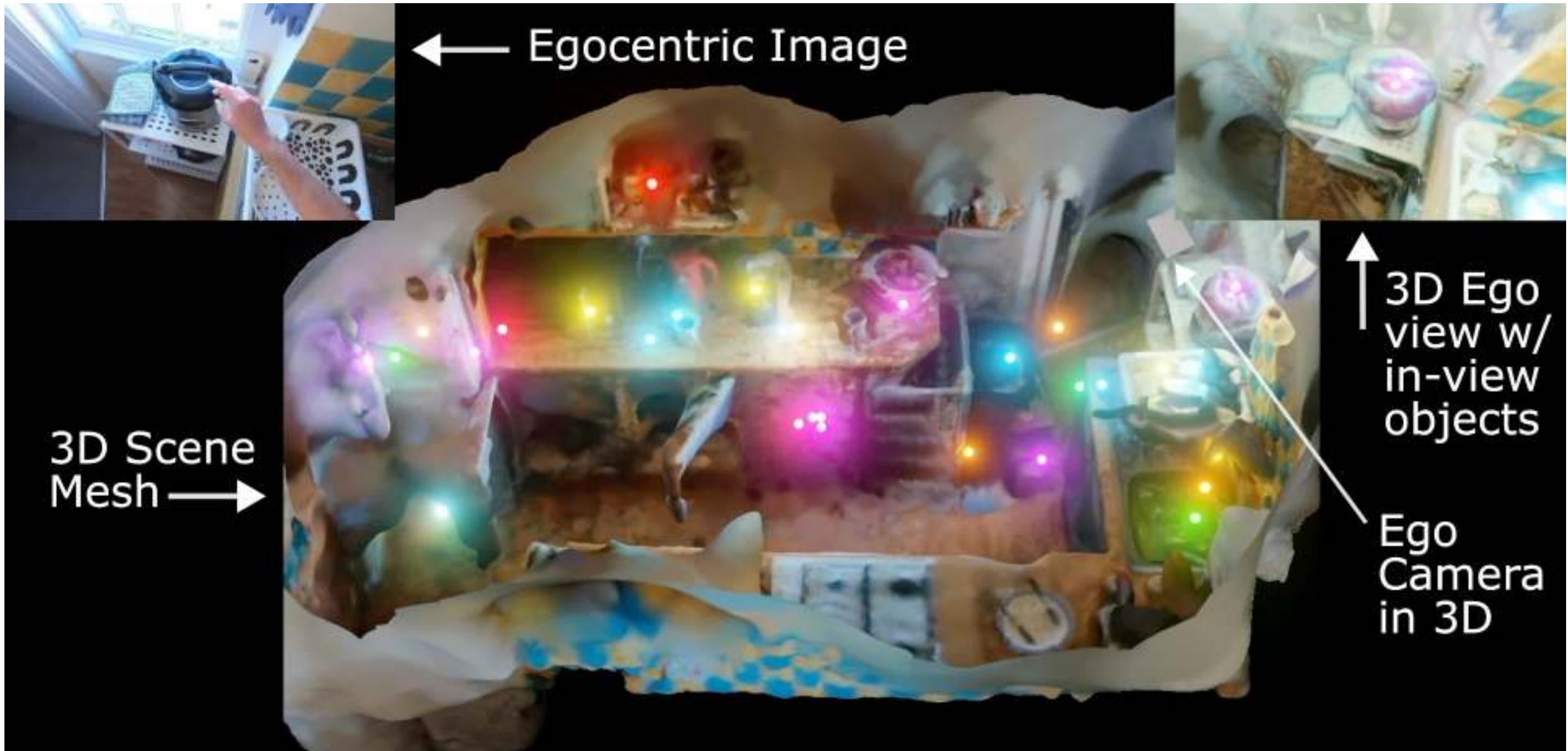
# Out of Sight, not Out of Mind

with: Chiara Plizzari  
Toby Perrett

Shubham Goel  
Angjoo Kanazawa



All active/moved objects in this video are represented by neon balls.  
Their initial positions are shown at the start of the video



All active/moved objects in this video are represented by neon balls.  
Their initial positions are shown at the start of the video

# Out of Sight, not Out of Mind

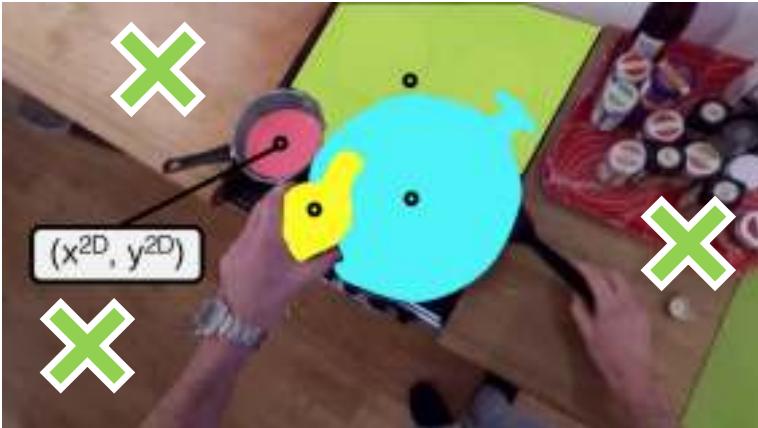
with: Chiara Plizzari  
Toby Perrett

Shubham Goel  
Angjoo Kanazawa

Lift

Match

Keep



0.0 ... 1.0

0.3m ... 1.8m



# Out of Sight, not Out of Mind

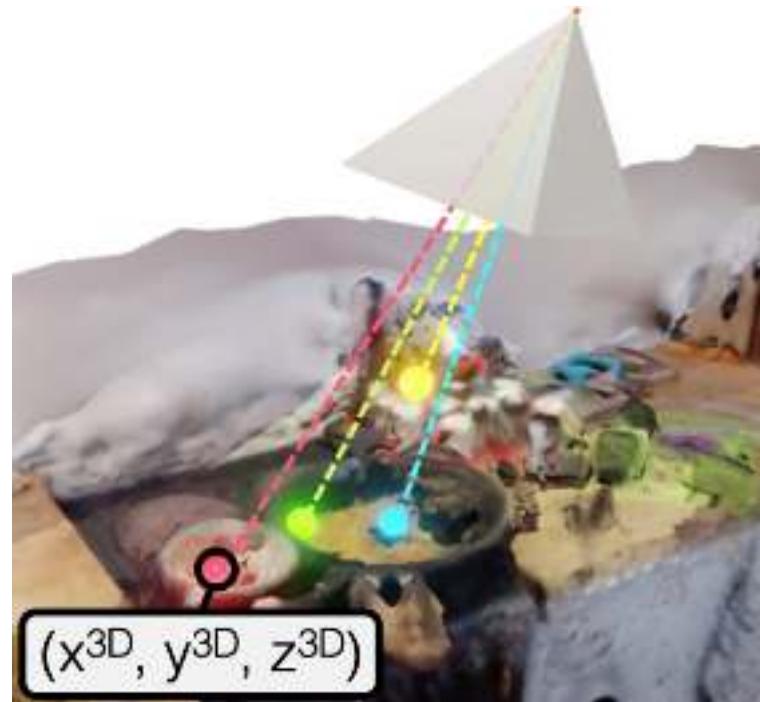
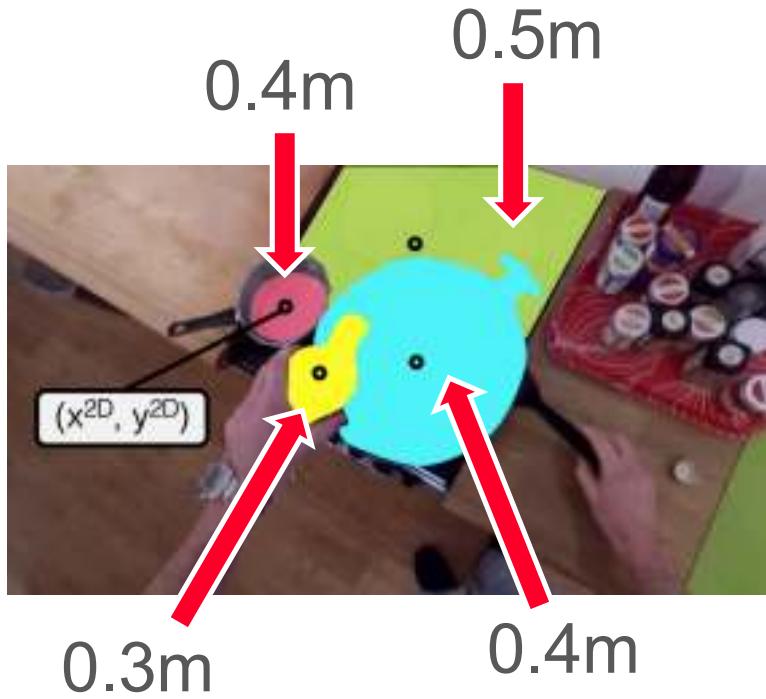
with: Chiara Plizzari  
Toby Perrett

Shubham Goel  
Angjoo Kanazawa

Lift

Match

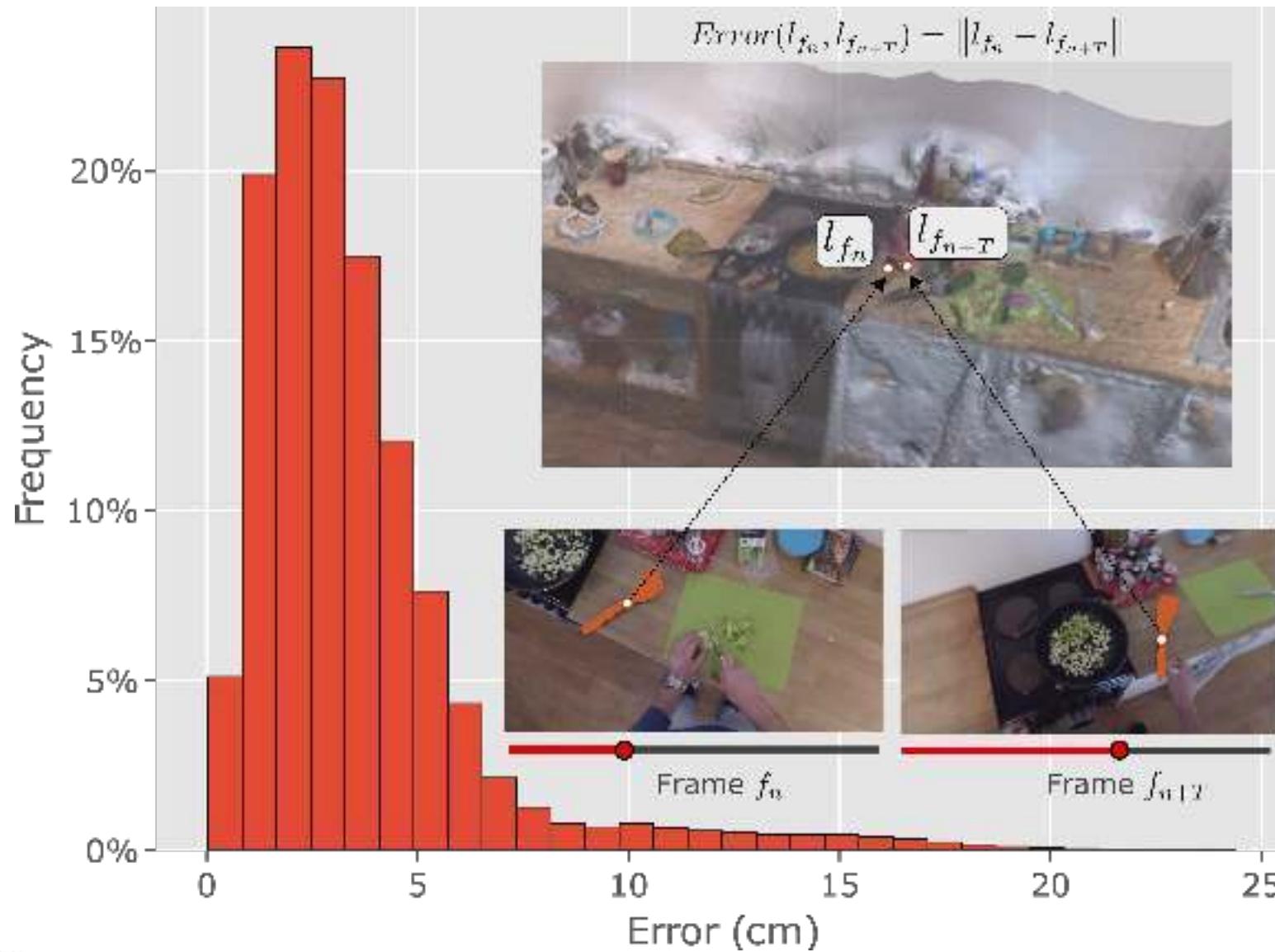
Keep



# Out of Sight, not Out of Mind

with: Chiara Plizzari  
Toby Perrett

Shubham Goel  
Angjoo Kanazawa



# Out of Sight, not Out of Mind

with: Chiara Plizzari  
Toby Perrett

Shubham Goel  
Angjoo Kanazawa

Lift

Match

Keep

Instead of tracking in 2D, we track in 3D, using combination of appearance and location distances

# Out of Sight, not Out of Mind

with: Chiara Plizzari  
Toby Perrett

Shubham Goel  
Angjoo Kanazawa

After we Lift, Match and Keep (LMK), we can reason about an object's visibility and position

- In-View vs Out-of-View
- In-Sight vs Out-of-Sight (Occluded)
- Within-Reach vs Out-of-Reach (defining the camera wearer's near space)



# Out of Sight, not Out of Mind

with: Chiara Plizzari  
Toby Perrett

Shubham Goel  
Angjoo Kanazawa

After we Lift, Match and Keep (LMK), we can reason about an object's visibility and position

- In-View vs Out-of-View
- In-Sight vs Out-of-Sight (Occluded)
- Within-Reach vs Out-of-Reach (defining the camera wearer's near space)



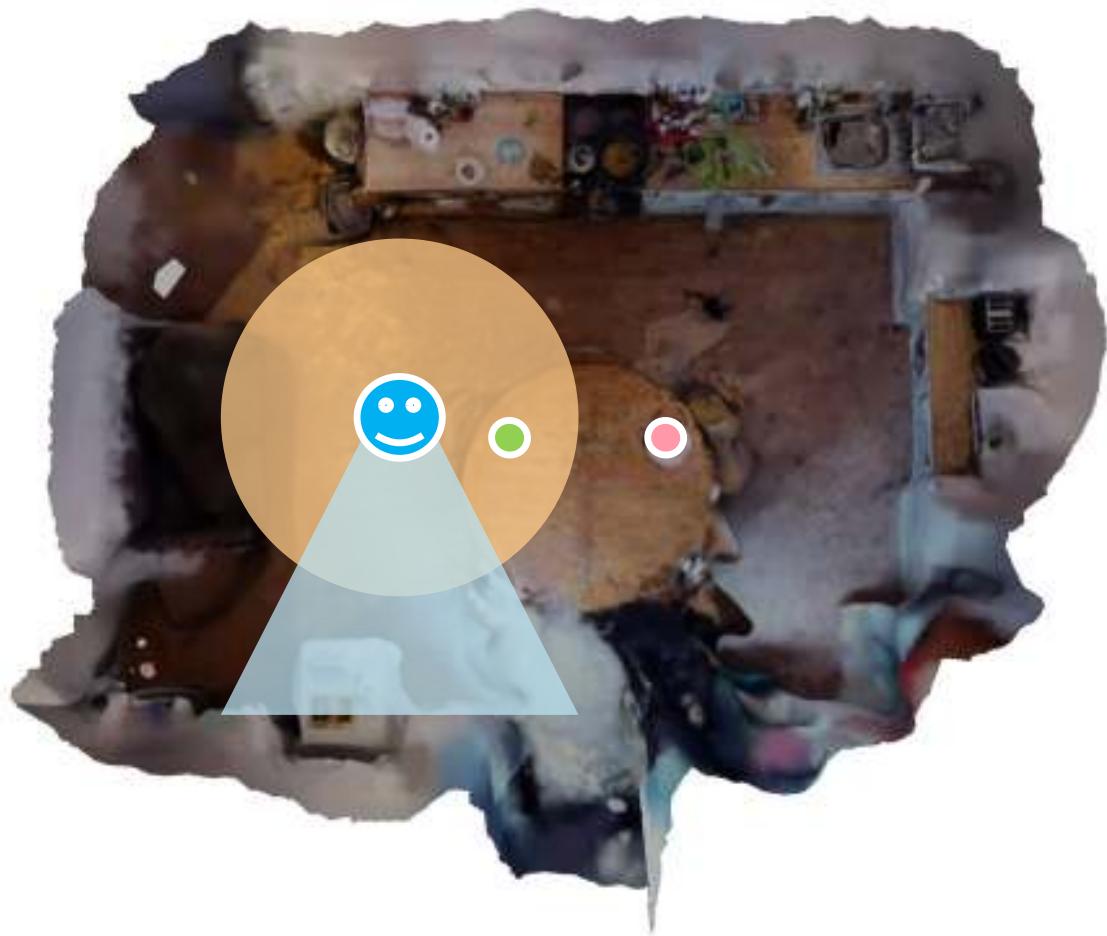
# Out of Sight, not Out of Mind

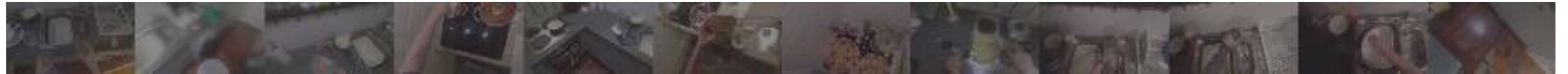
with: Chiara Plizzari  
Toby Perrett

Shubham Goel  
Angjoo Kanazawa

After we Lift, Match and Keep (LMK), we can reason about an object's visibility and position

- In-View vs Out-of-View
- In-Sight vs Out-of-Sight (Occluded)
- Within-Reach vs Out-of-Reach (defining the camera wearer's near space)





# Spatial Cognition from Egocentric Video: Out of Sight, Not Out of Mind

Chiara Plizzari

Shubham

Toby Perrett

Jacob Chalk

Angela

Dima Damen

Ground-Truth??

<http://dimadamen.github.io/OSNOM>



Politecnico  
di Torino

Berkeley  
UNIVERSITY OF CALIFORNIA



University of  
BRISTOL

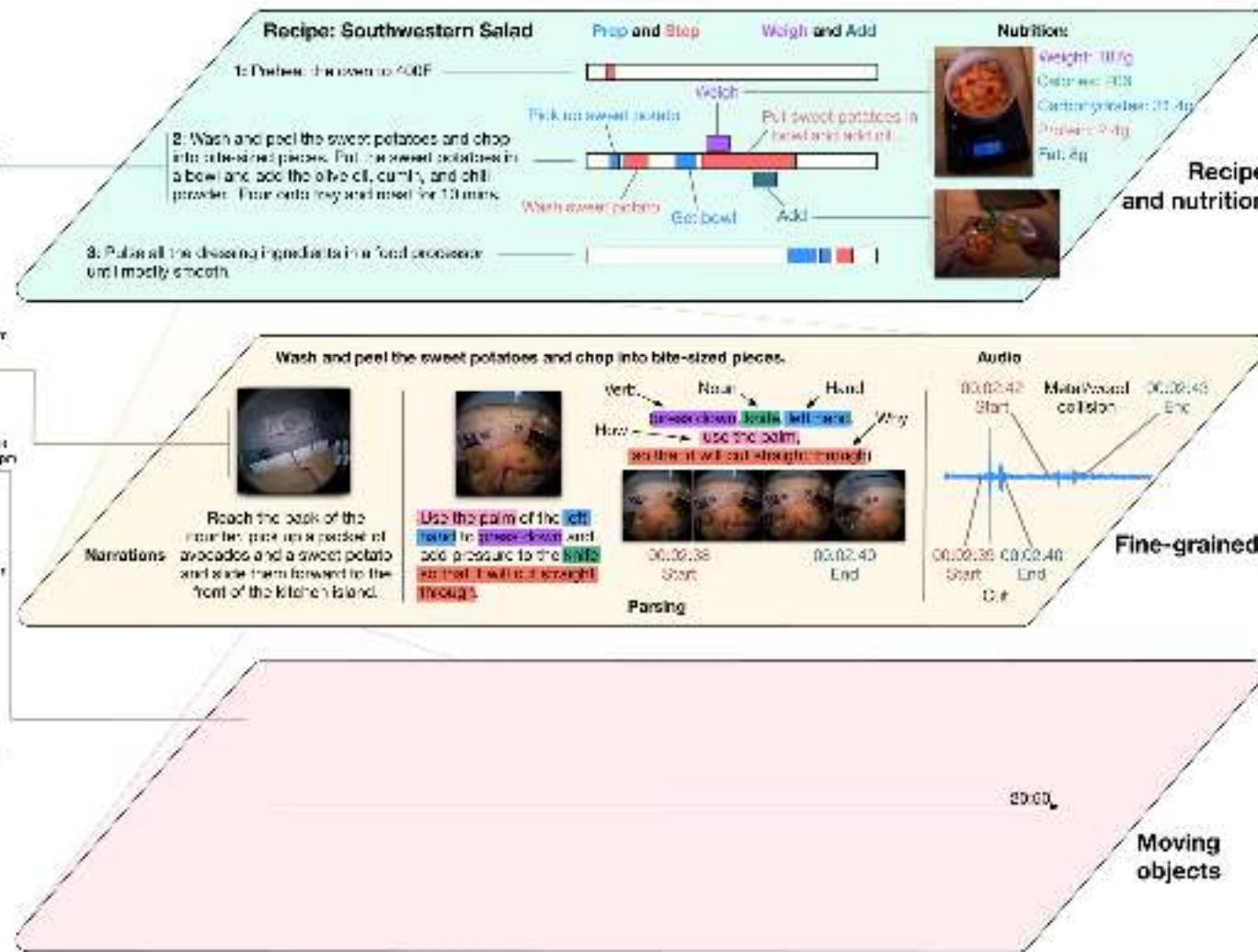


Plizzari et al (2025). Spatial Cognition from Egocentric Video: Out of Sight, Not Out of Mind. 3DV

...na Damen  
ICVSS2025

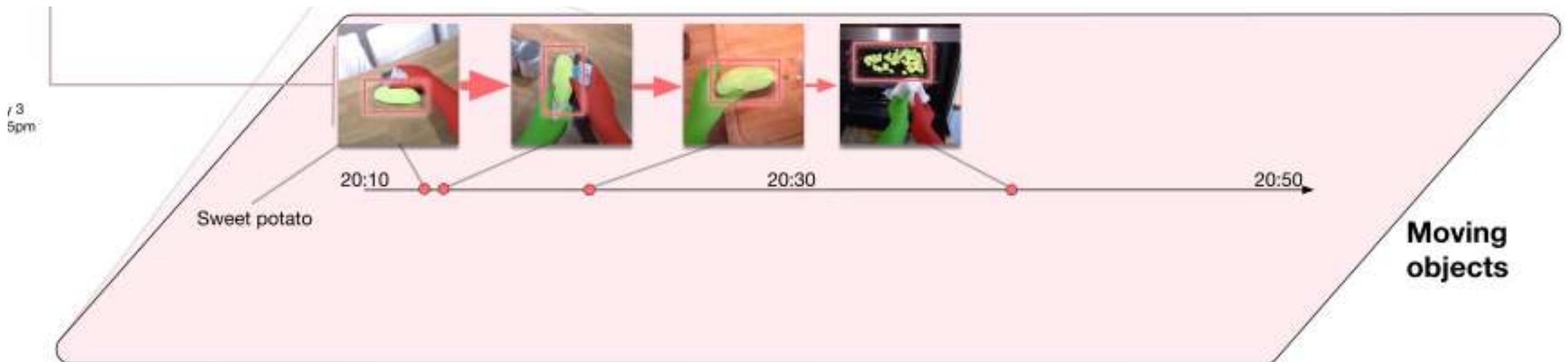


# HD-EPIC





# HD-EPIC

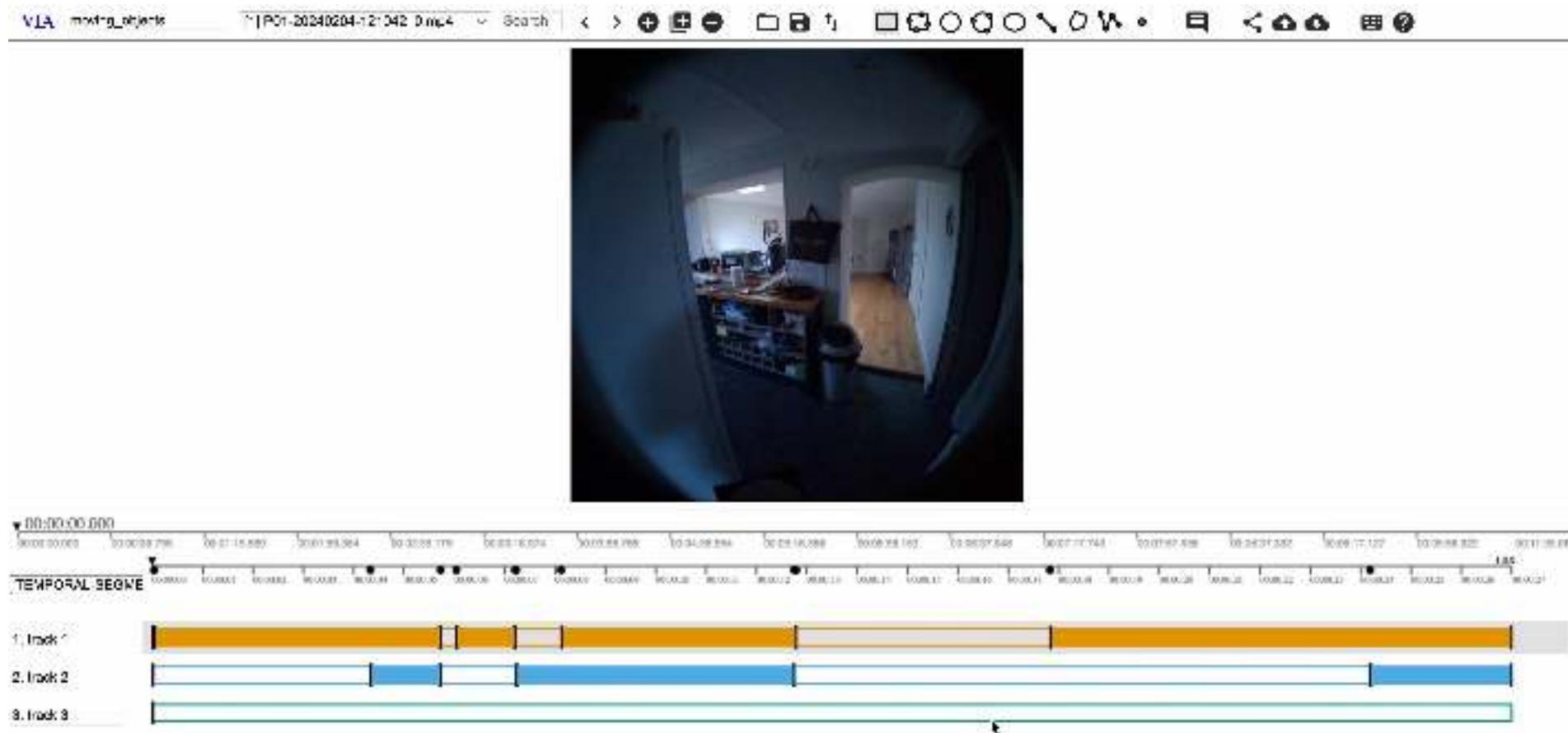




# HD-EPIC



- How to minimize the annotations for tracking objects...

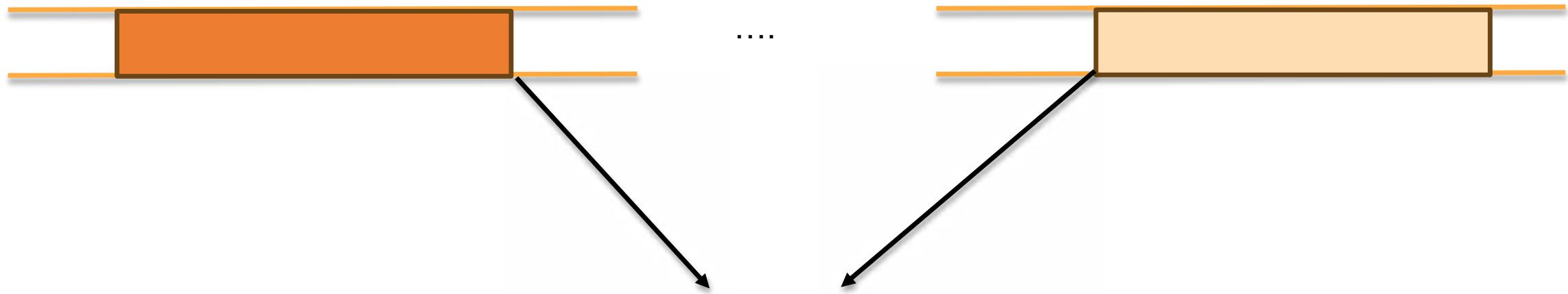




# HD-EPIC



- How to minimize the annotations for tracking objects...





# HD-EPIC



## Current Track

Choose Files - 201 files

04:51

12 / 199

← Previous      Next →      Undo

rubbish bin   box of chicken   wooden chopping board

Enter Track Name (optional)

Create New Track

Inconsistent Query

## Previous Tracks

Sort by Distance

Save Tracks

box of chicken (0.0m)

Add



plastic chopping board (0.3m)

Add



metal cooling rack (0.6m)

Add



plastic measuring cup (1.0m)

Add



hand washing liquid (1.3m)

Add



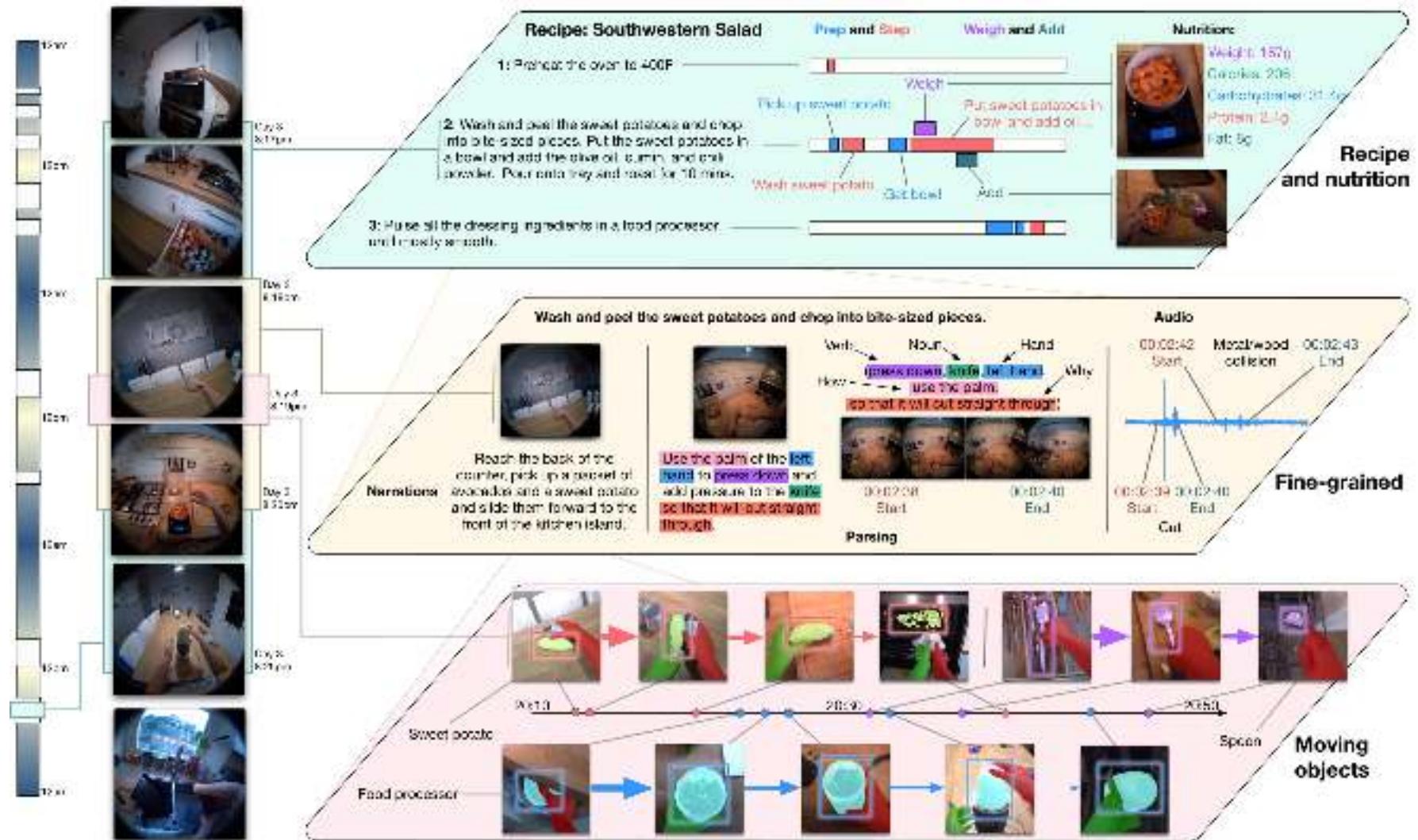
kitchen towel (1.5m)

Add





# HD-EPIC

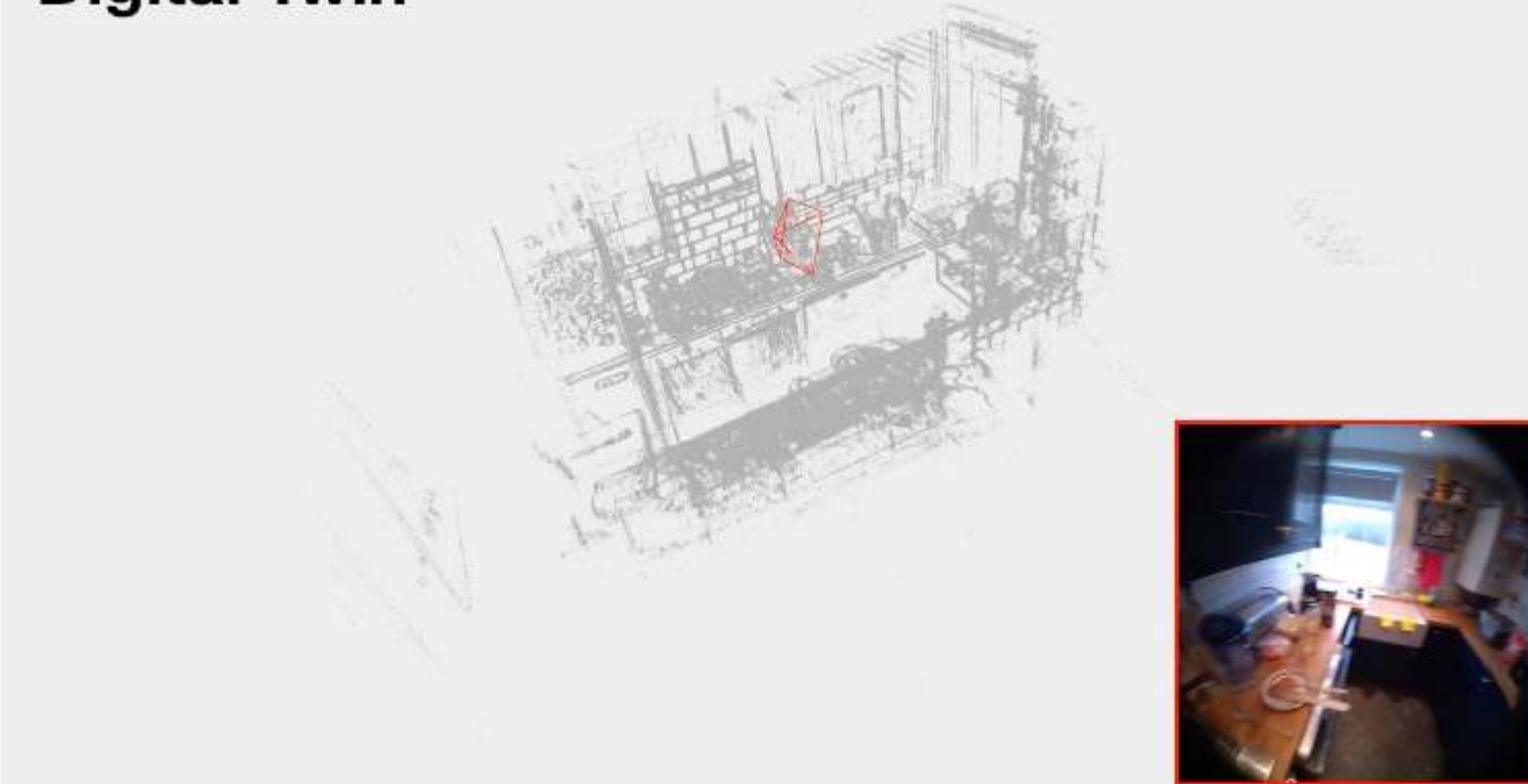




# HD-EPIC

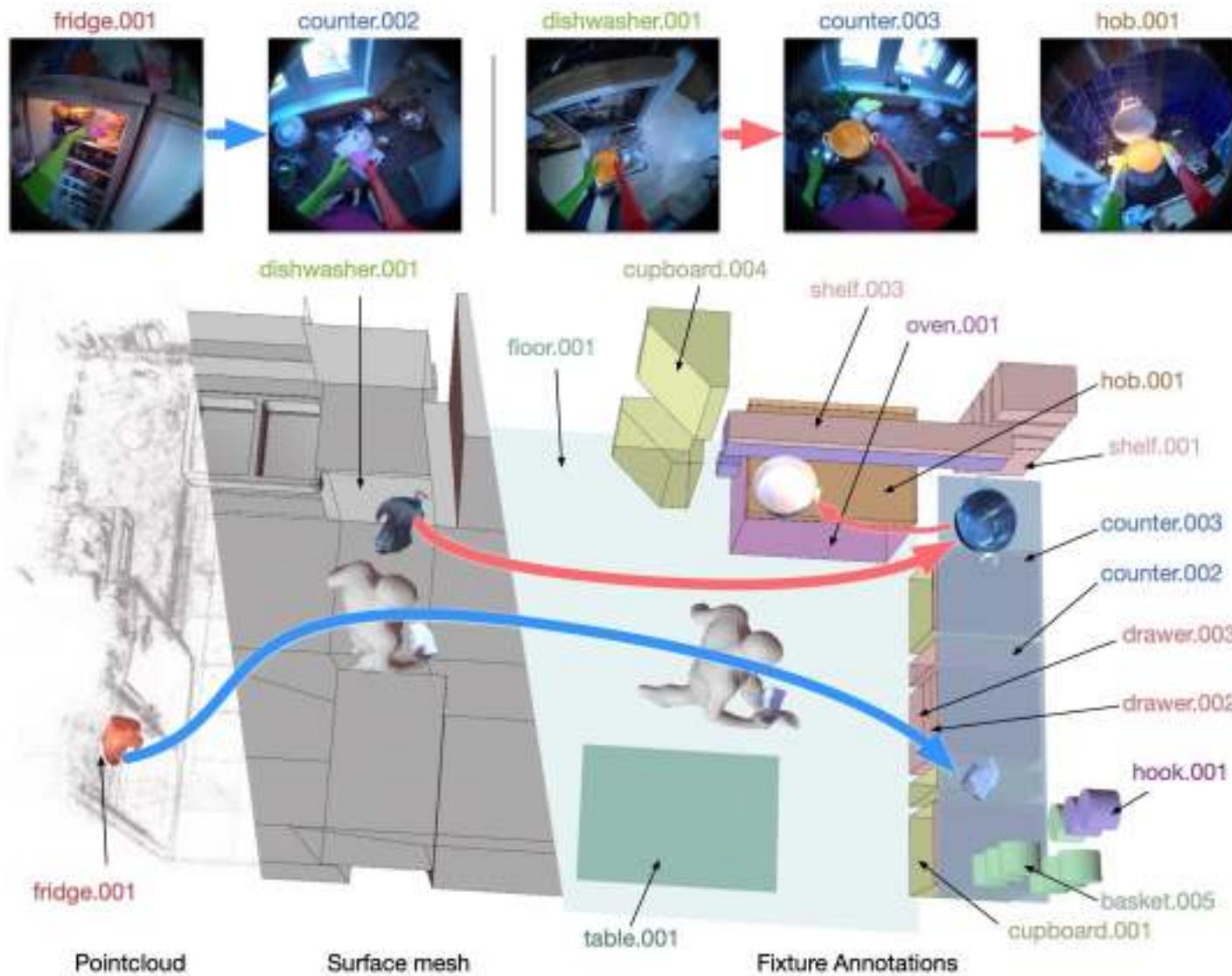


## Digital Twin



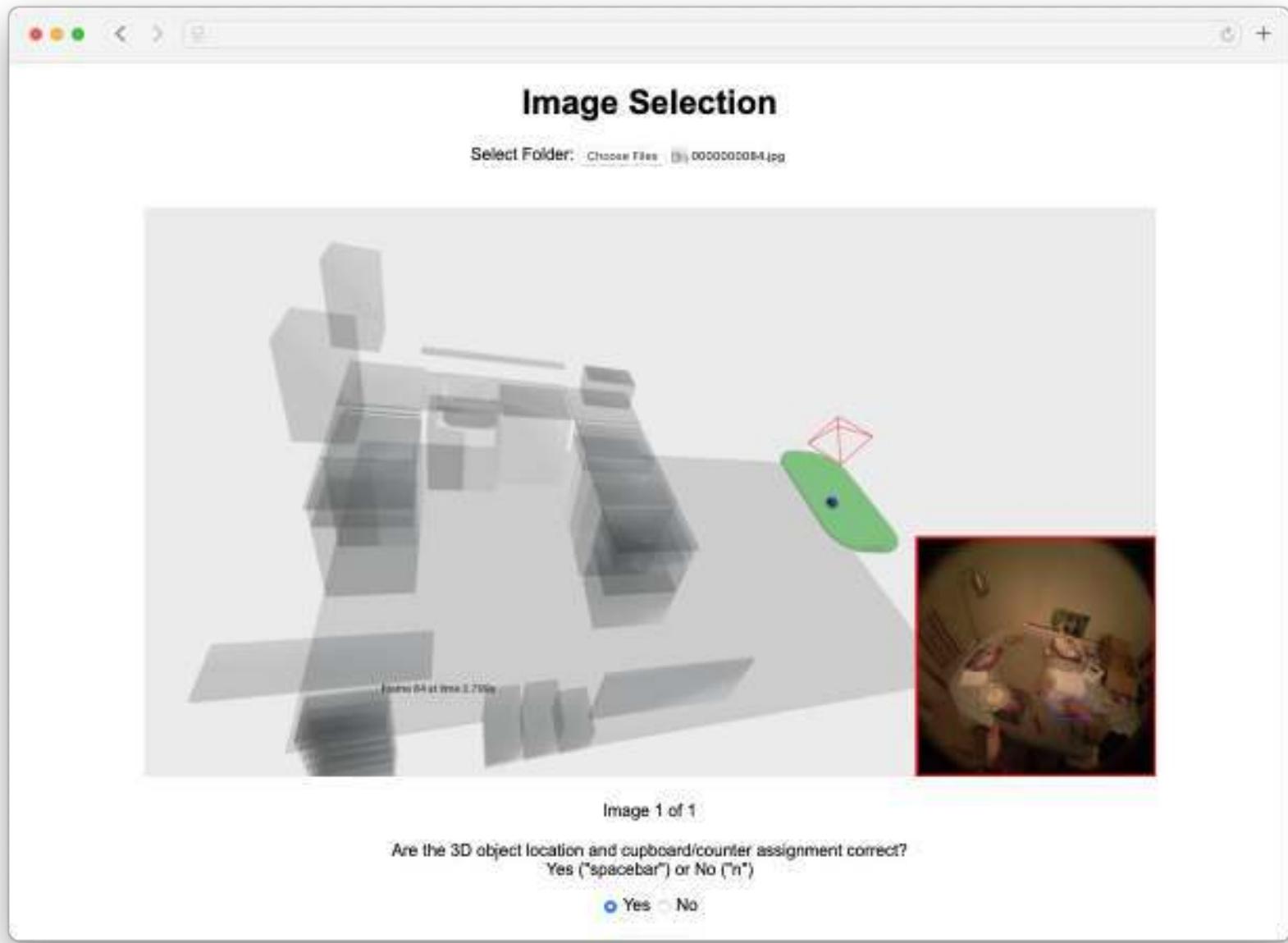


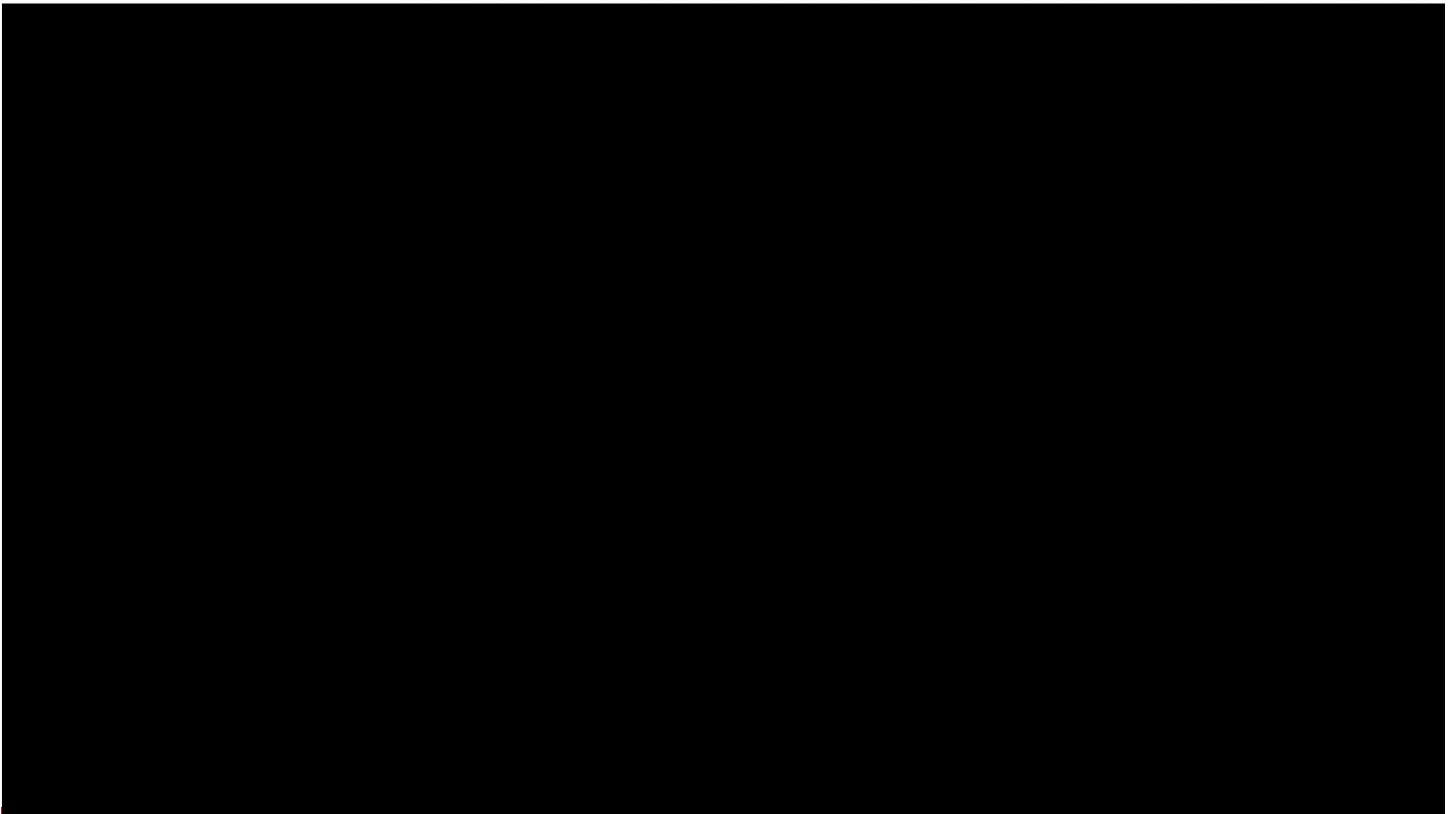
# HD-EPIC





# HD-EPIC







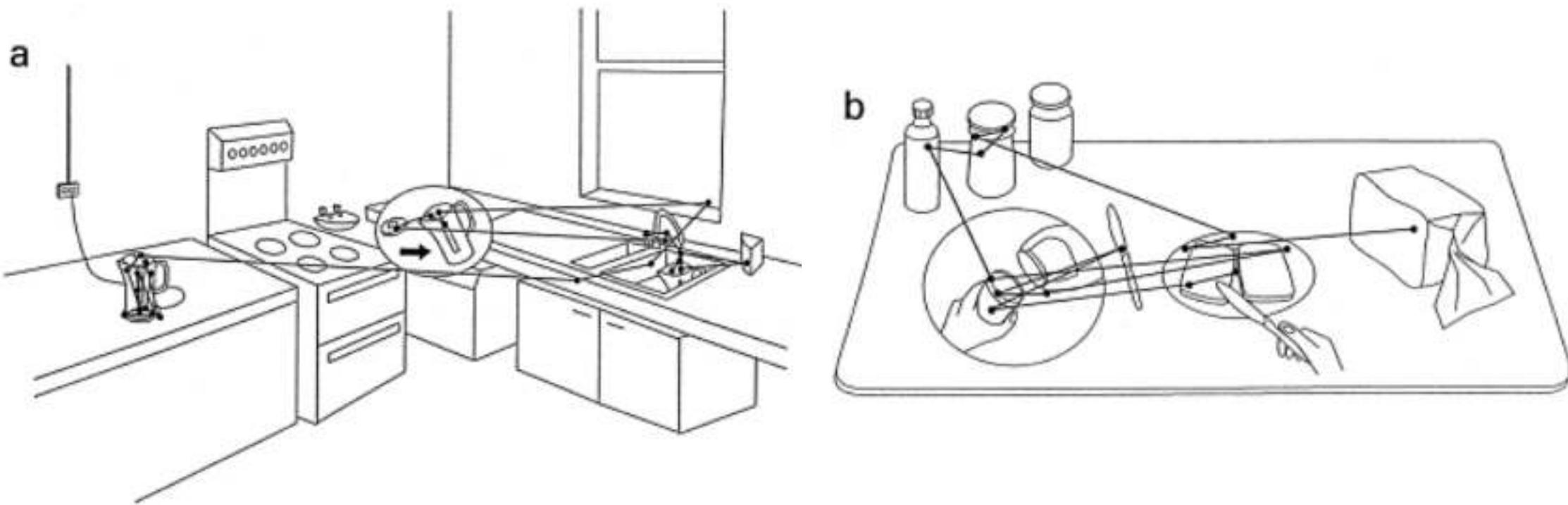
# Video Understanding Out of the Frame

What can we now do with these reconstructions:

- Point Tracking
- Object Tracking
- Gaze Estimation



# Gaze and Fixations

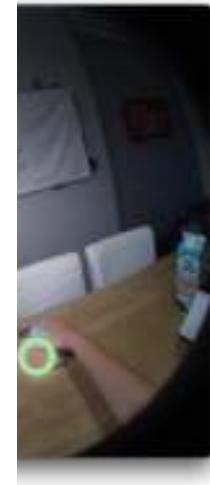
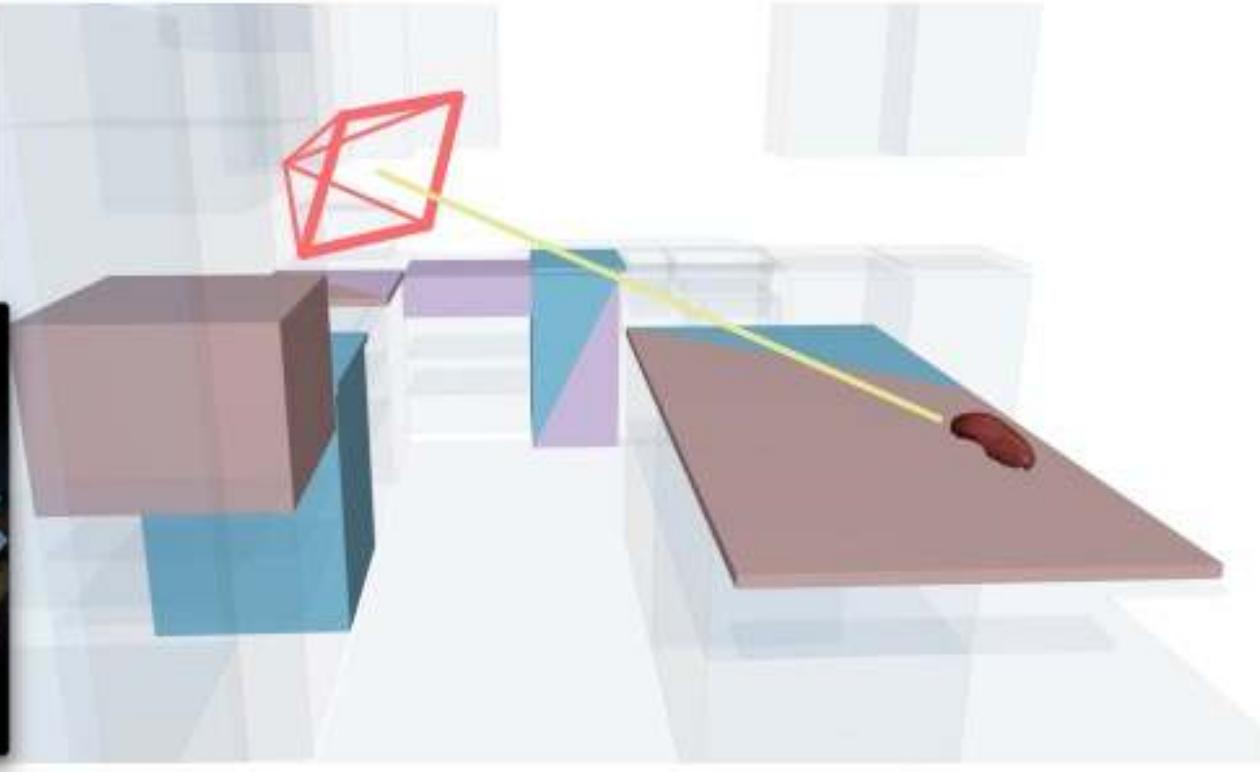




# HD-EPIC

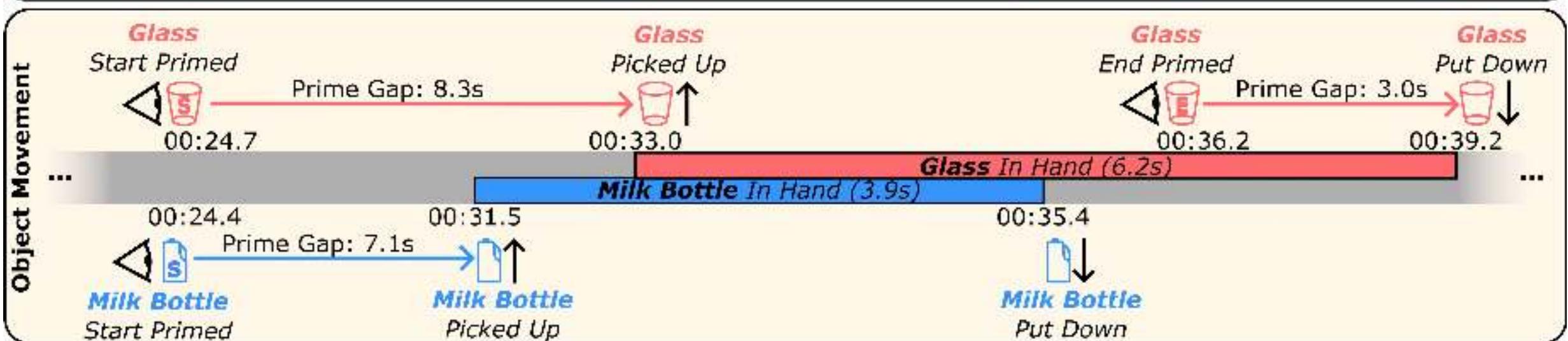


## Gaze priming



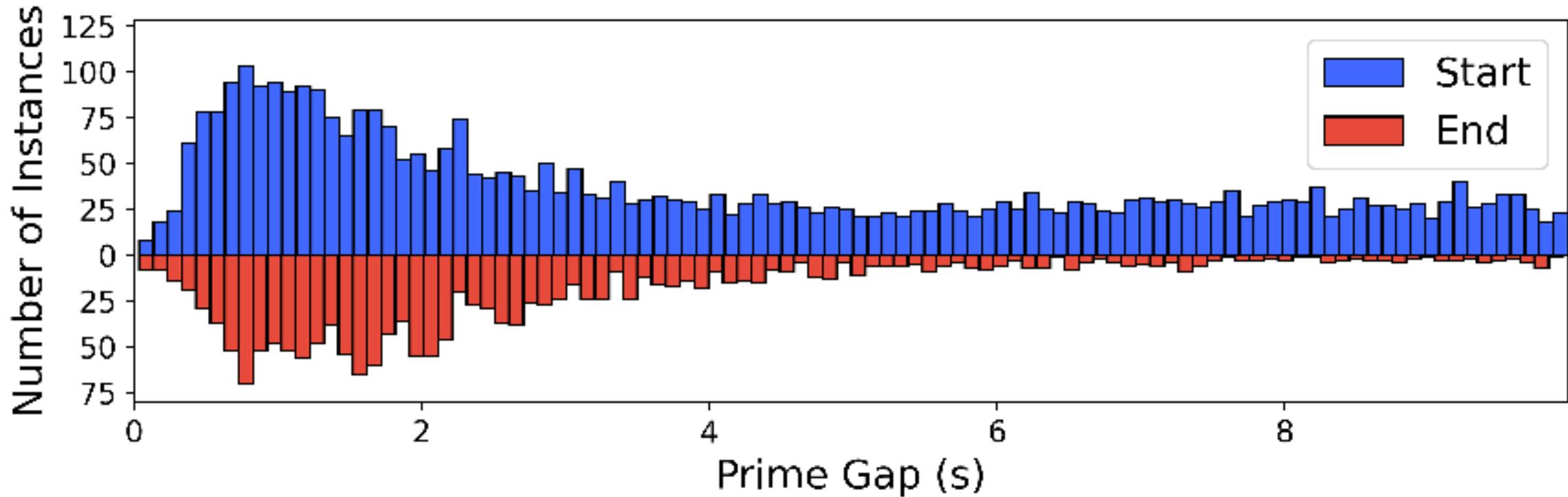


# HD-EPIC





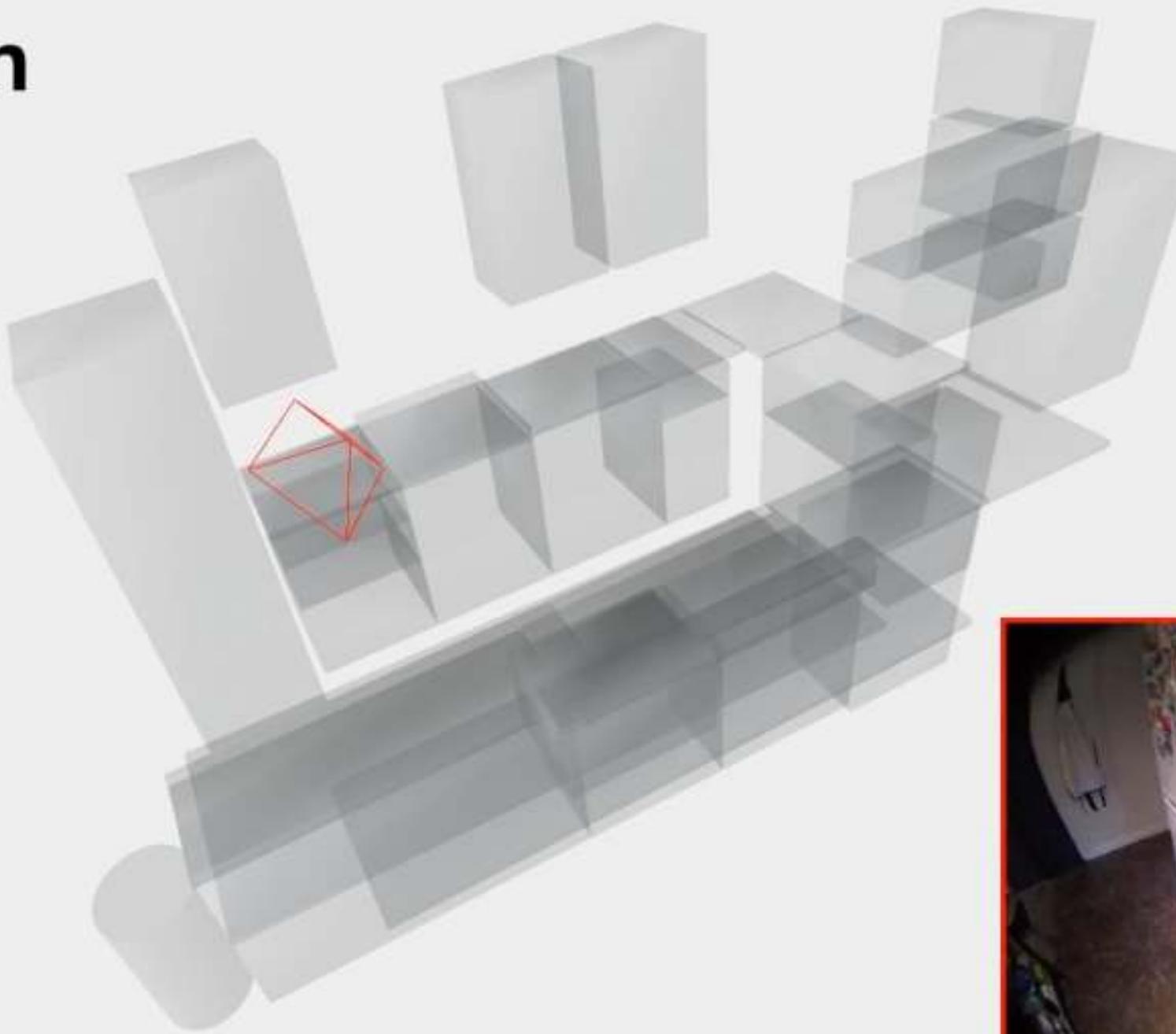
# HD-EPIC



# Digital Twin

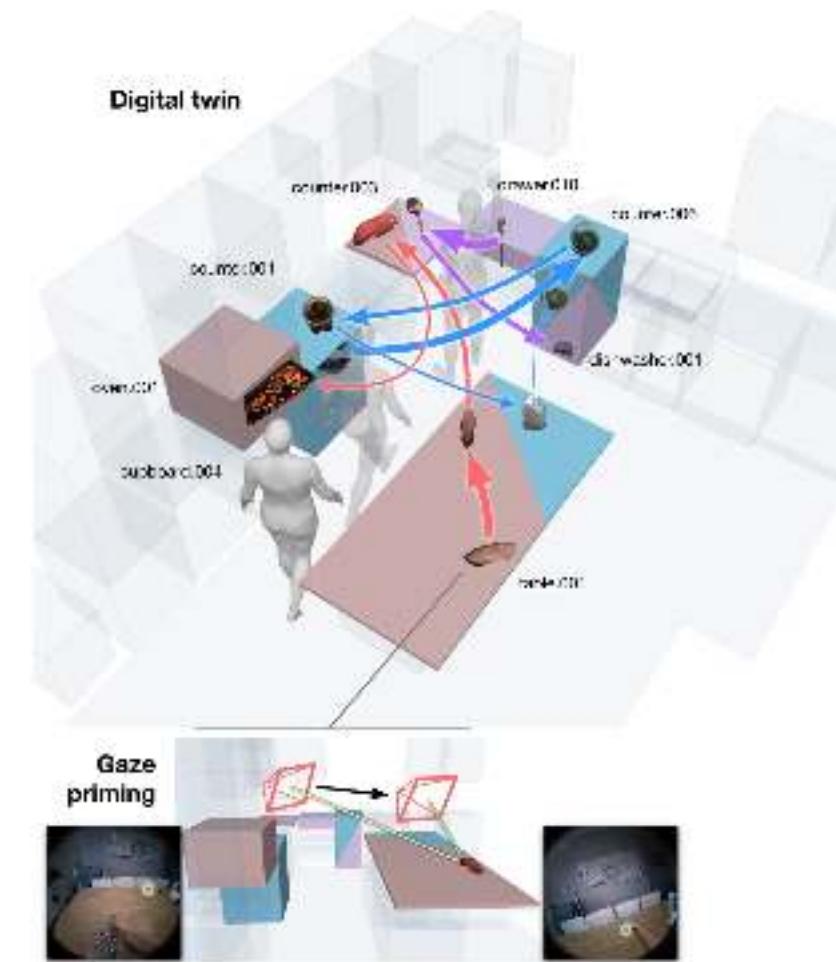
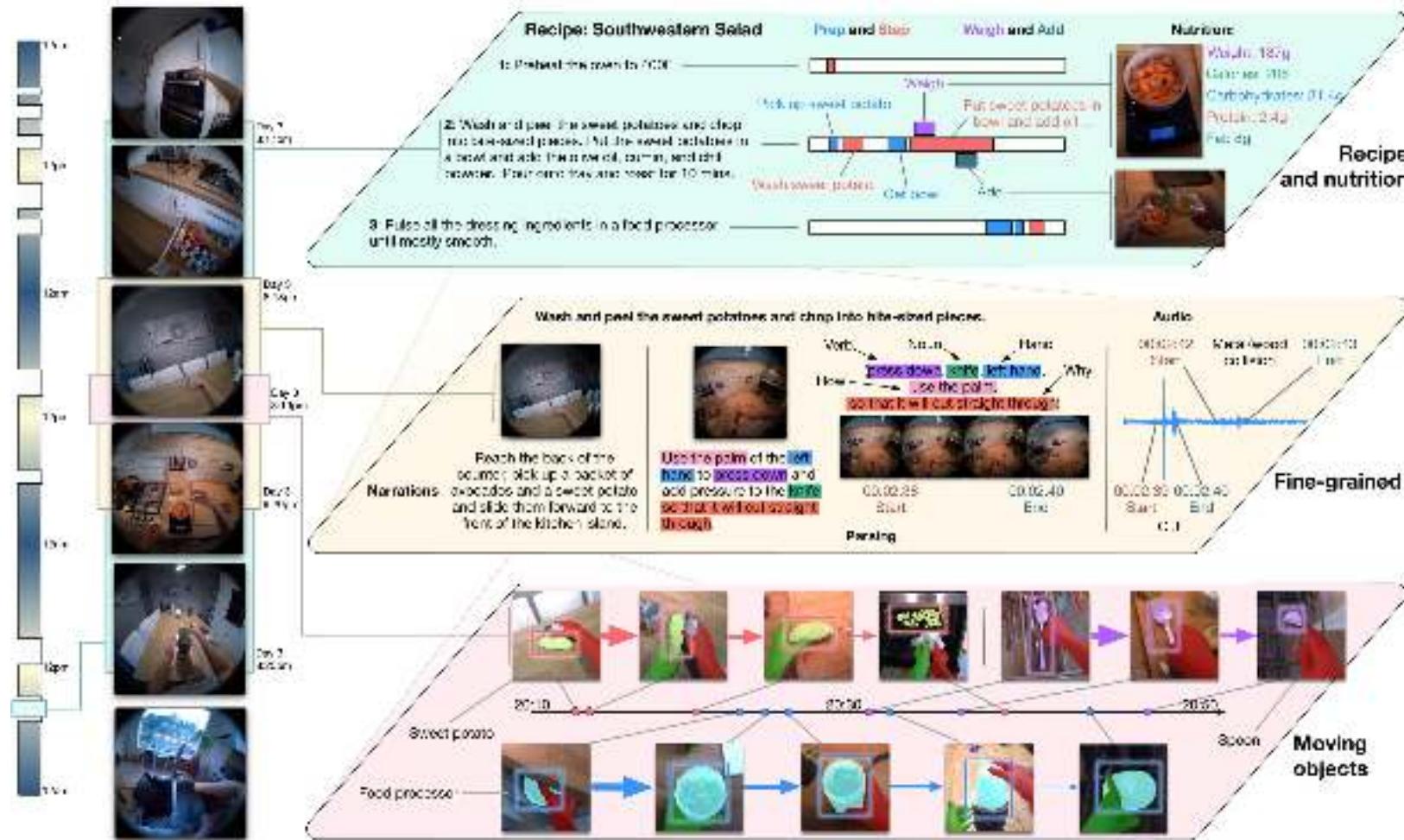
## Fixtures

Open drawer





# HD-EPIC





# HD-EPIC

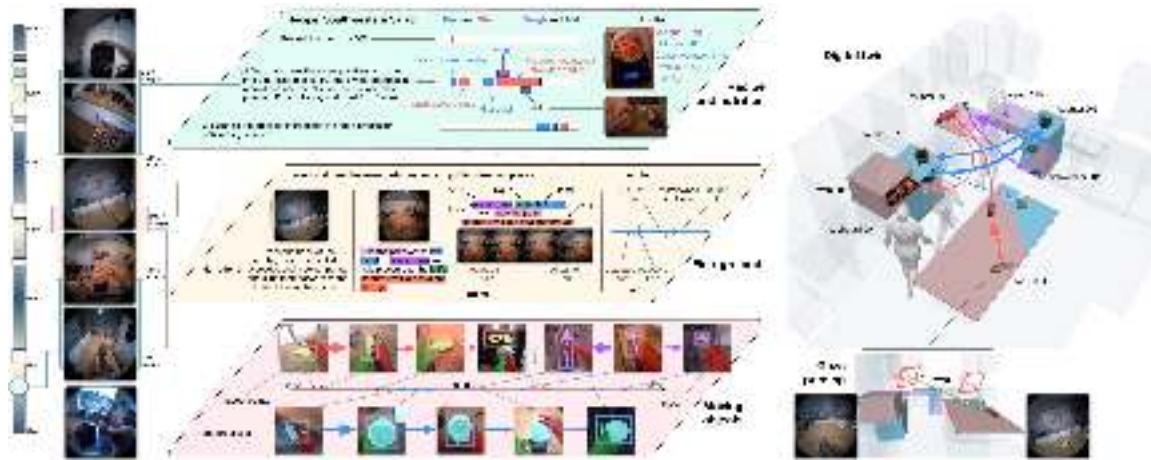


Annotation Type	Total annotations	Annotations/min
Narrations	59,454	24.0
Parsing (Verbs + Nouns + Hands + How + Why)	303,968	122.7
Recipes (Preps + Steps)	4,052	1.6
Sound	50,968	20.6
Action boundaries	59,454	24.0
Object Motion (Pick up + Put down + Fixtures + Bboxes + Masks)	153,480	62.0
Object Itinerary	4,881	2.0
Object Priming (Starts + Ends)	18,264	7.4
Total	263.2	

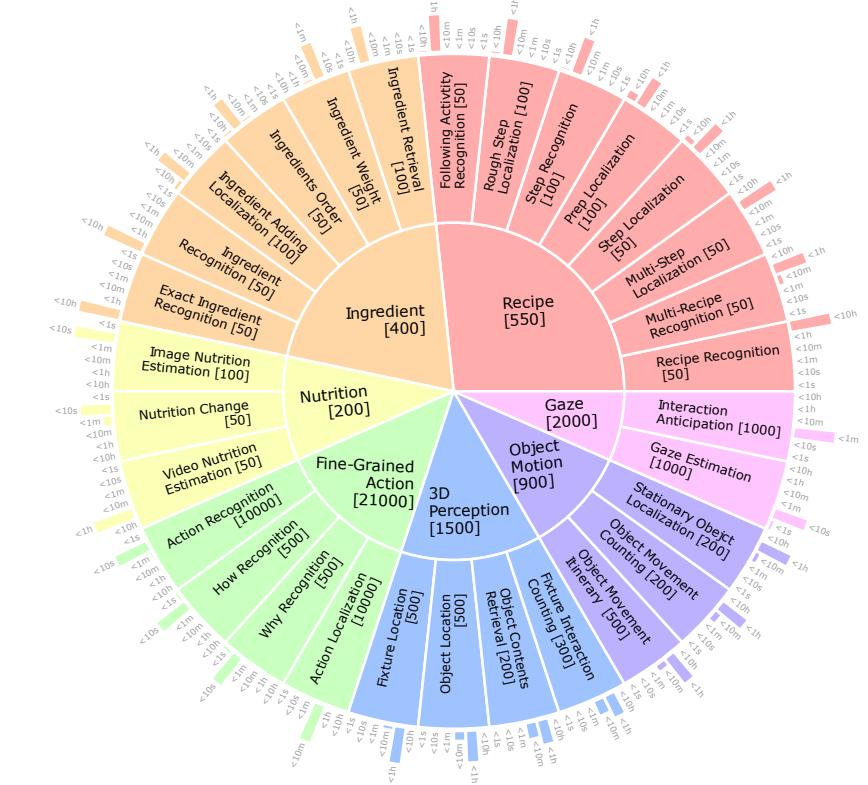
Table A3. HD-EPIC annotations per minute



# HD-EPIC



Sec 1: Highly-Detailed Dataset



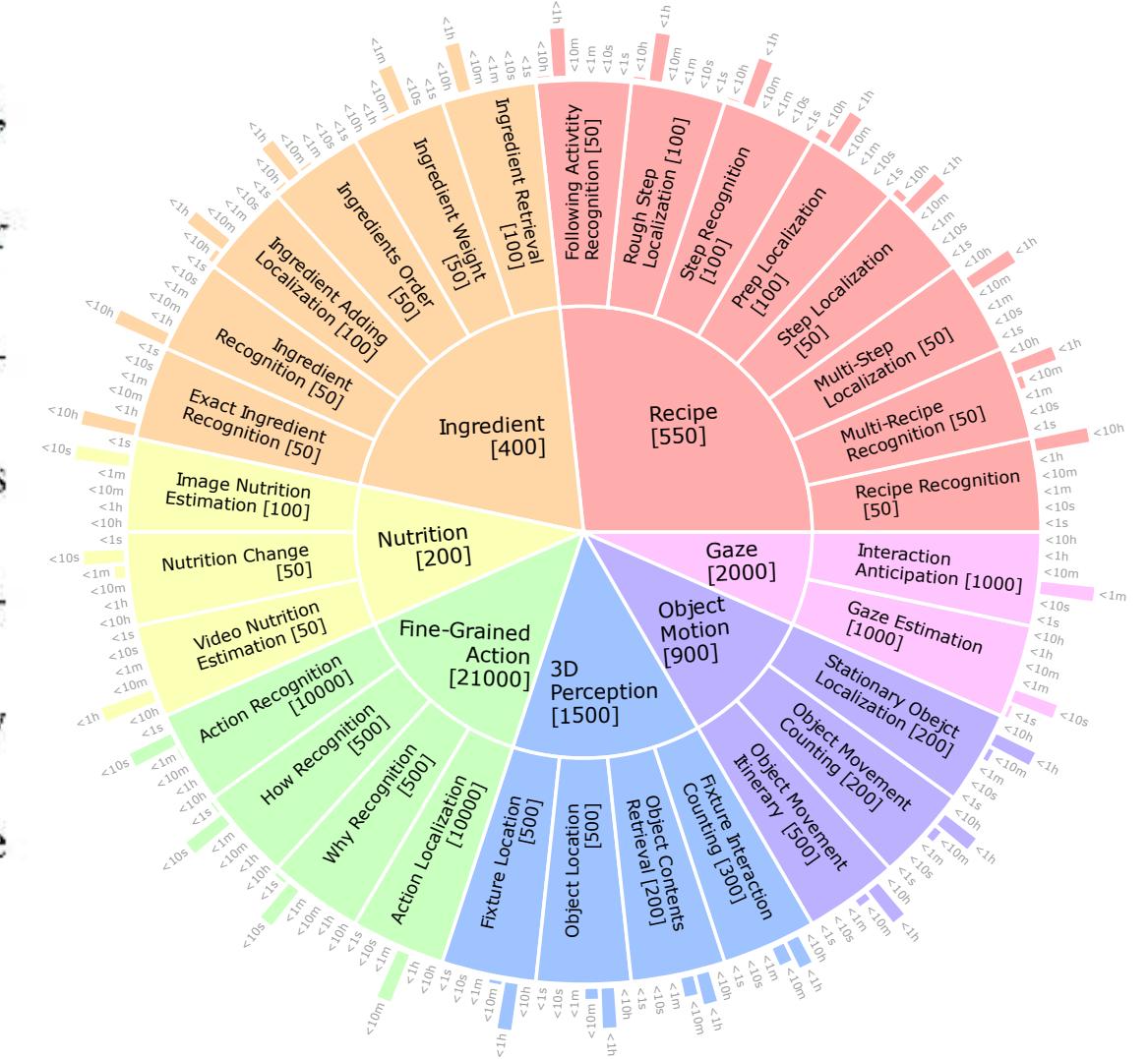
Sec 2: HD-EPIC VQA Benchmark



# HD-EPIC



1. **Recipe**. Questions on temporally localising, retrieving, or recognising recipes and their steps.
2. **Ingredient**. Questions on the ingredients used, their weight, their adding time and order.
3. **Nutrition**. Questions on nutrition of ingredients and nutritional changes as ingredients are added to recipes.
4. **Fine-grained action**. What, how, and why of actions and their temporal localisation.
5. **3D perception**. Questions that require the understanding of relative positions of objects in the 3D scene.
6. **Object motion**. Questions on where, when and how many times objects are moved across long videos.
7. **Gaze**. Questions on estimating the fixation on large landmarks and anticipating future object interactions.





# HD-EPIC

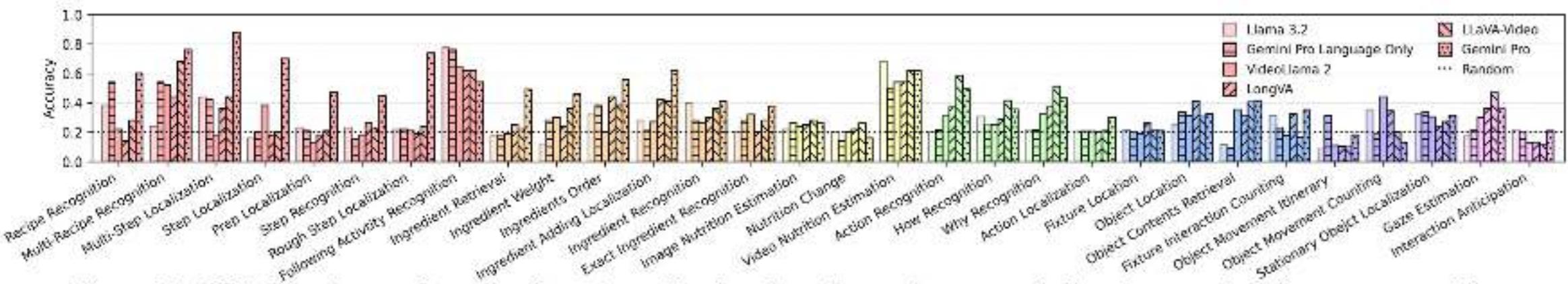
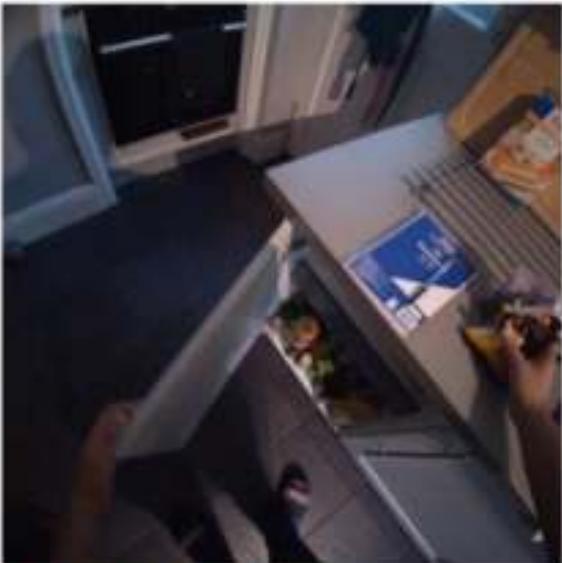
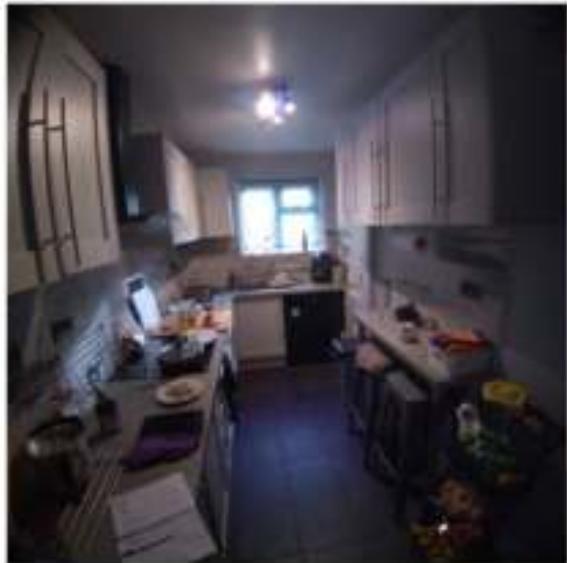


Figure 11. VQA Results per Question Prototype. Our benchmark contains many challenging questions for current models.

Model	Recipe	Ingredient	Nutrition	Action	3D	Motion	Gaze	Avg.
<b>Blind - Language Only</b>								
Llama 3.2	33.5	25.0	36.7	23.3	22.3	25.5	19.5	26.5
Gemini Pro	38.0	26.8	30.0	22.1	21.5	27.7	20.5	26.7
<b>Video-Language</b>								
VideoLlama 2	30.8	25.7	32.7	27.2	25.7	28.5	21.2	27.4
LongVA	29.6	30.8	33.7	30.7	32.9	22.7	24.5	29.3
LLaVA-Video	36.3	33.5	38.7	43.0	27.3	18.9	29.3	32.4
Gemini Pro	64.3	48.6	34.7	39.6	32.5	20.8	28.7	38.5
<i>Sample Human Baseline</i>	96.7	96.7	85.0	92.5	93.8	92.7	75.0	90.3



# HD-EPIC



How many times did I **open** the item at bounding box **(165, 452, 1408, 1408)** in  
**00:00:57?**

A. 3

B. 1

C. 4

D. 5

E. 2



# HD-EPIC





# HD-EPIC





HD-EPIC

## HD-EPIC

A Highly-Detailed Egocentric Video Dataset

[Paper \(ArXiv\)](#) ↗

[The Dataset](#)

[Explore Samples](#)

[Watch Video](#) ↗

[VQA benchmark](#)

[Explore VQA](#)

[Download](#)

[Team](#)



## News

- May 2025: Eye-Gaze Priming data has now been released! [Annotations link](#) ↗
- April 2026: VQA Challenge Benchmark is online now! [Challenge link](#) ↗.
- April 2025: Masks and object association annotations have now been released.
- Feb 2025: HD-EPIC ↗ accepted at CVPR 2025!



HD-EPIC



Try it Yourself

12 VQAs  
Focusing on 3D



[https://hd-epic.github.io/icvss\\_demo.html](https://hd-epic.github.io/icvss_demo.html)



HD-EPIC



Try it Yourself

Use Wise to Search  
through HD-EPIC

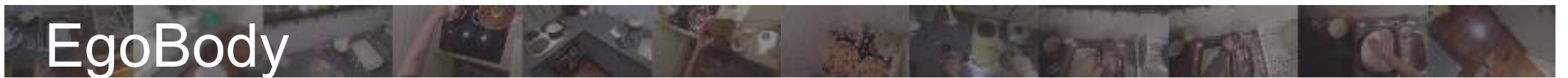


<https://meru.robots.ox.ac.uk/HD-EPIC/>



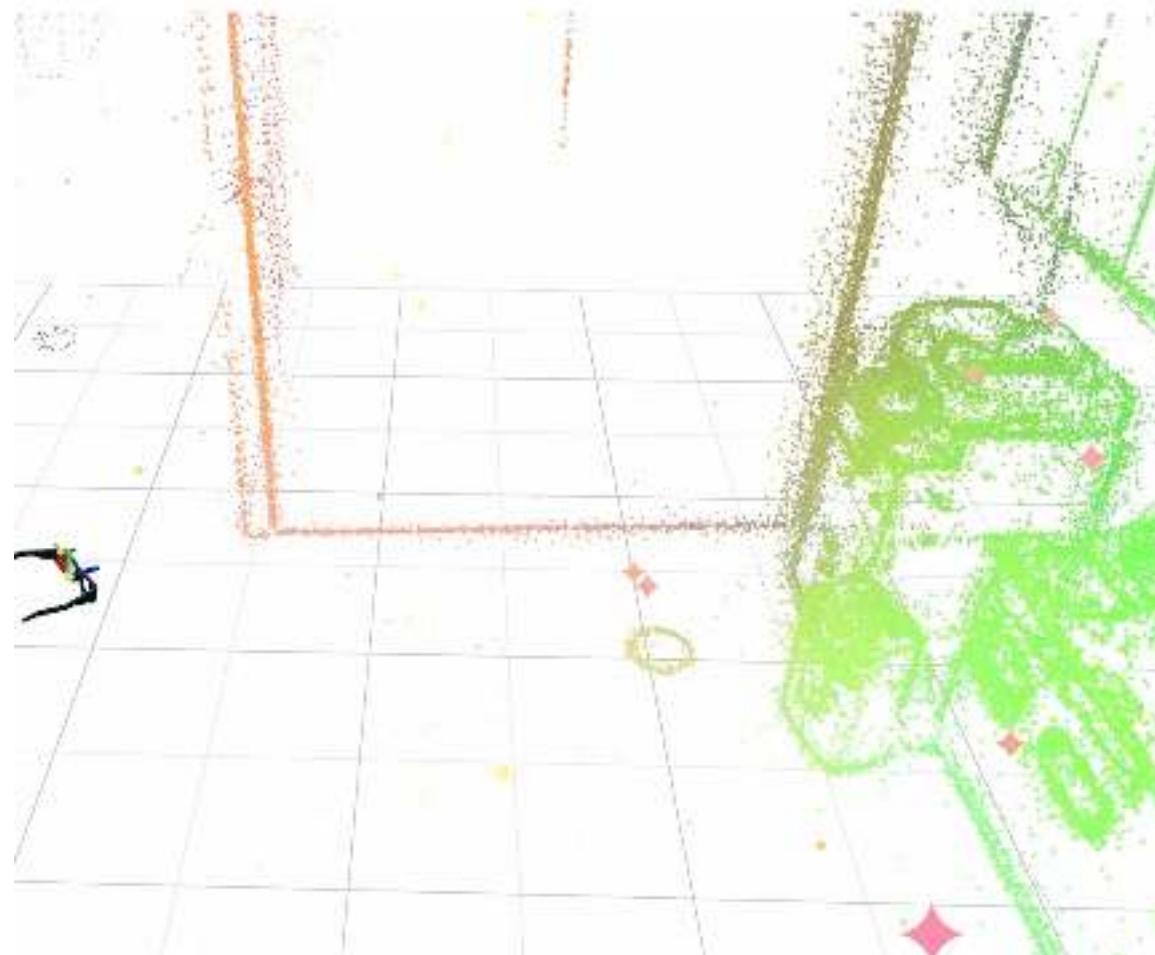
# Video Understanding Out of the Frame

Body and Hand Motion Estimation “out of the frame”



# EgoBody

EgoAllo uses egocentric ( 6d ) SLAM poses and images

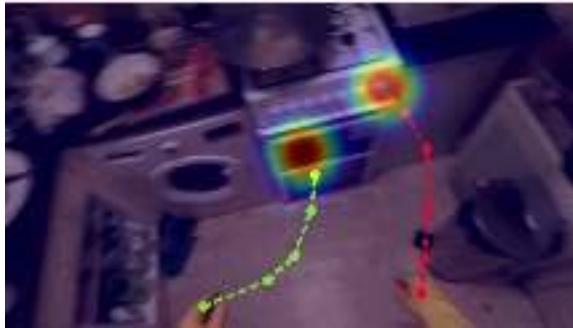


# EgoHand Forecasting – Previous Works

with: Masashi Hatano  
Zhifan Zhu  
Hideo Saito

## 2D Hand Forecasting

Given an egocentric video,  
forecast 2D hand positions of both hands  
→ Limited in 2D image plane



OCT [CVPR'22]



Diff-IP2D [IROS'25]

## 3D Hand Forecasting

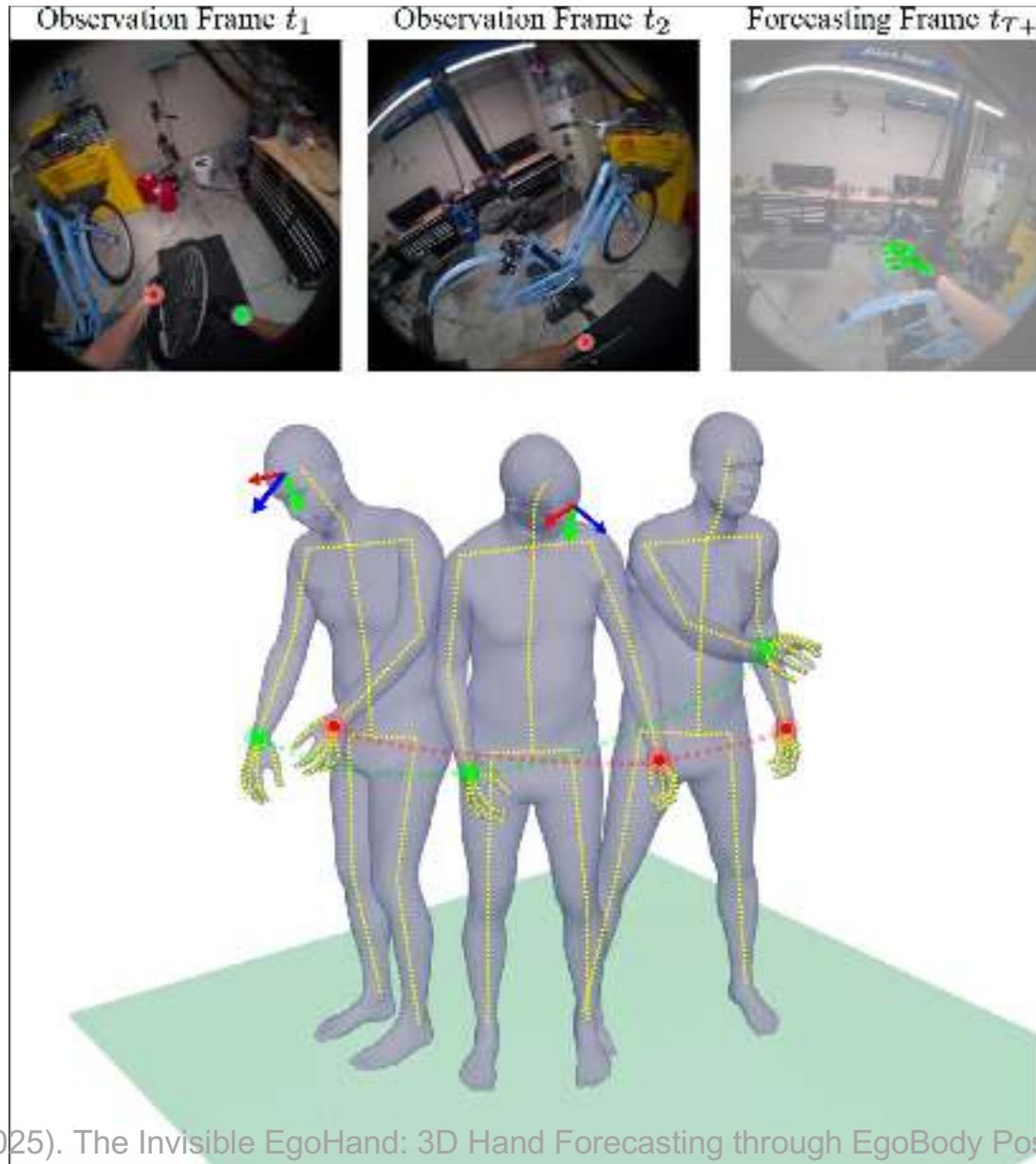
Given an egocentric video & 3D hand trajectory,  
forecast 3D hand positions of one hand



USST [ICCV'23]

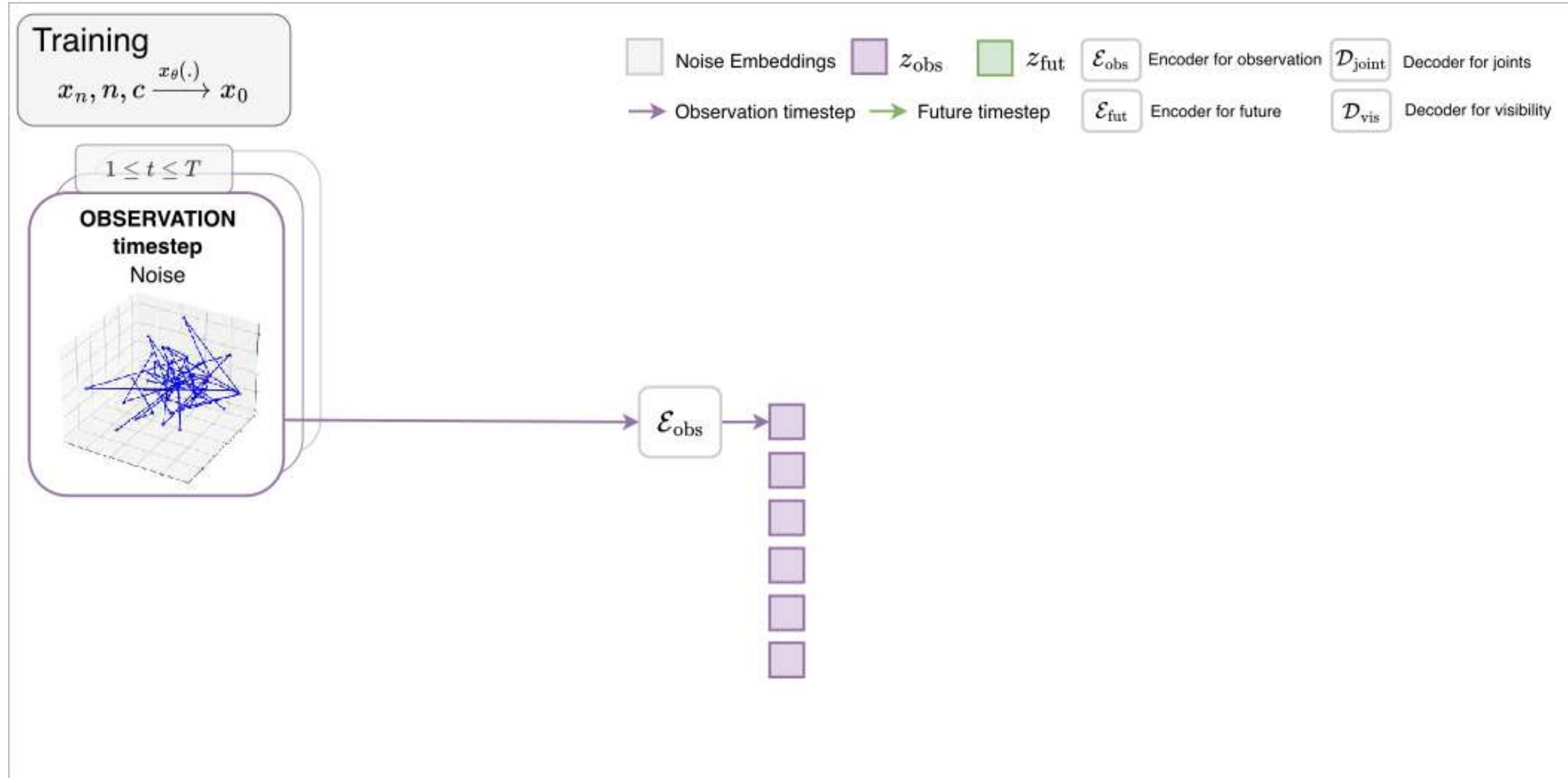
# The Invisible EgoHand

with: Masashi Hatano  
Zhifan Zhu  
Hideo Saito



# The Invisible EgoHand

with: Masashi Hatano  
Zhifan Zhu  
Hideo Saito



# The Invisible EgoHand

with: Masashi Hatano  
Zhifan Zhu  
Hideo Saito

Method	Hand Trajectory Forecasting				Hand Pose Forecasting			
			All				All	
	ADE	FDE	MPJPE	MPJPE-F				
Static	0.335	0.405	0.166	0.179				
CVM [61]	0.346	0.467	0.166	0.183				
EgoEgoForecast	0.295	0.352	0.166	0.177				
USST [3]	0.562	0.581	-	-				
Ours	<b>0.261</b>	<b>0.324</b>	<b>0.115</b>	<b>0.143</b>				

# The Invisible EgoHand

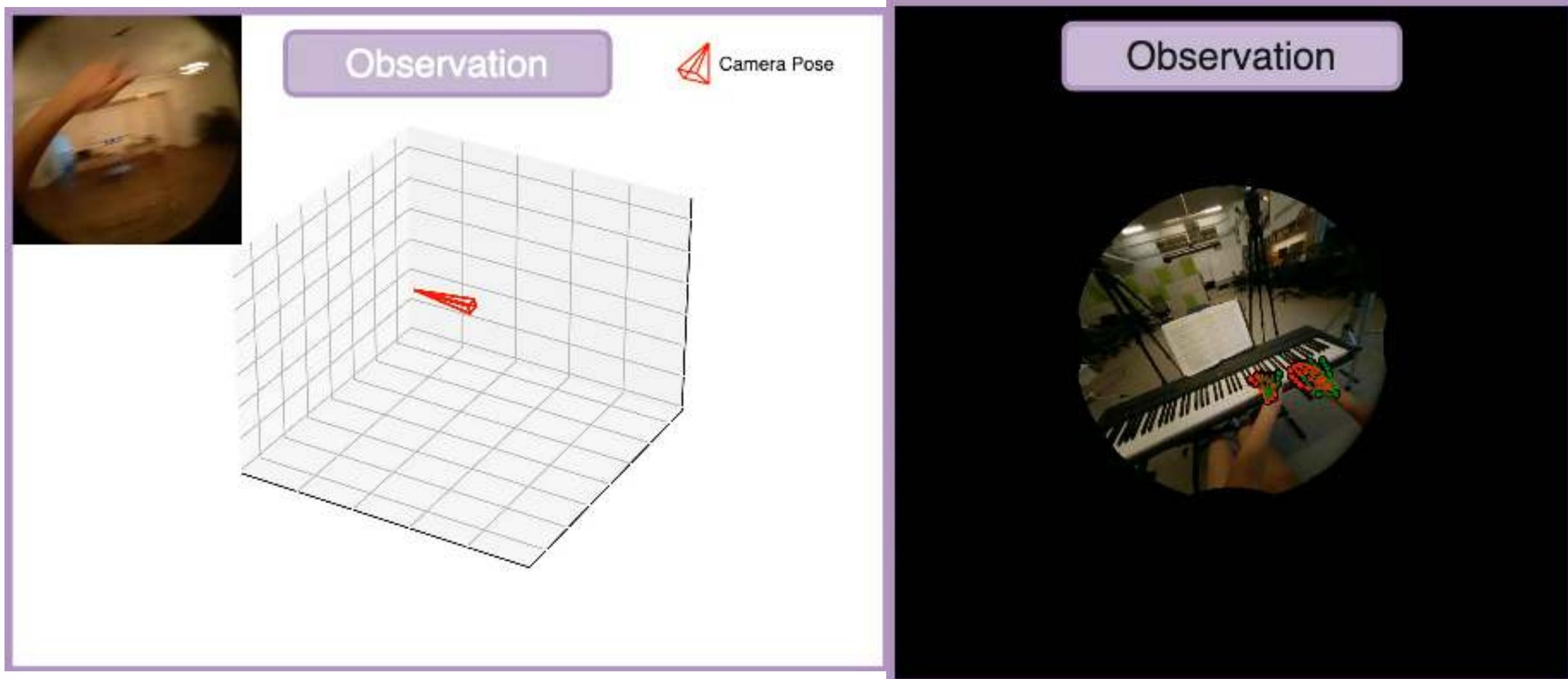
with: Masashi Hatano  
Zhifan Zhu  
Hideo Saito

Method	Hand Trajectory Forecasting			Hand Pose Forecasting		
	In-view	Out-of-view	All	In-view	Out-of-view	All
EgoEgoForecast	0.171	0.385	0.295	0.162	0.299	0.166
Ours w/o. 2D joint	0.151	0.377	0.282	0.139	0.269	0.142
Ours w/o. image	<b>0.116</b>	0.367	<b>0.261</b>	0.117	<b>0.234</b>	0.120
Ours w/o. $\mathcal{L}_{\text{reproj}}$	0.132	0.368	0.269	0.125	0.250	0.128
Ours w/o. $\mathcal{L}_{\text{vis}}$	0.127	0.377	0.272	0.121	0.240	0.124
Ours w/o. $\mathcal{L}_{\text{body}}$	0.129	0.385	0.277	0.120	0.258	0.123
Ours w/o. $\mathcal{L}_{\text{obs}}$	0.149	0.390	0.289	0.139	0.250	0.142
<b>Ours</b>	<b>0.116</b>	<b>0.366</b>	<b>0.261</b>	<b>0.112</b>	0.240	<b>0.115</b>

- Without visible 2D joints, significant performance drops can be seen
- 2D reprojection loss serves as effective regularization
- Visibility loss & Body joints loss contribute for out-of-view scenario

# The Invisible EgoHand

with: Masashi Hatano  
Zhifan Zhu  
Hideo Saito





# Video Understanding Out of the Frame

From First-Point View to Second- and Third-

# FPV with SPV

Input: paired egocentric videos

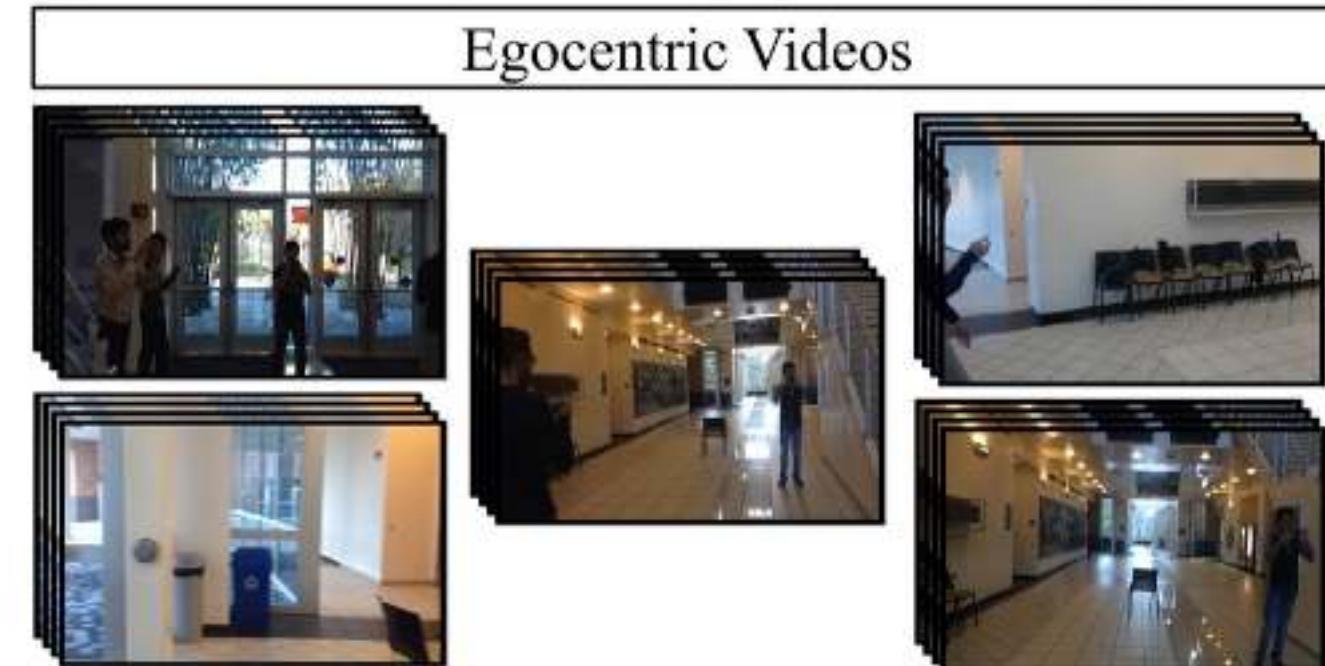


Egocentric video of person A

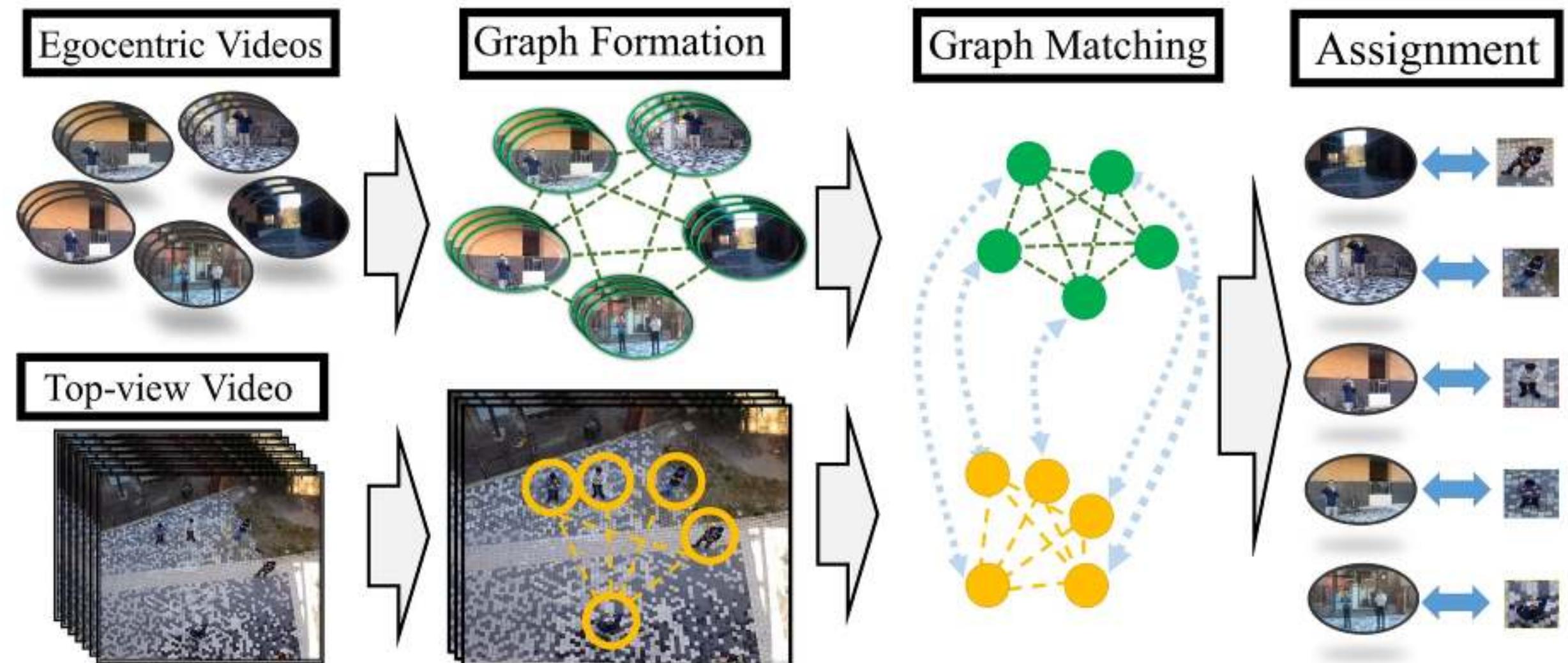


Egocentric video of person B

# FPV with TPV (top-view)



# FPV with TPV (top-view)



# Ego-Exo4D

with: Kristen Grauman  
+102 authors



# Ego-Exo4D

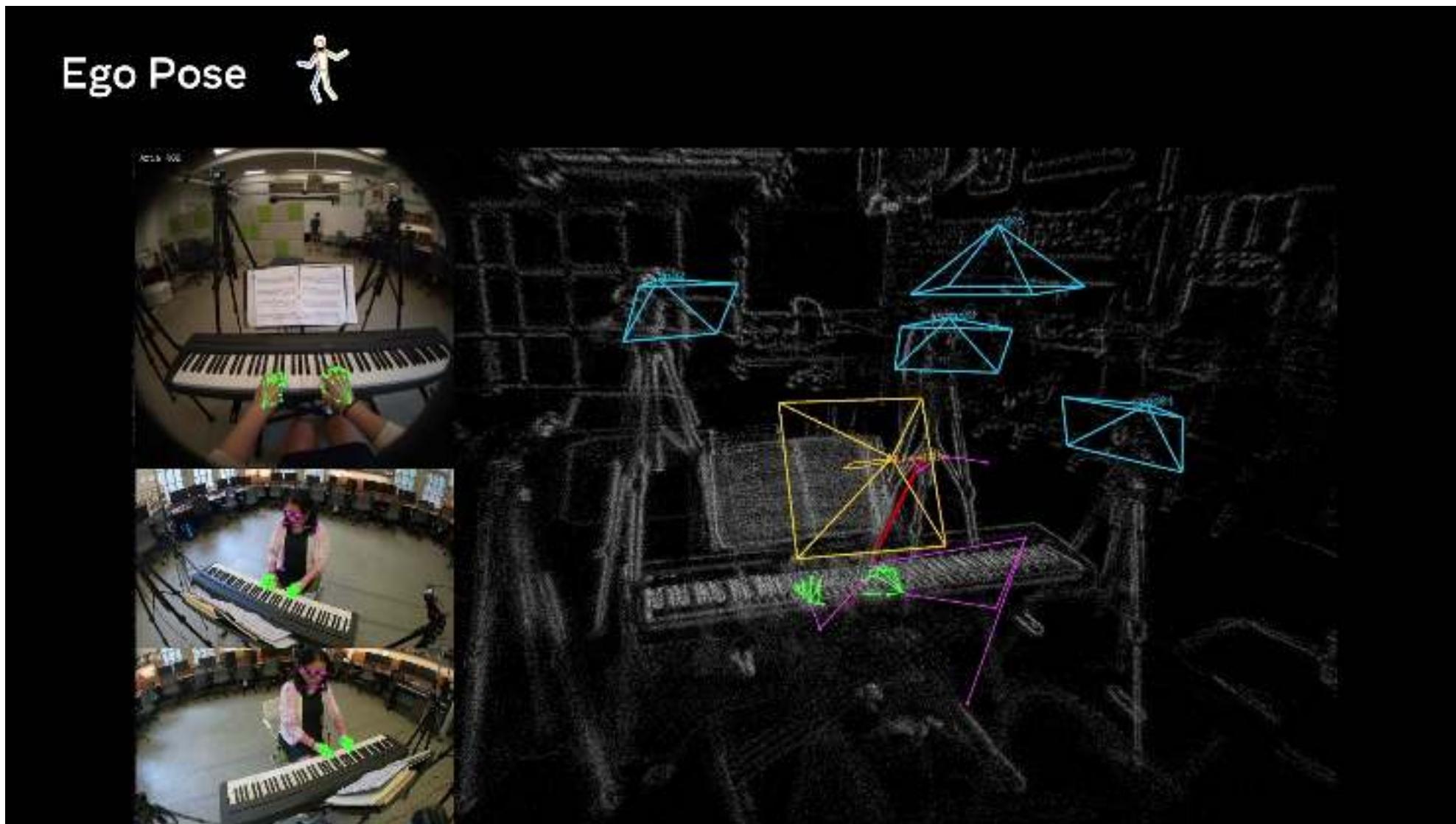
with: Kristen Grauman  
+102 authors

## Ego-Exo Relation

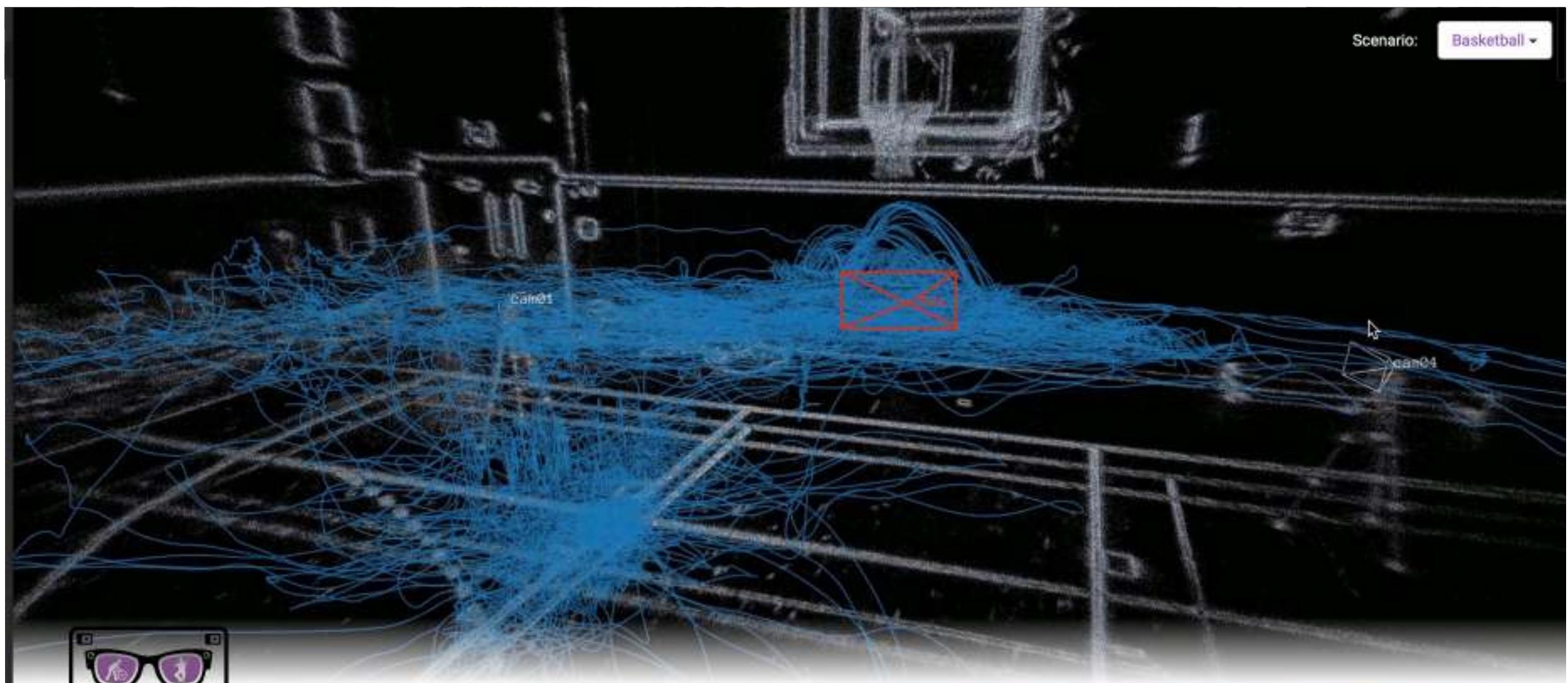


# Ego-Exo4D

with: Kristen Grauman  
+102 authors



Scenario: Basketball



## EGO-EXO4D

A diverse, large-scale multi-modal, multi-view, video dataset and benchmark collected across 13 cities worldwide by 839 camera wearers, capturing 1422 hours of video of skilled human activities.

*Hover your mouse over scene cameras above to see a sample video for the chosen scenario.*

[Learn More](#) ↓

[Watch Video](#) ↗

[Start Here](#) ↗

# In today's tutorial



Motivation and Datasets in  
Egocentric Video Understanding



Video Understanding  
Out of the Frame



Video Understanding:  
Data and Tasks



Teaser: The Wizard of Oz  
at the Sphere



Videos are Multimodal



Outlook into the Future of  
Egocentric Vision



Connected Videos of One's Life



Conclusion

# The Wizard of Oz at the Sphere

Coming in August 2025

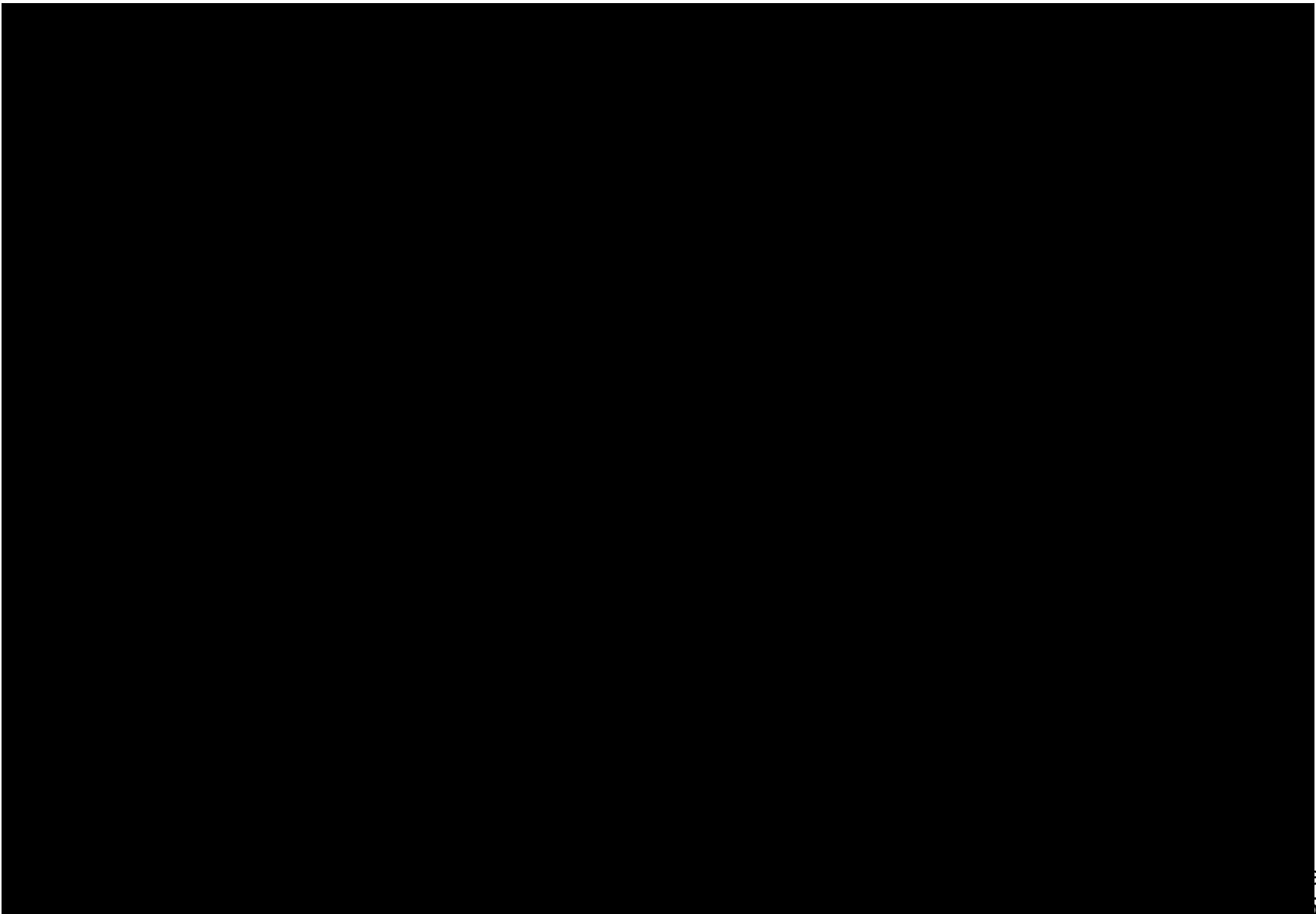


<https://behindthecurtain.withgoogle.com>

Dima Damen  
ICVSS2025

# The Wizard of Oz @ The Sphere

- The Movie (1939)
- Technicolour pioneer
- Iconic characters



# The Wizard of Oz @ The Sphere



**Ralph Winte**

Head of Physical Production -  sphere

Dima Damen  
ICVSS2025

# The Wizard of Oz @ The Sphere



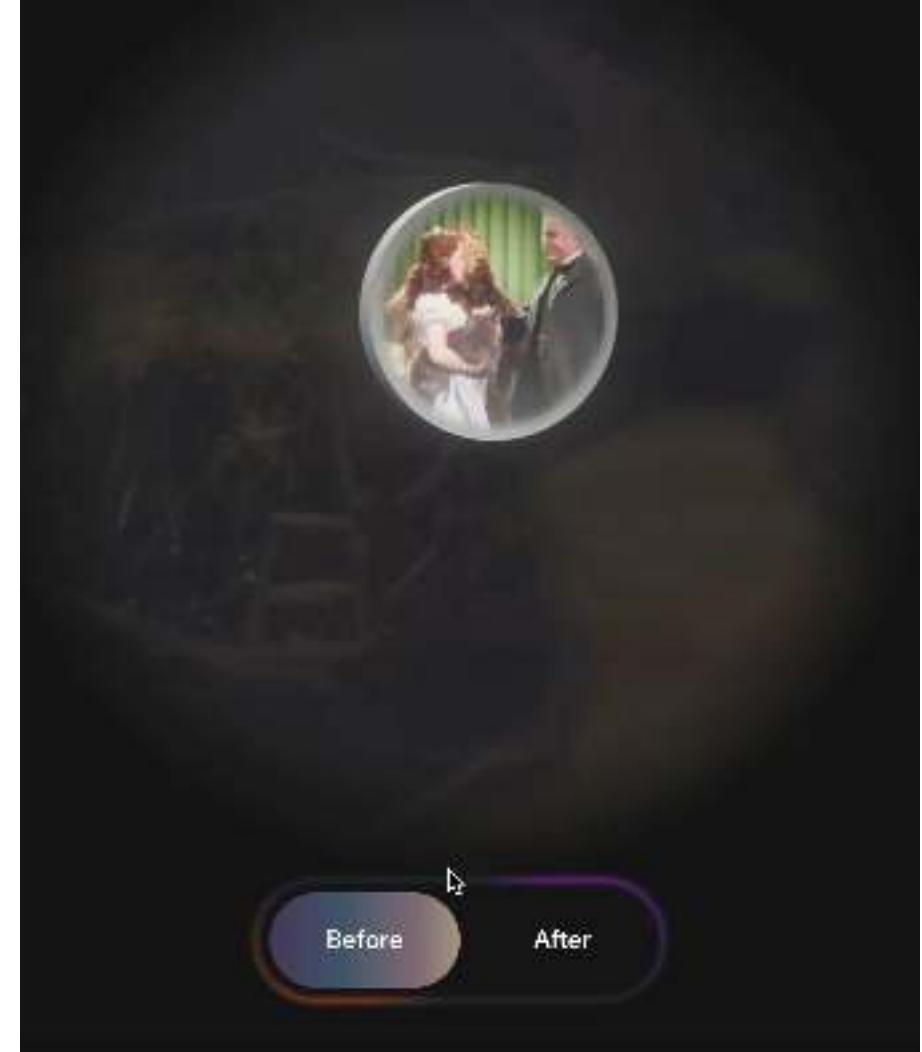
Dima Damen  
ICVSS2025

# The Wizard of Oz @ The Sphere

- Super-resolution,
- Outpainting...



<https://behindthecurtain.withgoogle.com>



Dima Damen  
ICVSS2025

# The Wizard of Oz @ The Sphere

- Performance Interpolation,



One of the most ambitious challenges was addressing what the teams called the "performance gap" -

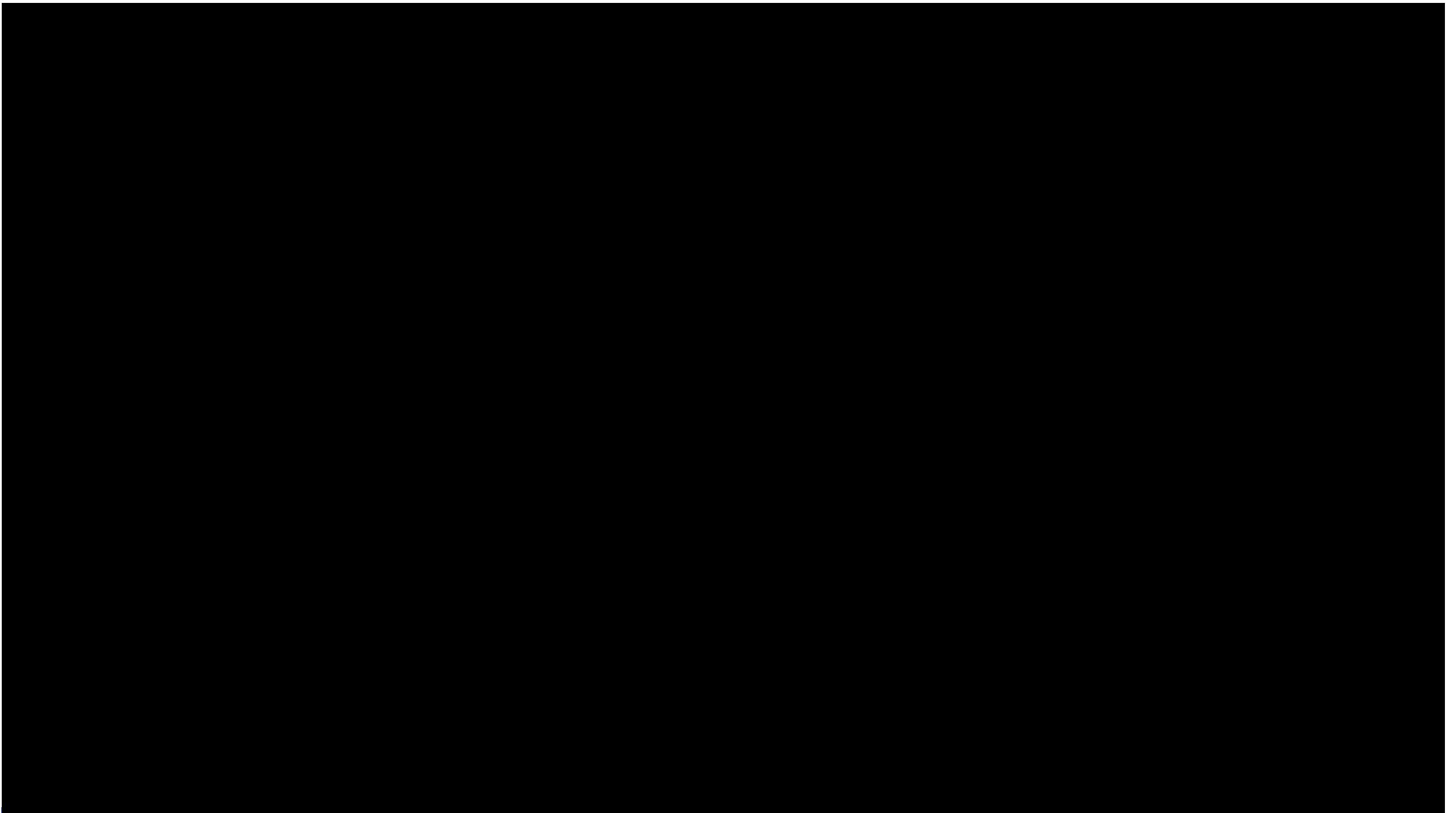
# The Wizard of Oz @ The Sphere

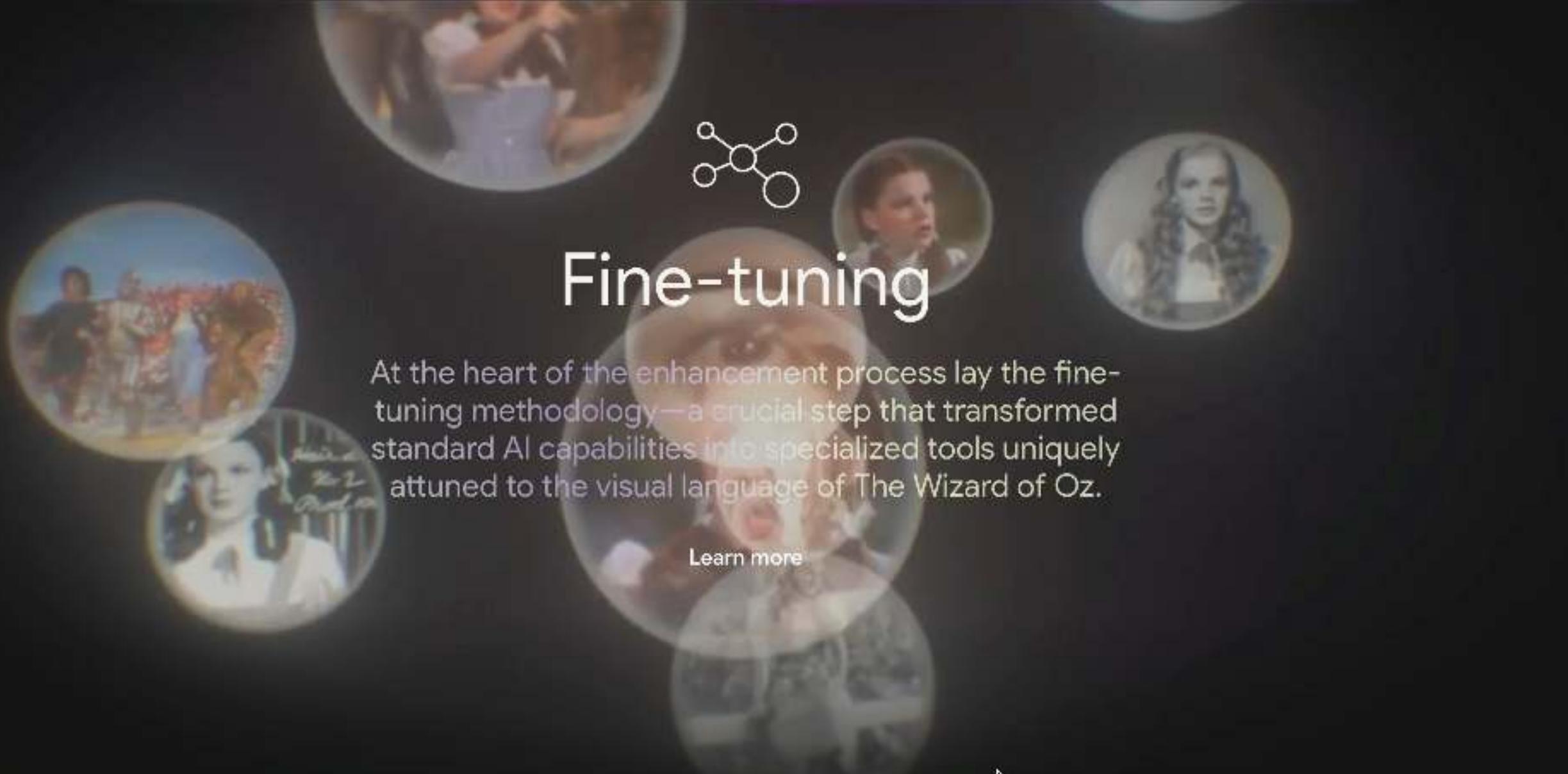
- Auto-Director



<https://behindthecurtain.withgoogle.com>

Dima Damen  
ICVSS2025





# Fine-tuning

At the heart of the enhancement process lay the fine-tuning methodology—a crucial step that transformed standard AI capabilities into specialized tools uniquely attuned to the visual language of *The Wizard of Oz*.

[Learn more](#)



<https://behindthecurtain.withgoogle.com>

Dima Damen  
ICVSS2025

# In today's tutorial



Motivation and Datasets in  
Egocentric Video Understanding



Video Understanding  
Out of the Frame



Video Understanding:  
Data and Tasks



Teaser: The Wizard of Oz  
at the Sphere



Videos are Multimodal



Outlook into the Future of  
Egocentric Vision



Connected Videos of One's Life



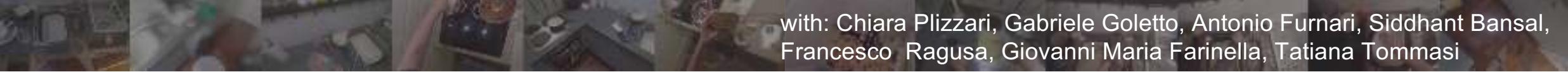
Conclusion



# An Outlook into the Future of Egocentric Vision

Chiara Plizzari\*, Gabriele Goletto\*, Antonino Furnari\*, Siddhant Bansal\*, Francesco Ragusa\*, Giovanni Maria Farinella<sup>†</sup>, Dima Damen<sup>†</sup>, Tatiana Tommasi<sup>†</sup>





with: Chiara Plizzari, Gabriele Goletto, Antonio Furnari, Siddhant Bansal,  
Francesco Ragusa, Giovanni Maria Farinella, Tatiana Tommasi

# Envisioning an Ambitious Future and Analysing the Current Status of Egocentric Vision

How did we do this?

We imagined a device – *EgoAI* and envisioned its utility in multiple scenarios



**EGO-Designer**



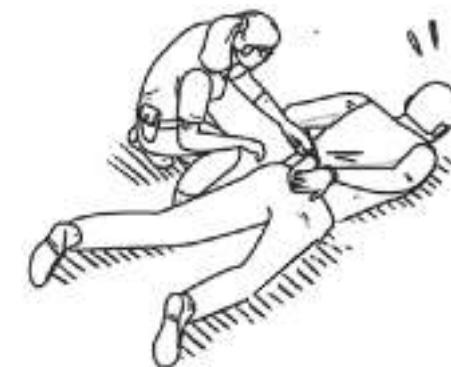
**EGO-Worker**



**EGO-Tourist**



**EGO-Home**

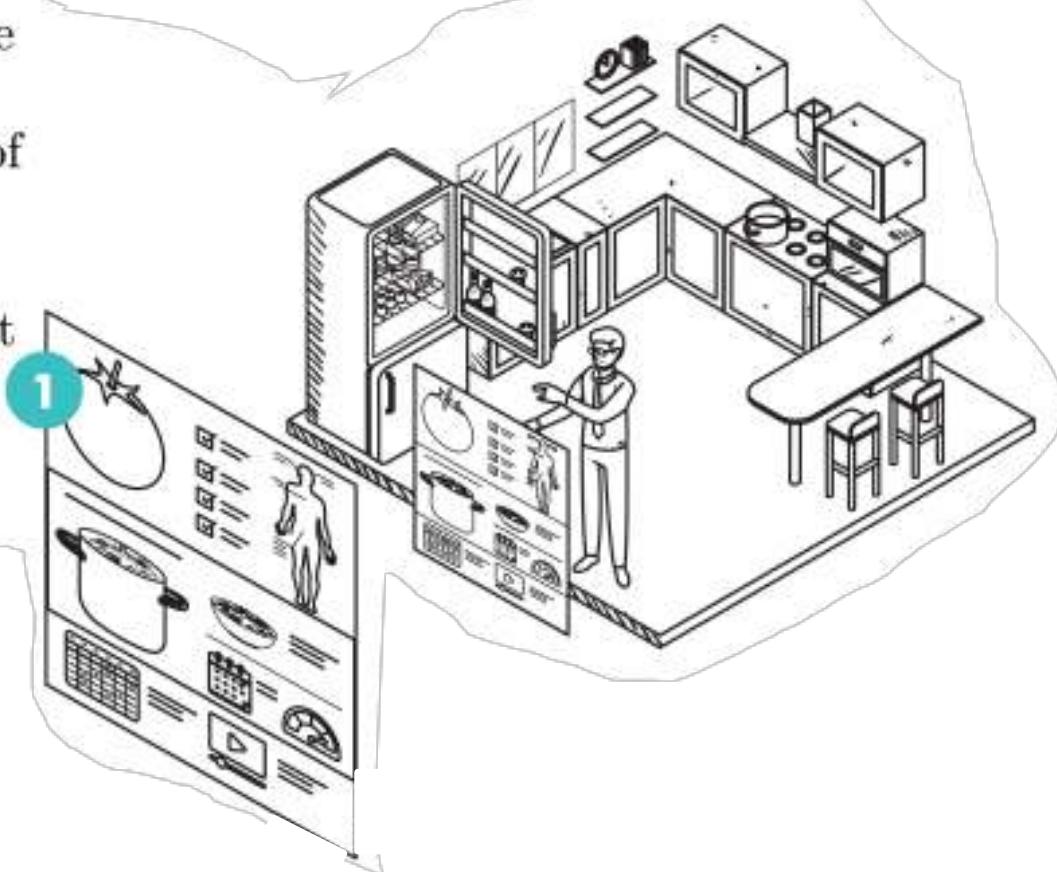


**Ego-Police**

# EGO-Home

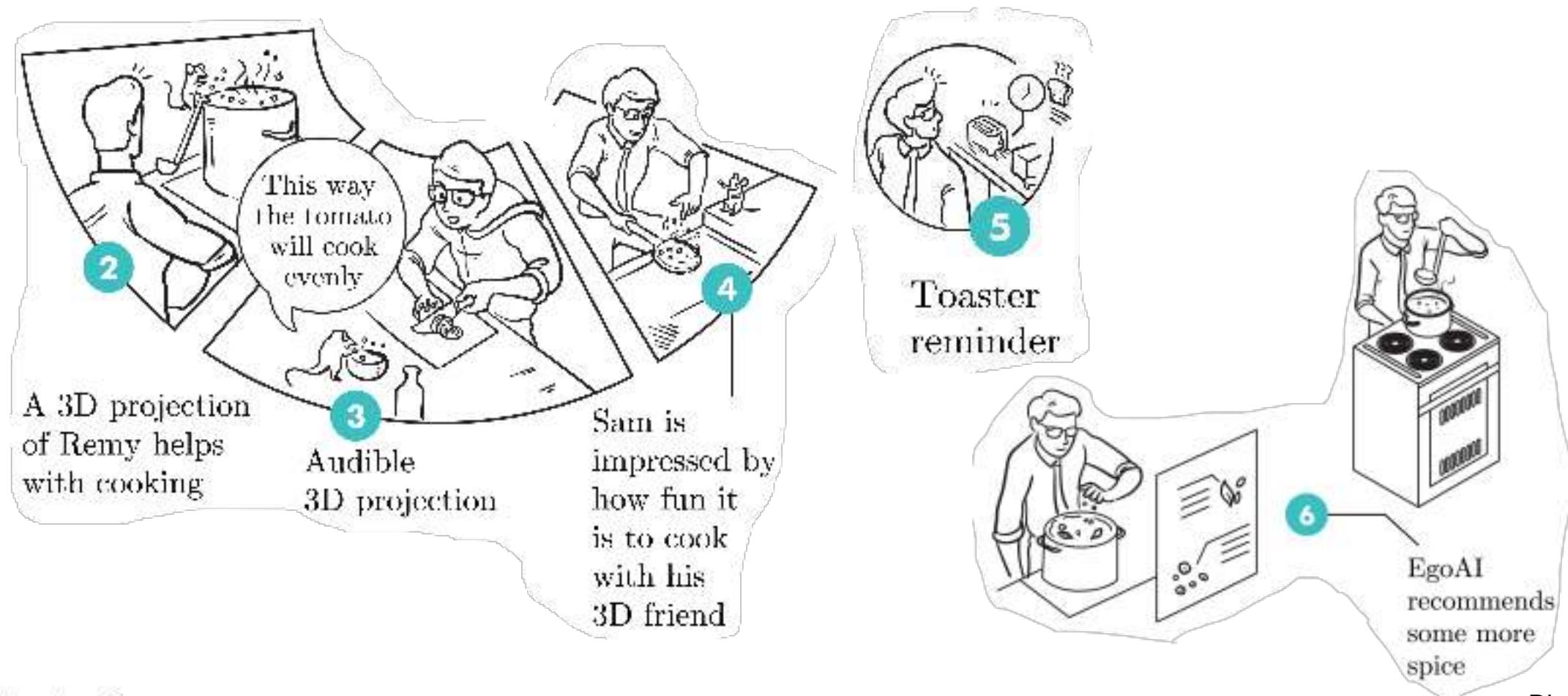
with: Chiara Plizzari, Gabriele Goletto, Antonio Furnari, Siddhant Bansal,  
Francesco Ragusa, Giovanni Maria Farinella, Tatiana Tommasi

Sam is finally home after a long day.  
EgoAI kept track of Sam's food intake and a tomato soup sounds like the best complementary nutrition



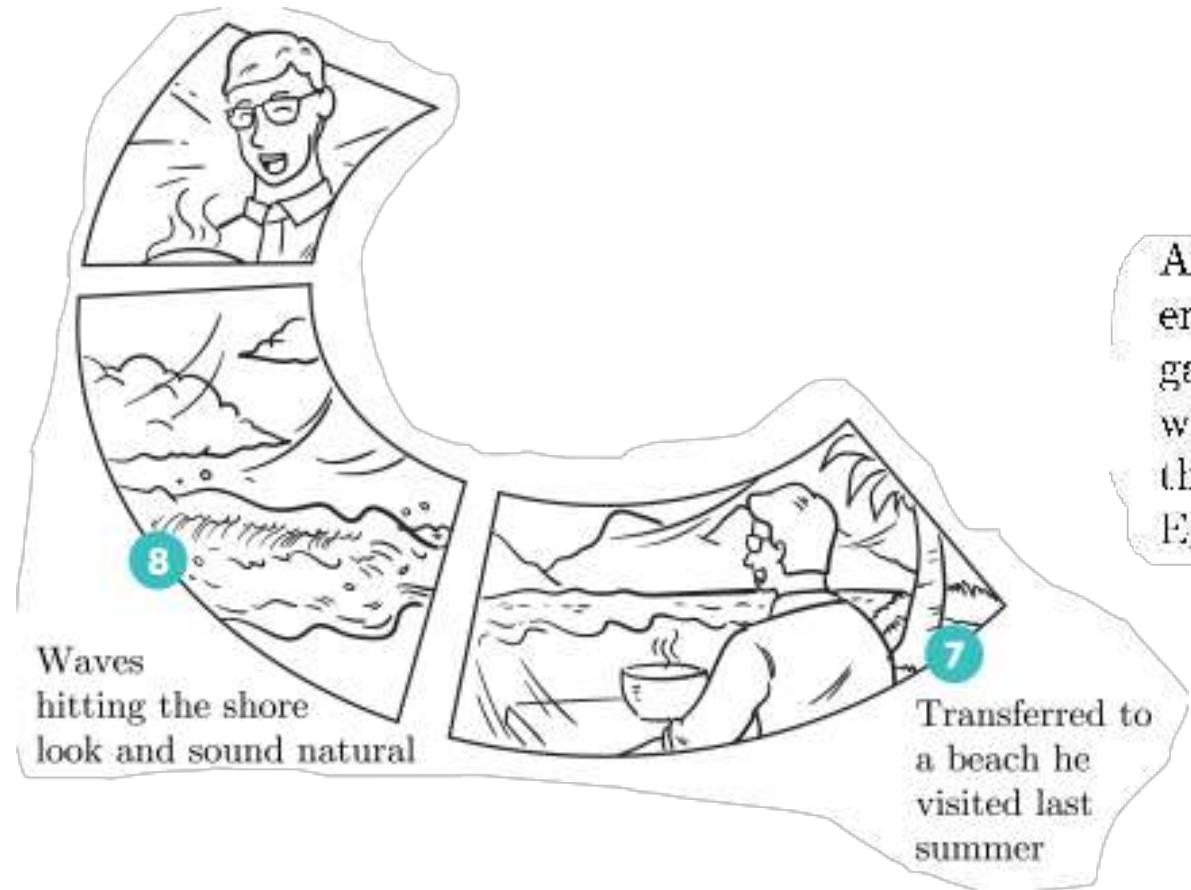
# EGO-Home

with: Chiara Plizzari, Gabriele Goletto, Antonio Furnari, Siddhant Bansal, Francesco Ragusa, Giovanni Maria Farinella, Tatiana Tommasi



# EGO-Home

with: Chiara Plizzari, Gabriele Goletto, Antonio Furnari, Siddhant Bansal, Francesco Ragusa, Giovanni Maria Farinella, Tatiana Tommasi

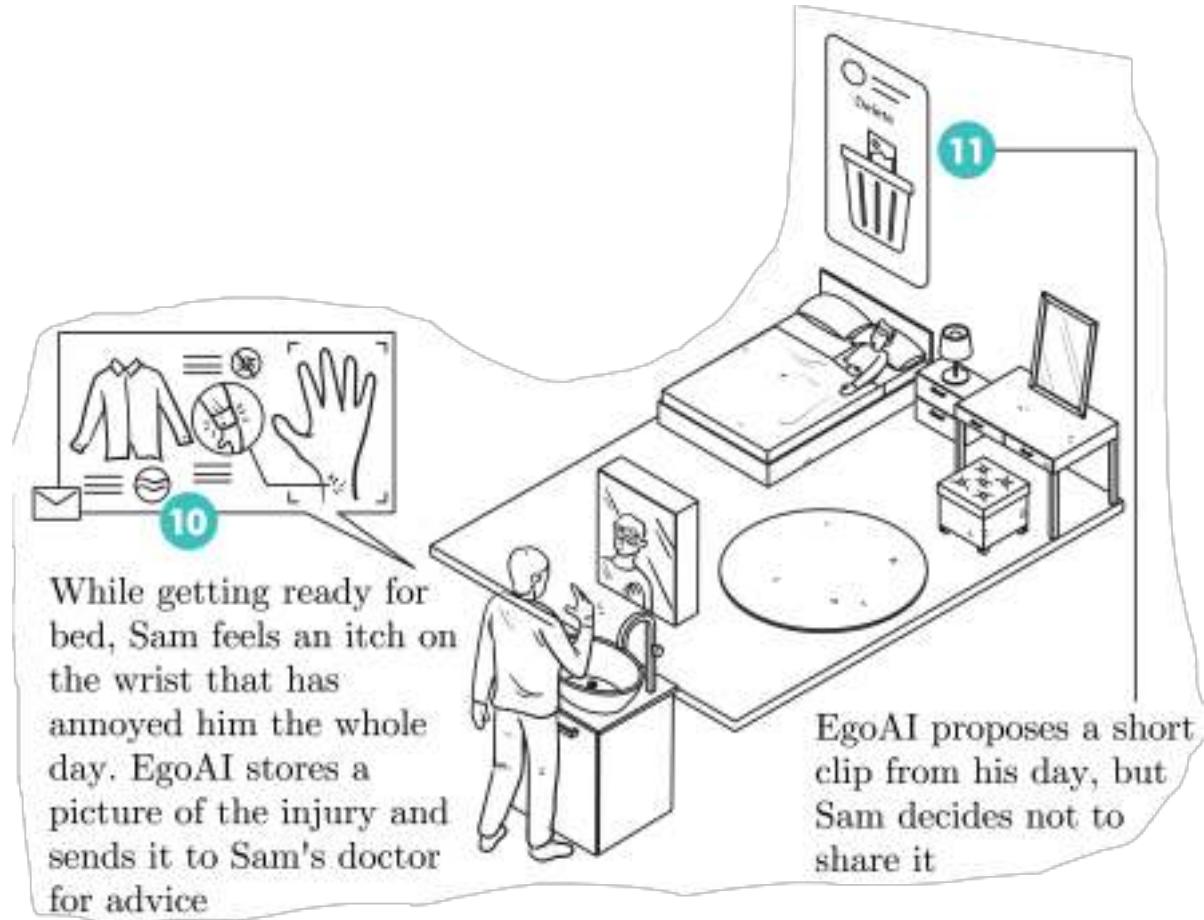


After dinner, Sam enjoys a group card game with his friends, who are connected through their own EgoAI



# EGO-Home

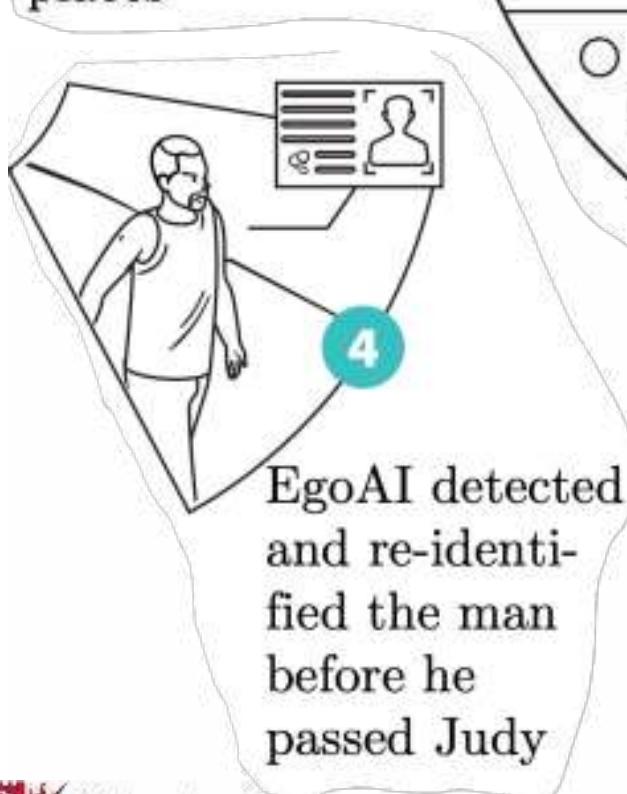
with: Chiara Plizzari, Gabriele Goletto, Antonio Furnari, Siddhant Bansal, Francesco Ragusa, Giovanni Maria Farinella, Tatiana Tommasi



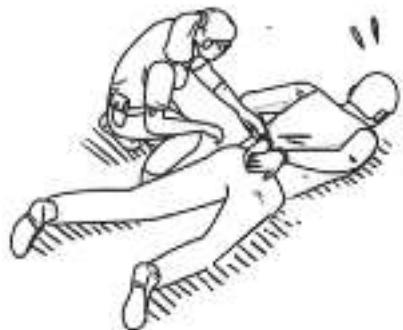
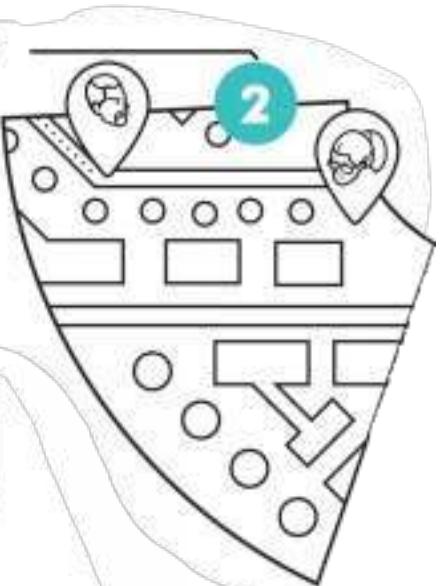
# From Stories to Tasks

with: Chiara Plizzari, Gabriele Goletto, Antonio Furnari, Siddhant Bansal, Francesco Ragusa, Giovanni Maria Farinella, Tatiana Tommasi

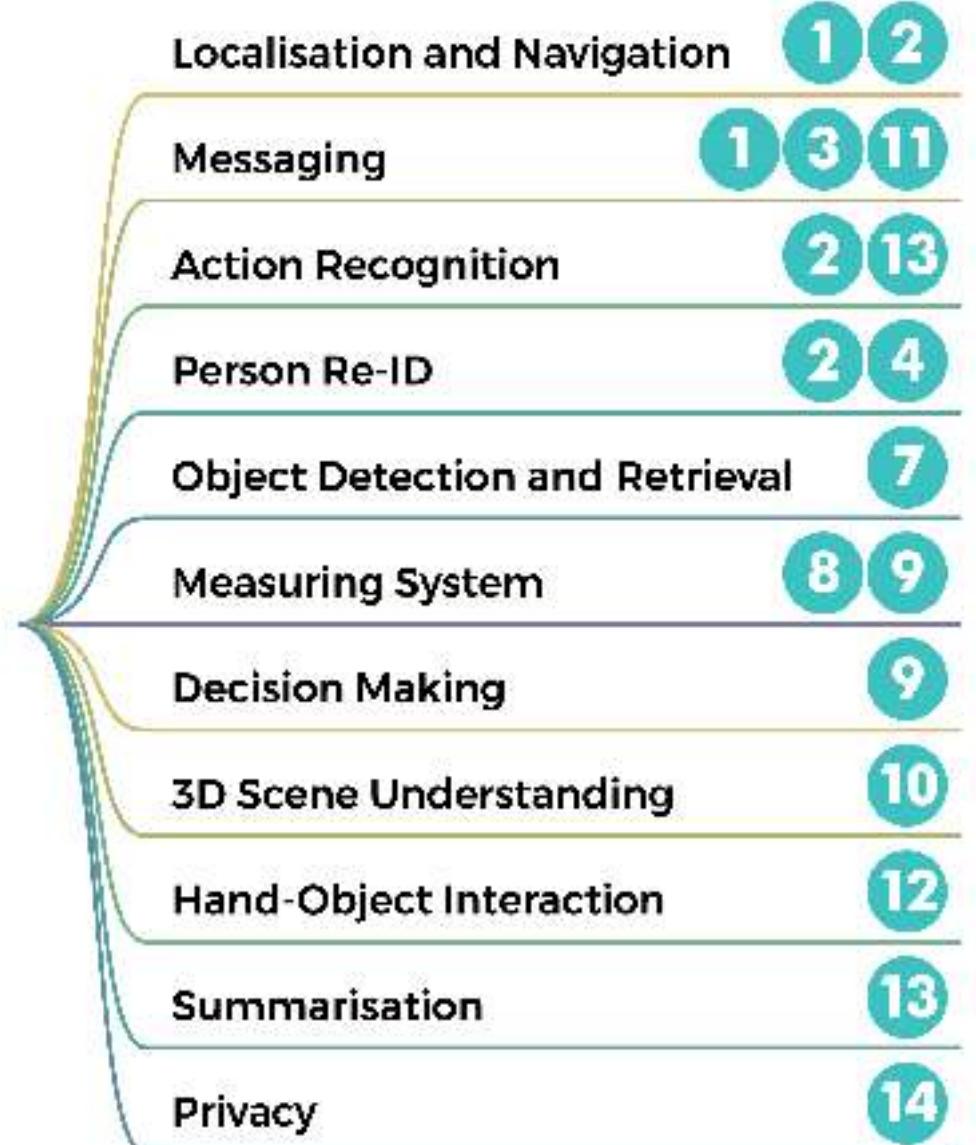
EgoAI helps Judy navigate through the shortest safe path to target places



EgoAI detected and re-identified the man before he passed Judy

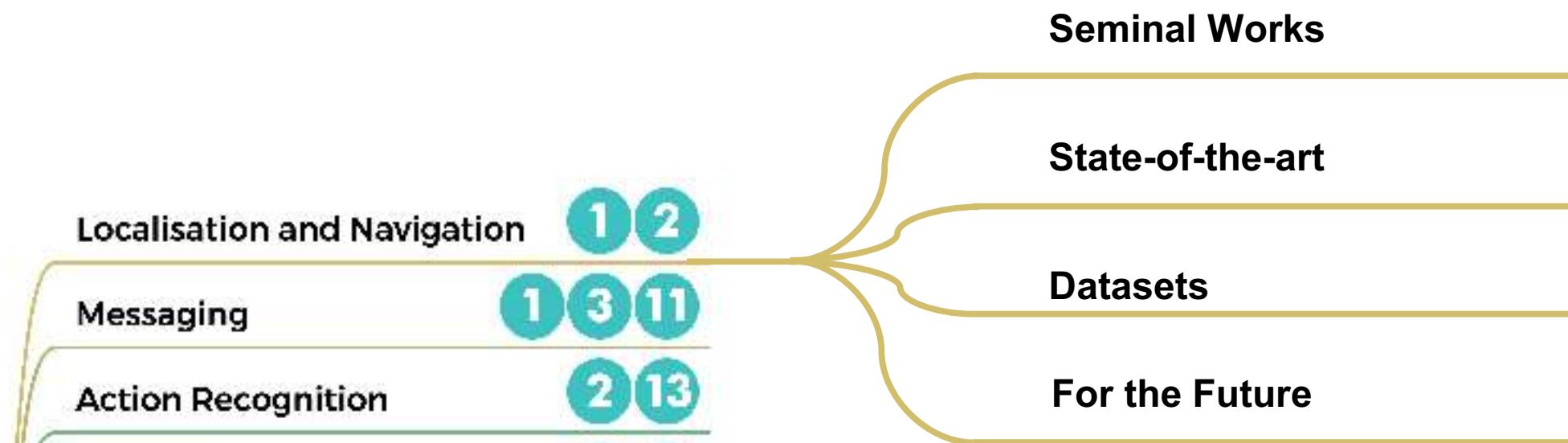


**EGO-Police**



# The Survey Part

with: Chiara Plizzari, Gabriele Goletto, Antonio Furnari, Siddhant Bansal,  
Francesco Ragusa, Giovanni Maria Farinella, Tatiana Tommasi



# The Survey Part

with: Chiara Plizzari, Gabriele Goletto, Antonio Furnari, Siddhant Bansal,  
Francesco Ragusa, Giovanni Maria Farinella, Tatiana Tommasi

- 12 tasks
- 46 pages (excluding references)
- 462 references

# In today's tutorial



Motivation and Datasets in  
Egocentric Video Understanding



Video Understanding  
Out of the Frame



Video Understanding:  
Data and Tasks



Teaser: The Wizard of Oz  
at the Sphere



Videos are Multimodal



Outlook into the Future of  
Egocentric Vision



Connected Videos of One's Life



Conclusion

# The Team

*grateful*



2017



2018



2019



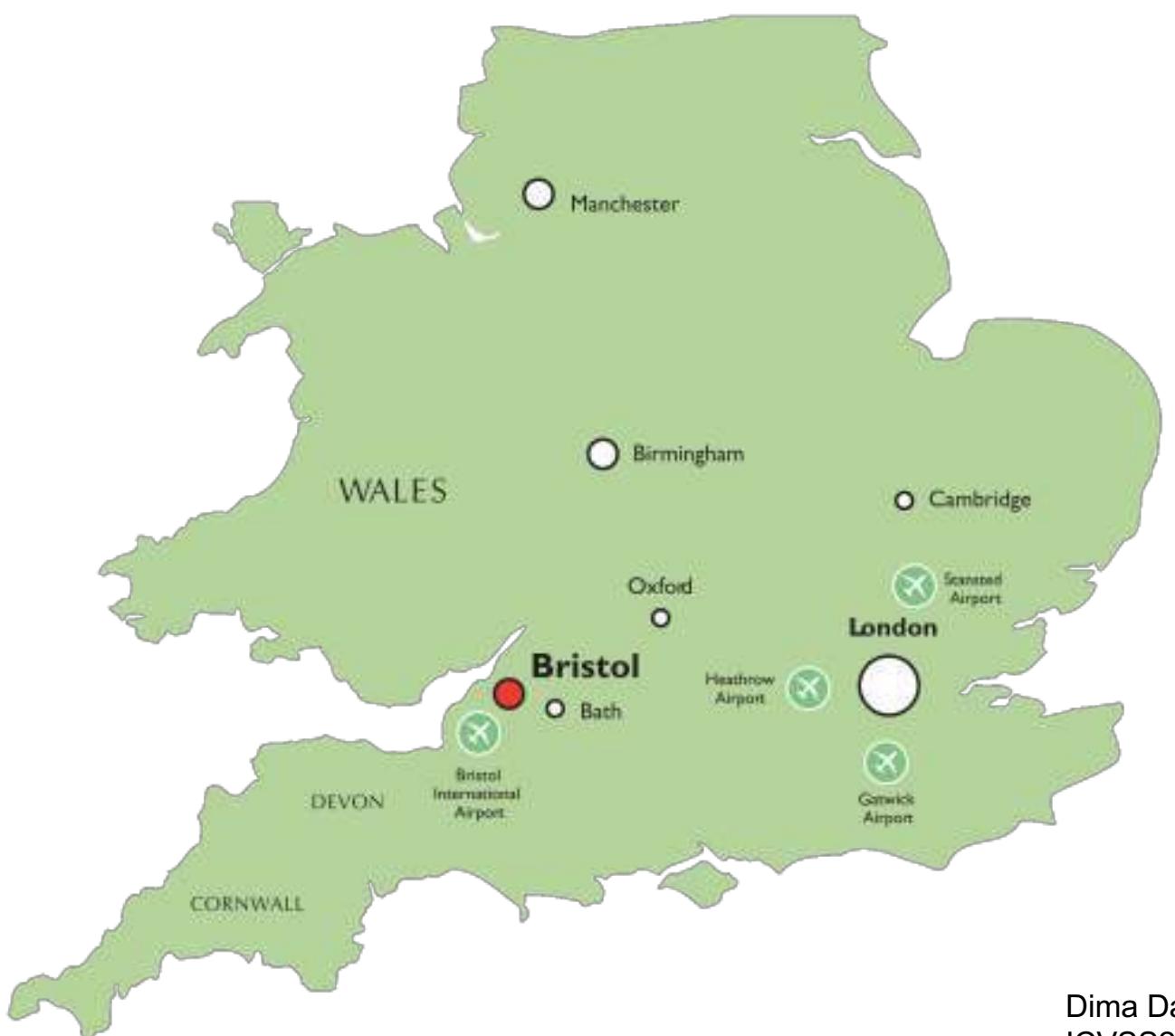
My research team...

*grateful*



Dima Damen  
ICVSS2025

# Visit us...



# Thank you

For further info, datasets, code, publications...

<http://dimadamen.github.io>



@dimadamen



@dimadamen.bsky.social



<http://www.linkedin.com/in/dimadamen>

## Q&A