



# Video Understanding

# Definitions...

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

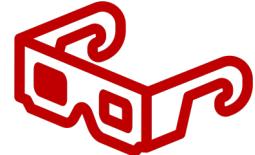


# Definitions...

**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

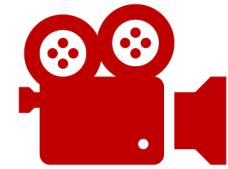




Motivation and Datasets in  
Egocentric Video Understanding



Video Understanding  
Out of the Frame



Video Understanding:  
Data and Tasks



Teaser: The Wizard of Oz  
& Genie 3



Videos are Multimodal



Outlook into the Future of  
Egocentric Vision



Connected Videos of One's Life



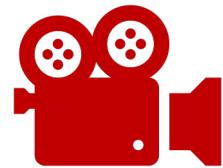
Conclusion



Motivation and Datasets in  
Egocentric Video Understanding



Video Understanding  
Out of the Frame



Video Understanding:  
Data and Tasks



Teaser: The Wizard of Oz  
& Genie 3



Videos are Multimodal



Outlook into the Future of  
Egocentric Vision



Connected Videos of One's Life

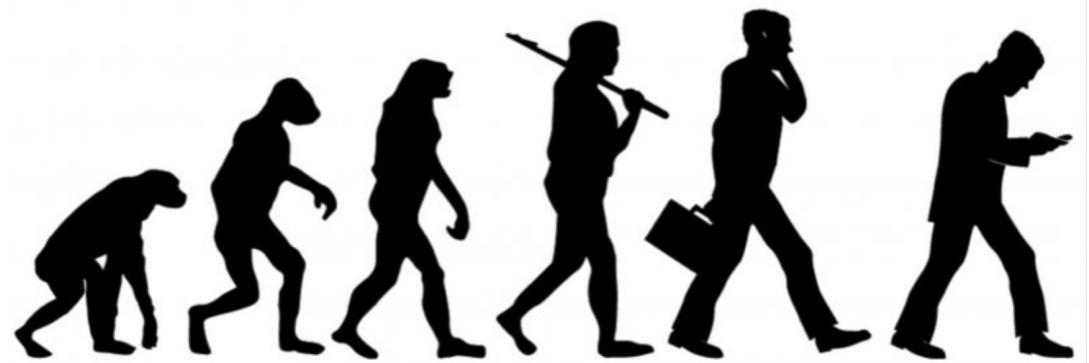


Conclusion

# The present...



Photo \*Illustration\* by Pelle Cass





# The future...

## HoloLens 2

A new vision for computing

[See pricing and options >](#) [Watch the HoloLens 2 video](#)

[FACEBOOK](#) Who We Are What We Build Our Actions Our Community Resources

**PROJECT ARIA GLASSES**

The goal of Project Aria is learning 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 when and where to use the device. We'll be asking people of diverse backgrounds to participate in the program to create an accurate and varied view of the world.

Project Aria glasses are not a consumer product, nor are they an AR glasses prototype. The glasses do not include a display and research participants cannot directly view video or listen to audio captured by the device, but participants can view low-resolution thumbnails via a companion app installed on their phone for the purpose of deleting segments of data. We'll use encryption to store the data on the Aria device and a secure ingestion system to upload data from the research devices to Facebook's separate, designated back-end storage space.

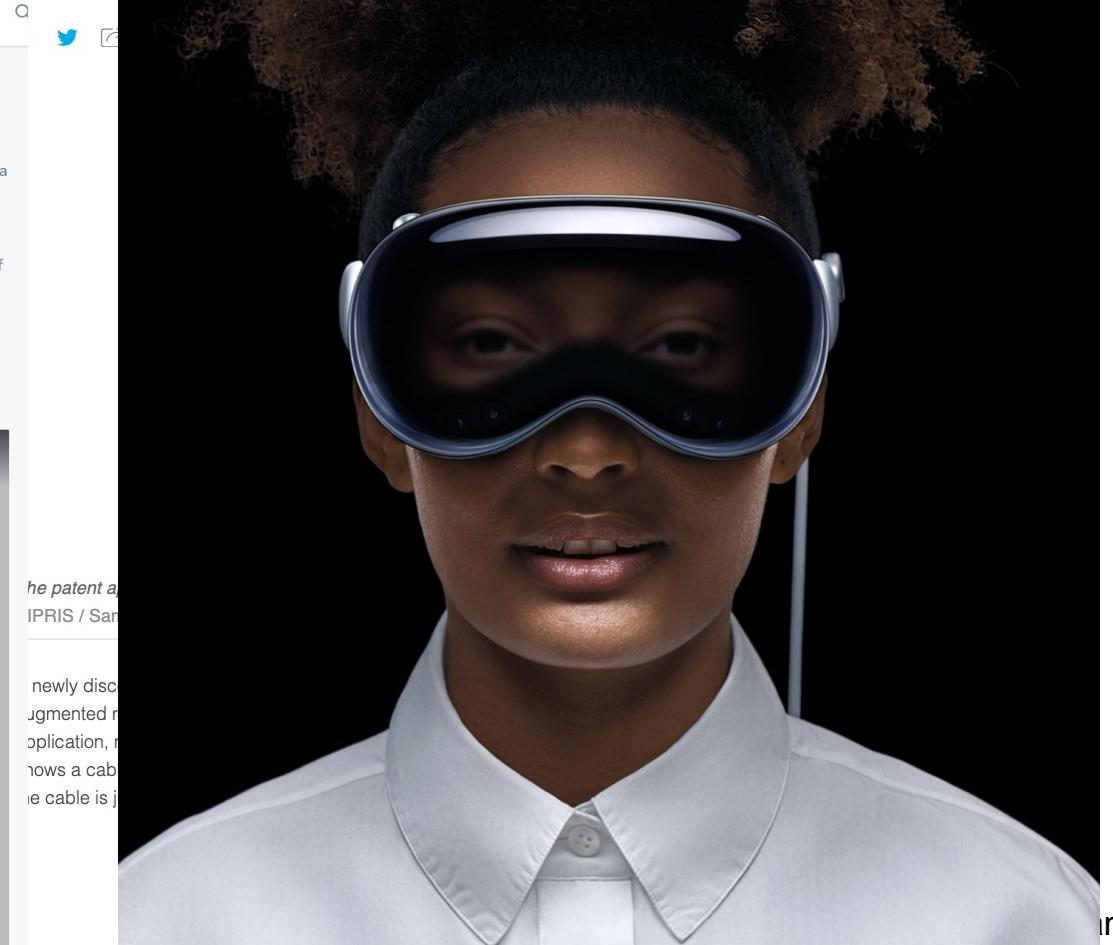
 **Facebook Reality Labs**  
Tech@Facebook • Follow



## Samsung patent application reveals augmented reality headset design

*It comes as the Gear VR slowly fades away*

By Jon Porte | Oct. 1, 2019 8:11 a.m. EDT



# The future is here...

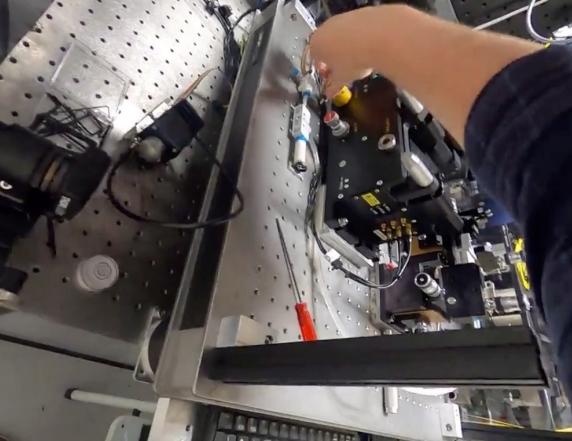


Dima Damen  
PAISS 2025

# The future can be imagined...



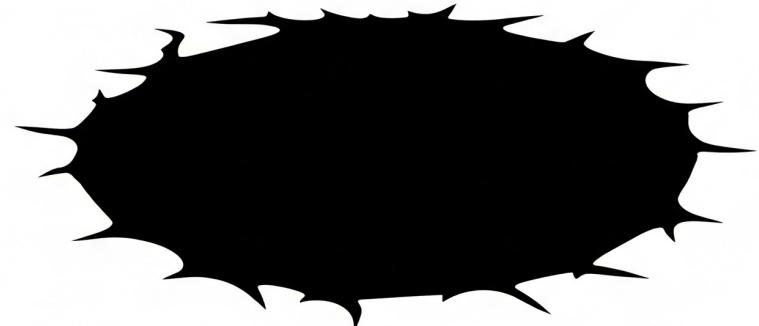
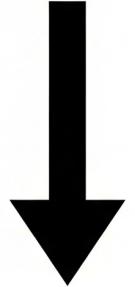
# Egocentric Videos?

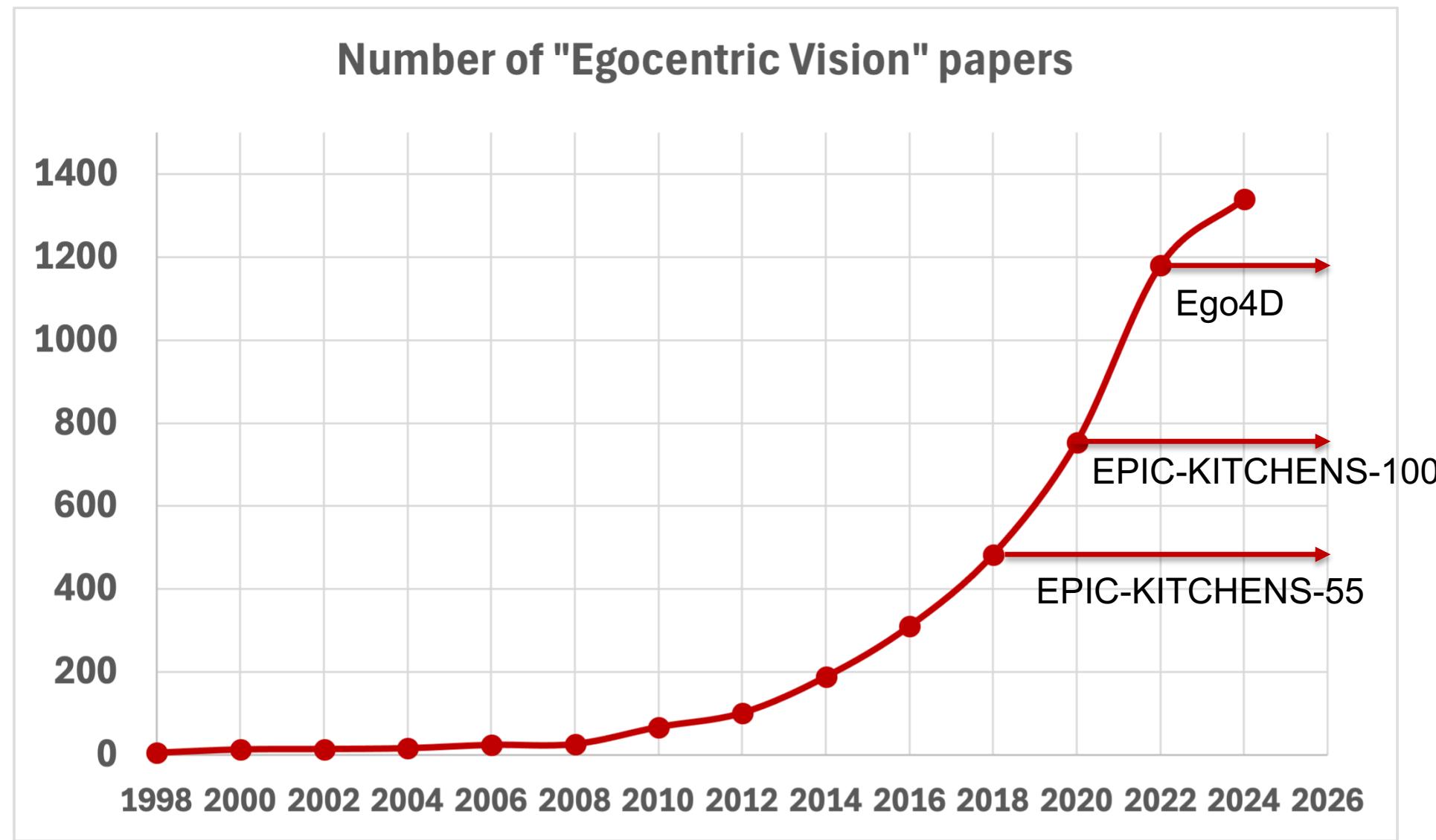


# Machine Learning in Practice

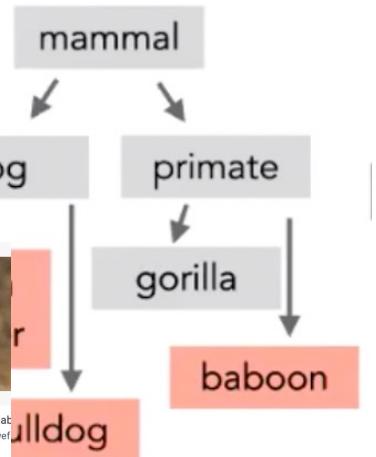
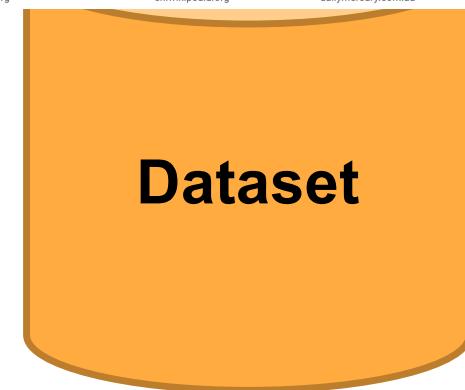
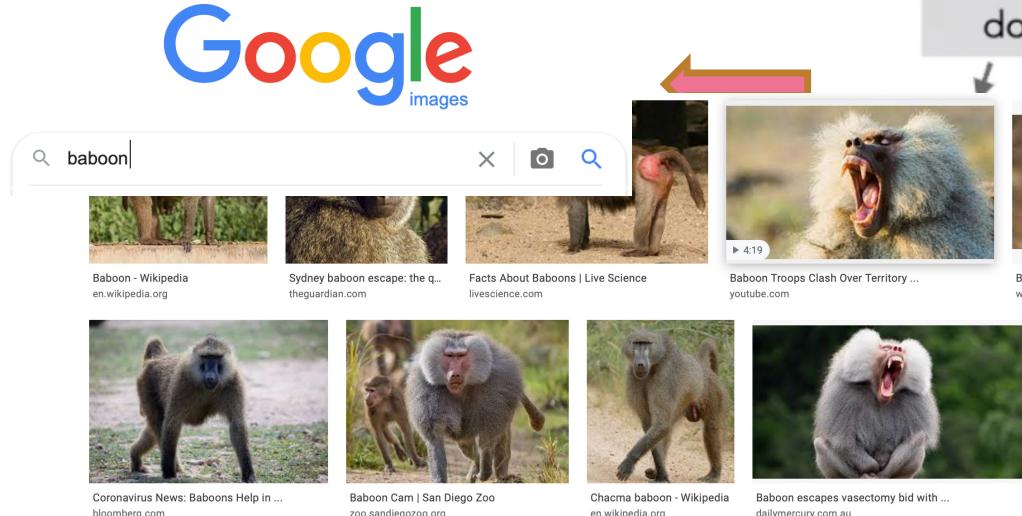
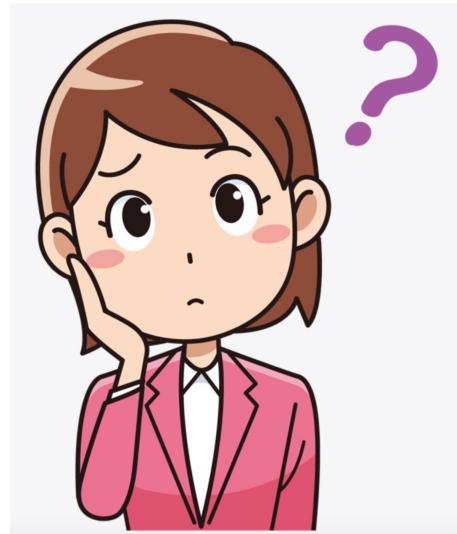


**no data  
no machine  
learning research**

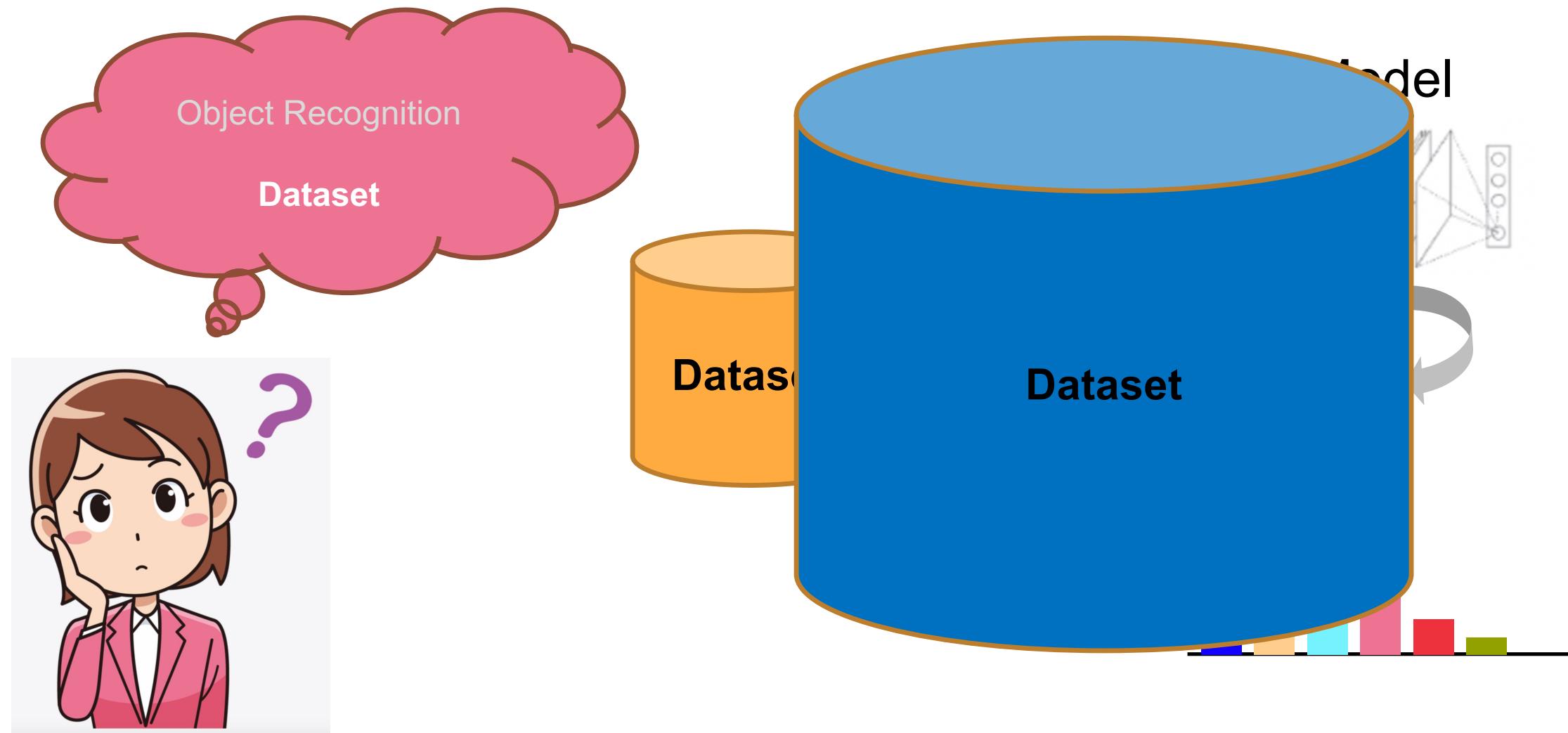




# Machine Learning in Practice



# 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
- One Exception**

# 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.

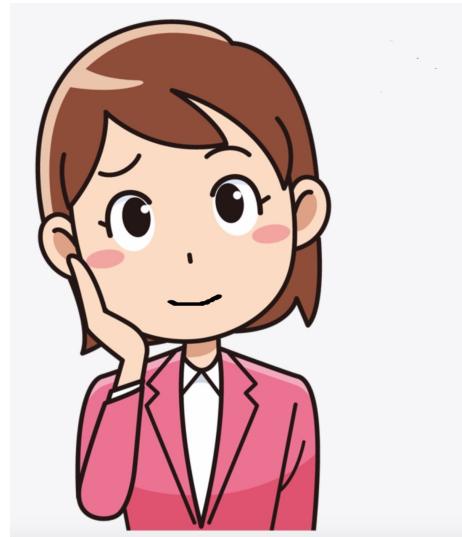
 Share

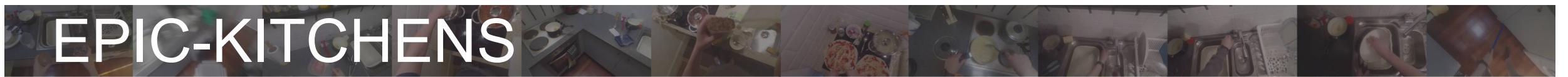


To get started, grab a cup of your favorite beverage and watch our video trailer (5 minutes):

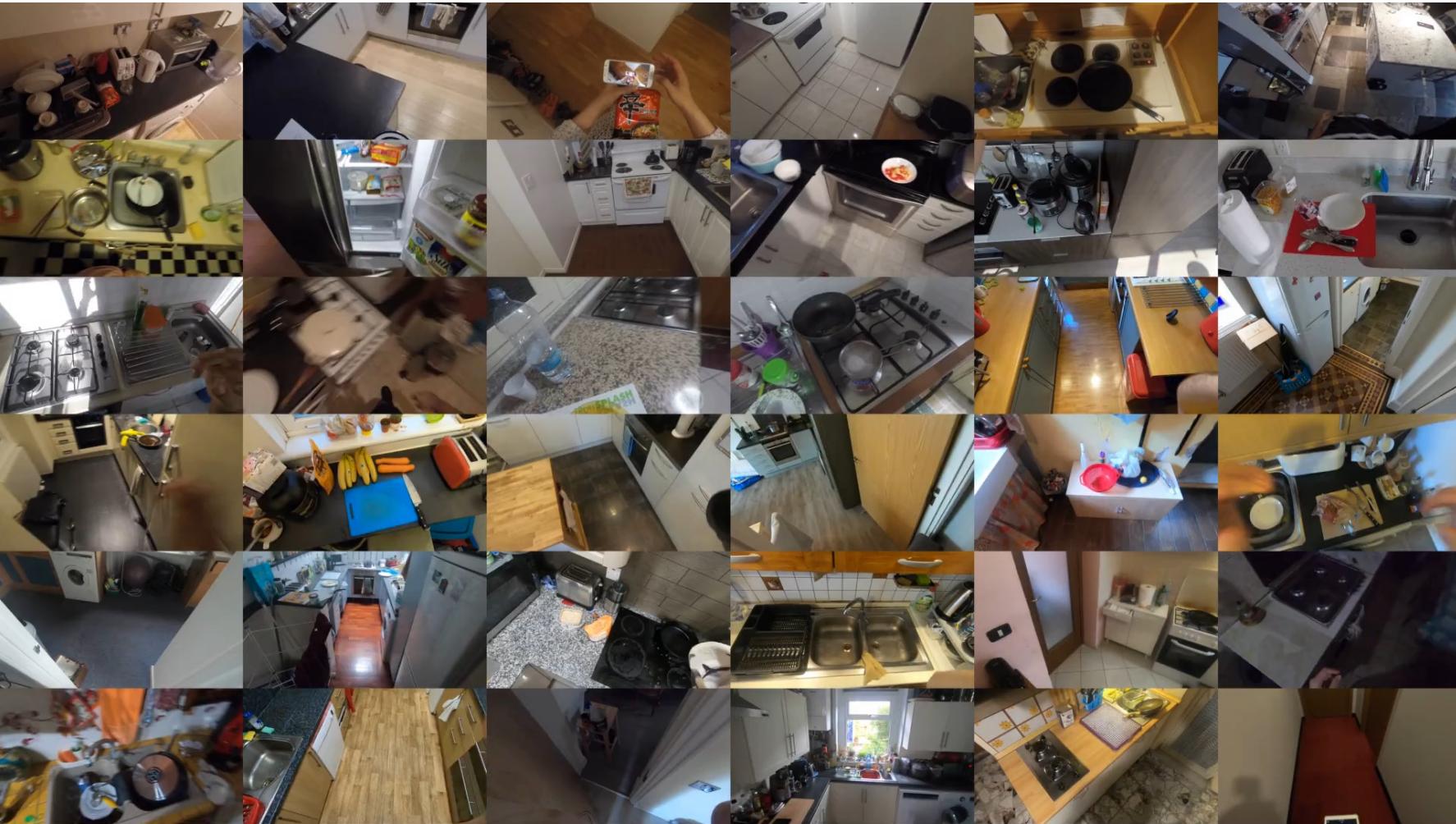
stereo flow sceneflow depth odometry object tracking road semantics raw data

# Machine Learning in Practice

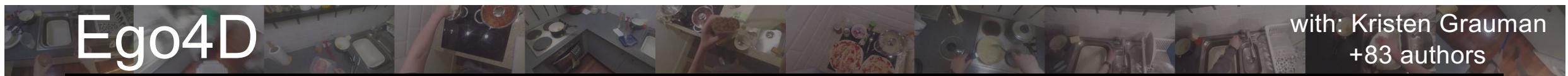




# EPIC-KITCHENS







## Narration

C: camera wearer

13.2 sentences/min  
3.8 M sentences

1,772 verbs

put pick move  
take hold drop  
hit pour lift  
clean stretch

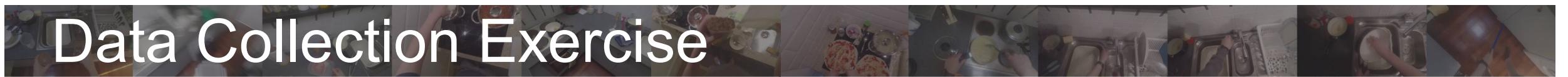
4,336 nouns

card cloth  
bottle hand  
paper wood  
container phone tray screw cover

#C C scraps off wood filler from one putty knife with the other putty knife  
#C C picks up another putty knife from the white board



# Data Collection Exercise



2017 - now

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



2020 - now

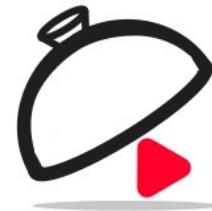
13 universities  
3670 hours  
923 participants  
74 locations  
9 countries  
Short-term recording  
All daily activities



**EGO-EXO4D**

2024 - now

1286 hours  
740 camera wearers  
Skilled activities

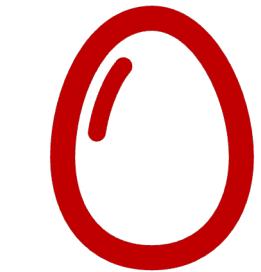


**HD-EPIC**

2025 - now

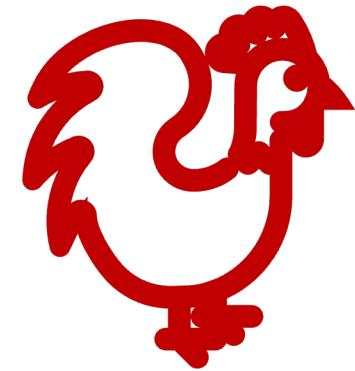
Validation dataset  
41 hours  
9 participants  
Highly-Detailed  
Digital Twin

# Data Collection Exercise

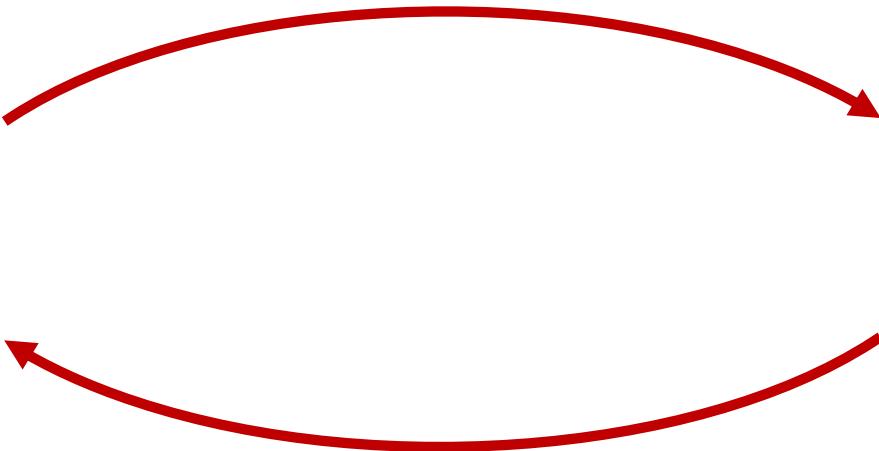


**Labels**

Pascal VOC  
ImageNet  
Kinetics  
Something-Something



**Data**

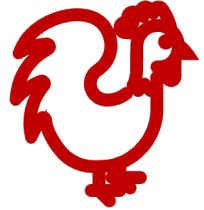


EPIC-KITCHENS  
Ego4D  
Ego-Exo4D  
HD-EPIC  
...  
KITTI

# The chicken or the egg...



## Data



Naturally unbalanced

Harder to label (exposes ambiguity)

Closer to application

Multiple tasks

## Labels

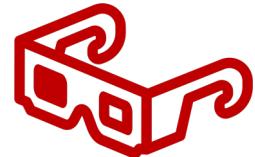


Unnaturally balanced (or nearly)

Easier to label (hides ambiguity)

Can be expanded

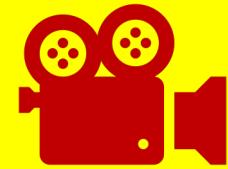
Single task



Motivation and Datasets in  
Egocentric Video Understanding



Video Understanding  
Out of the Frame



Video Understanding:  
Data and Tasks



Teaser: The Wizard of Oz  
& Genie 3



Videos are Multimodal



Outlook into the Future of  
Egocentric Vision



Connected Videos of One's Life



Conclusion

# The history of VIDEO

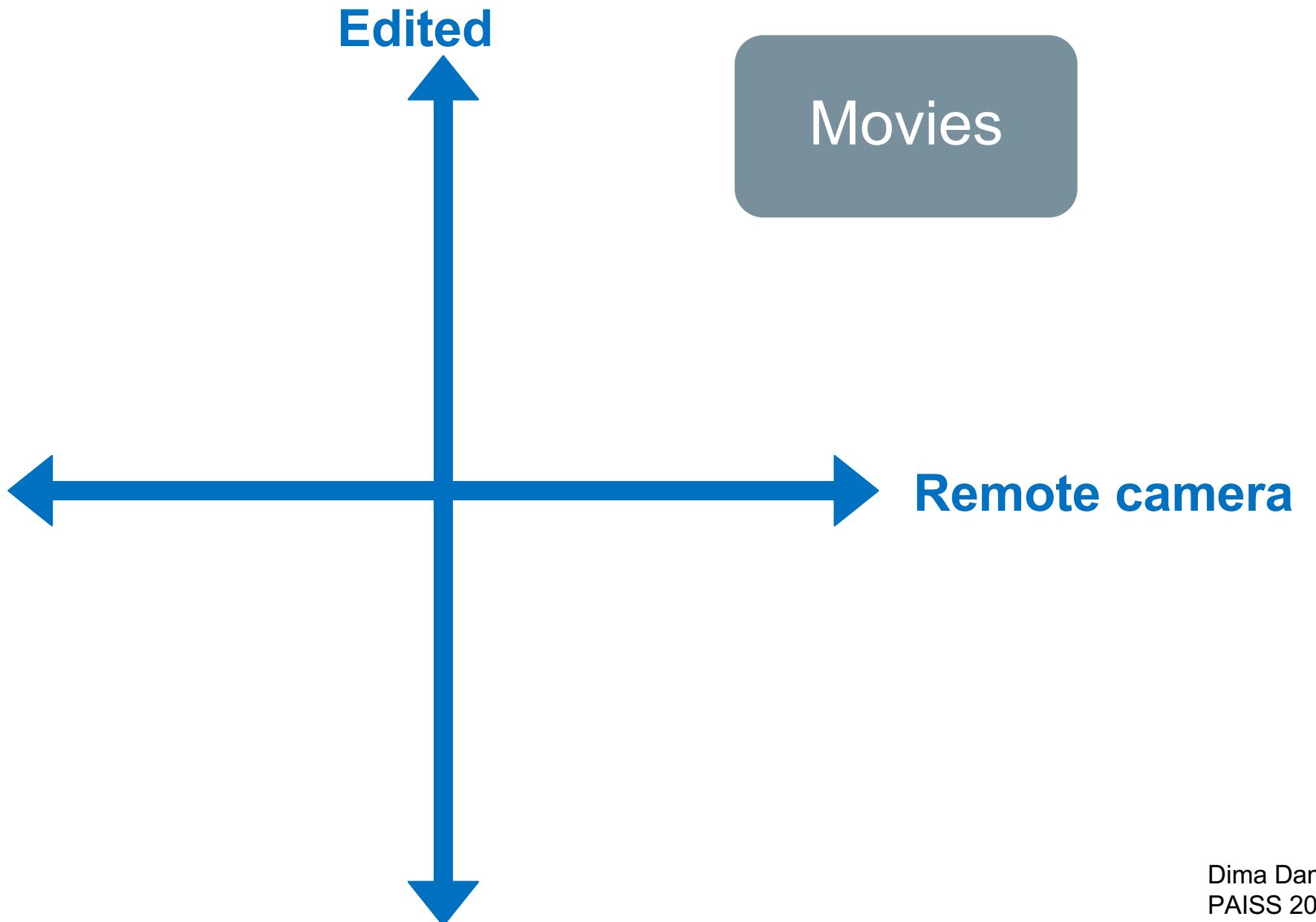


# The history of VIDEO





# 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-

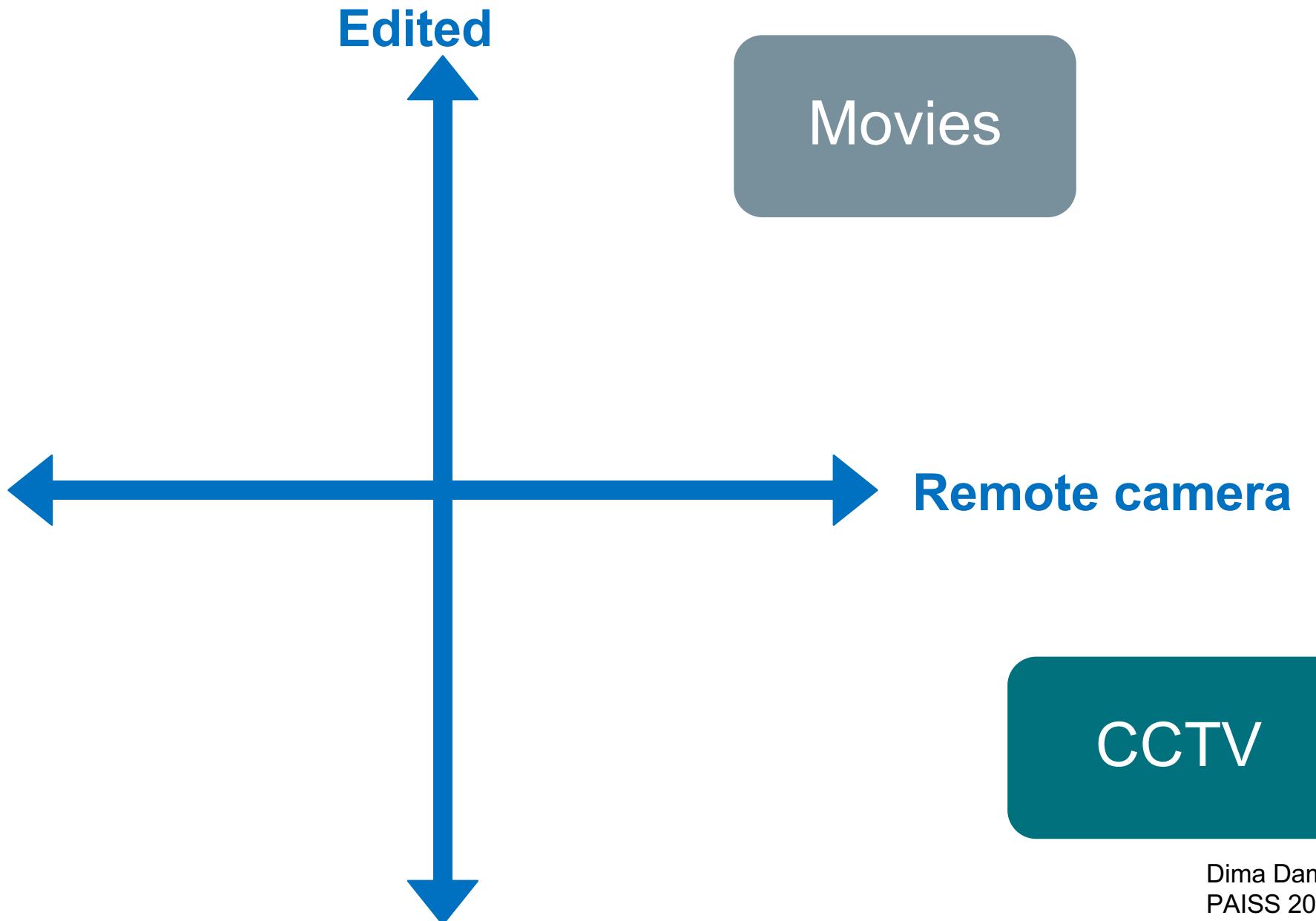
Laptev and Perez (2007)

# The history of VIDEO understanding





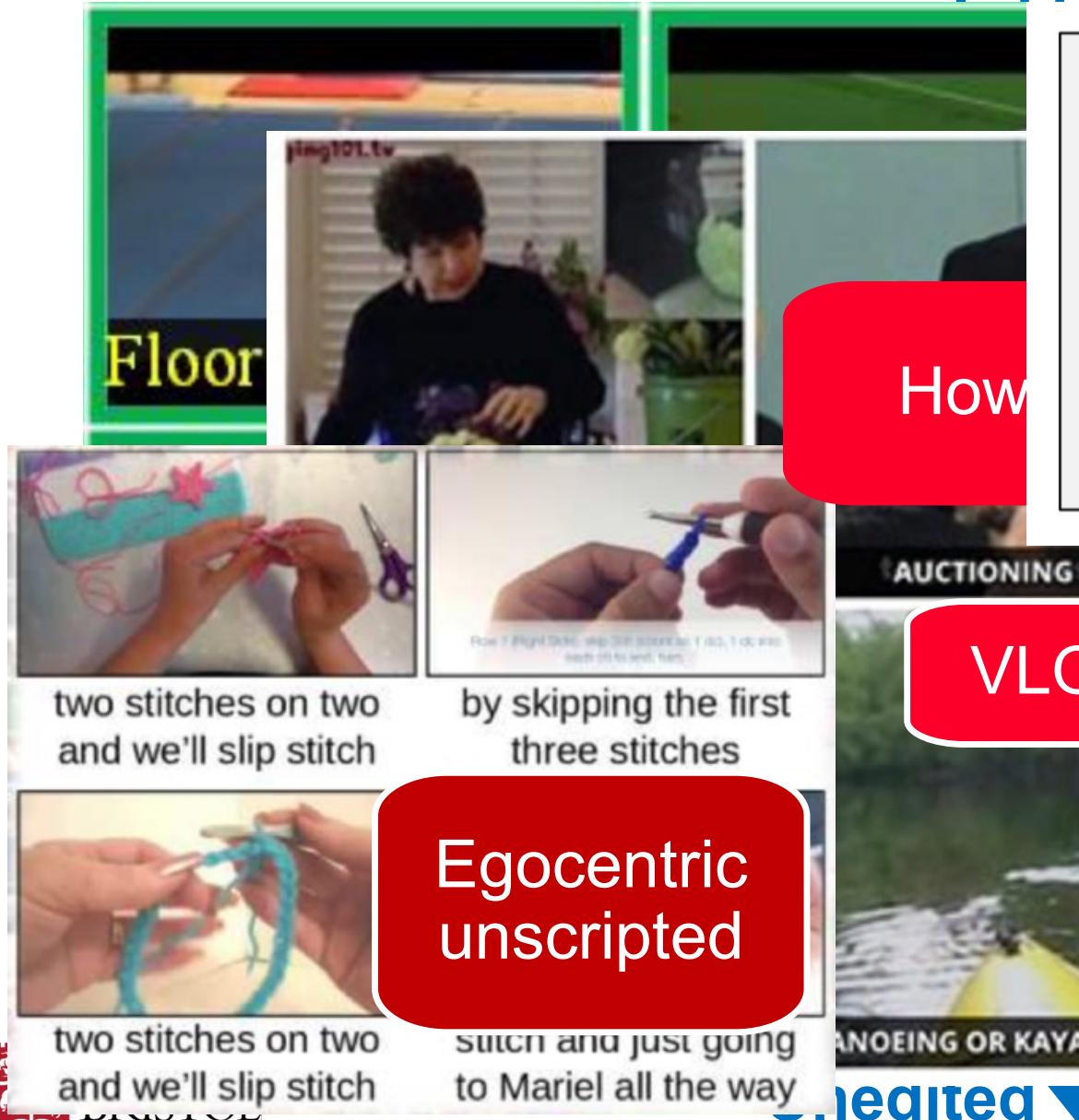
# The history of **VIDEO** understanding



# The history of **VIDEO** understanding



# The history of **VIDEO** understanding



**Templated,  
Multilingual Domain  
Queries:**

"Morning routine",  
"realistic ditl 2015",  
"mijn realistische  
routine", "Ma routine  
d'apres-midi", ...

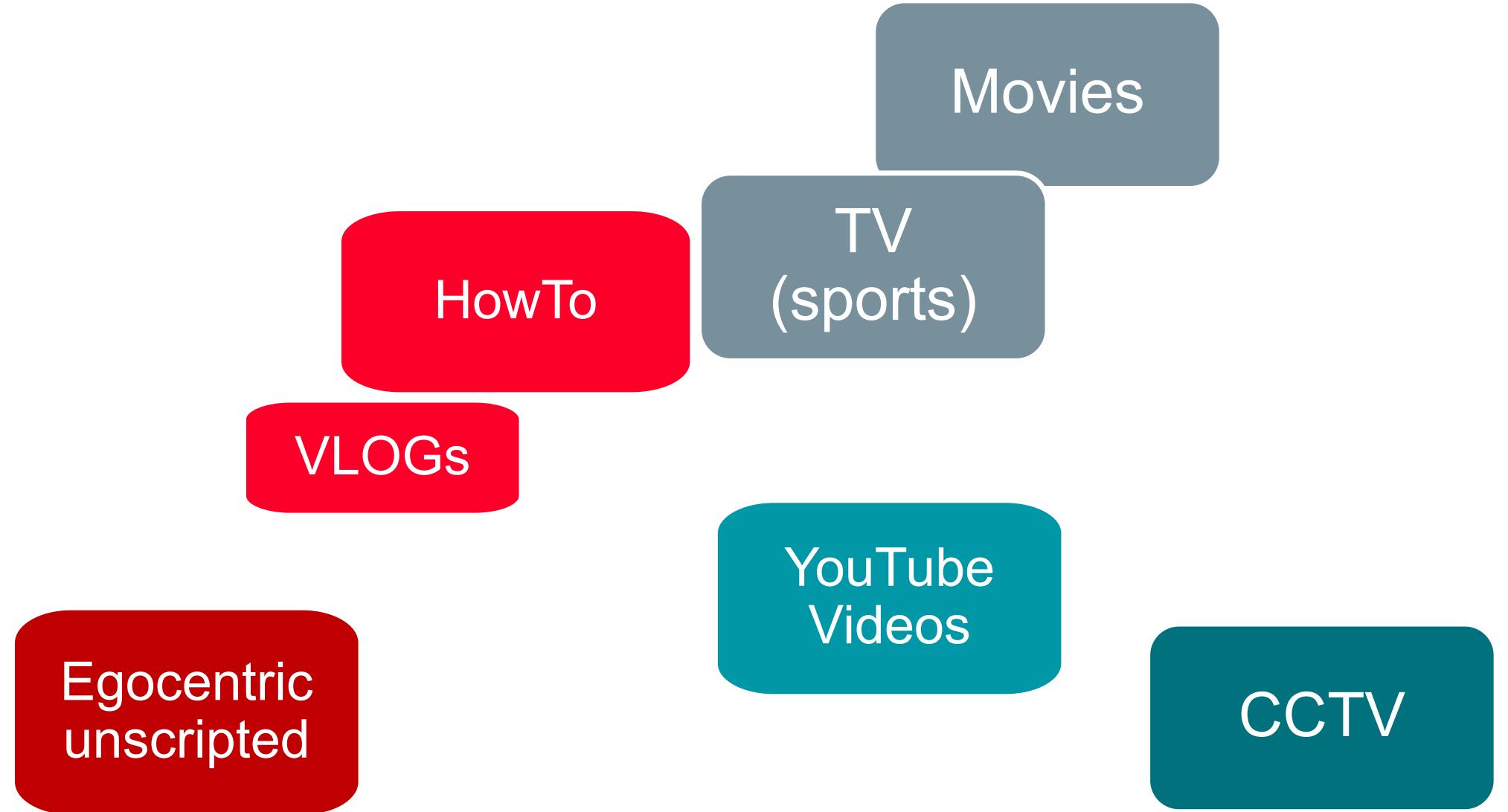
216K Video Candidates (2.5 Years)  
Low *Video-level Purity*

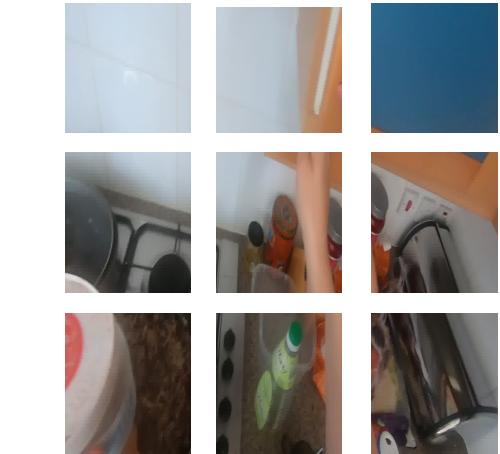


**Remote camera**

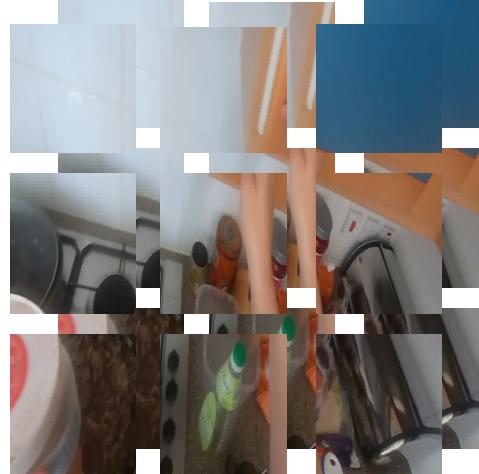
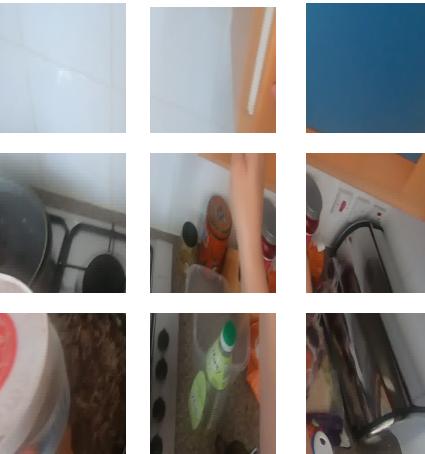
**CCTV**

# The history of **VIDEO** understanding

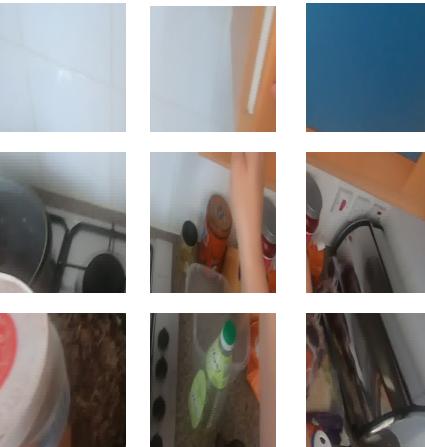




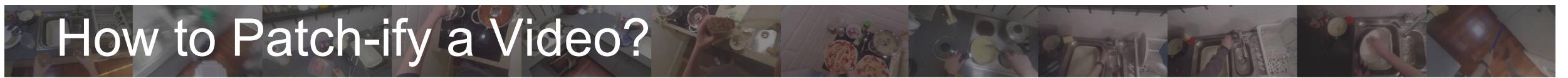
$H \times W \times C$



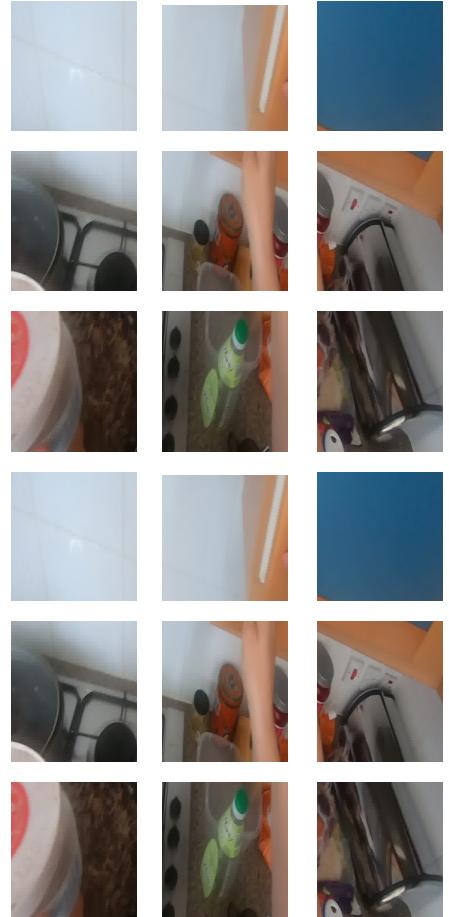
$H \times W \times [CT]$



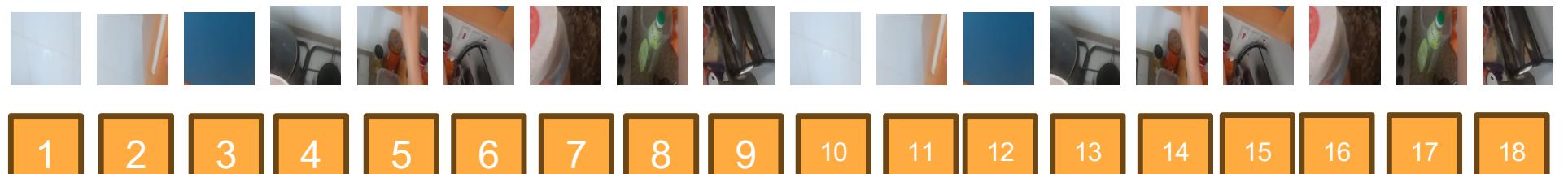
$[TH] \times W \times C$



# How to Patch-ify a Video?

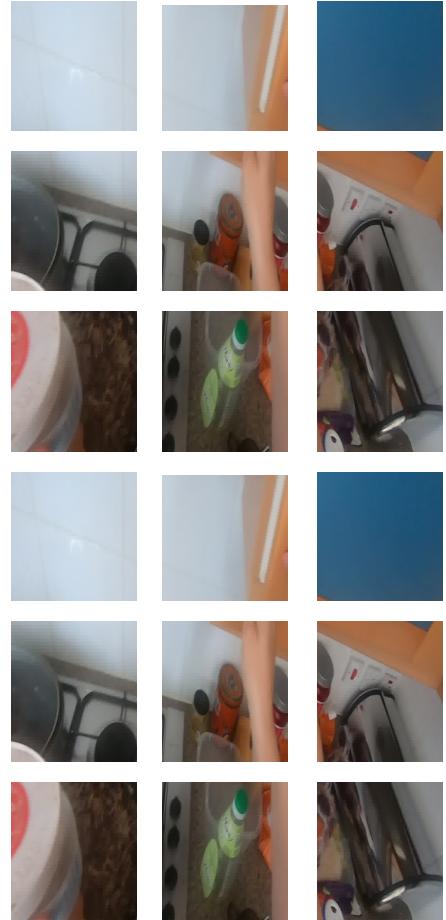


Flatten

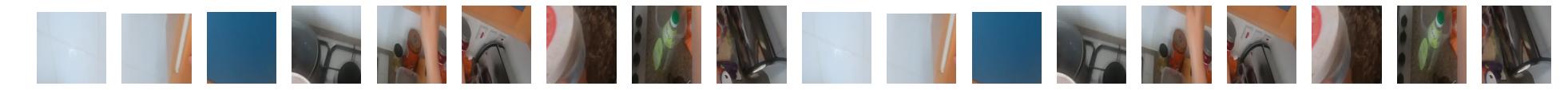


1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18

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

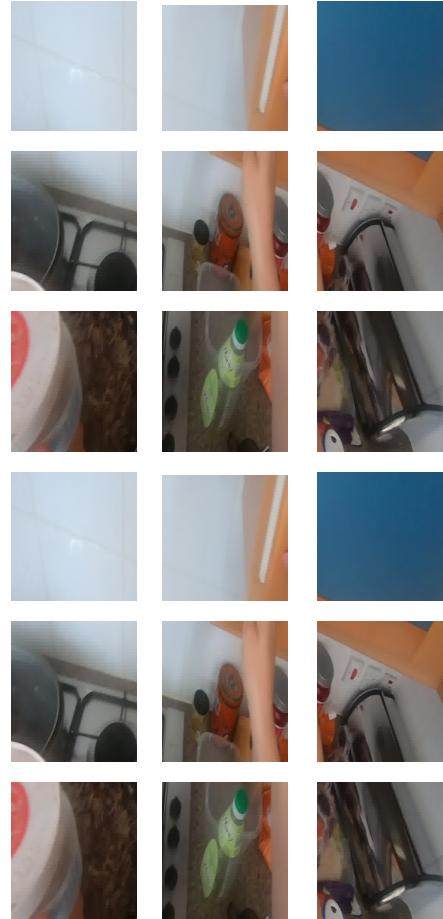
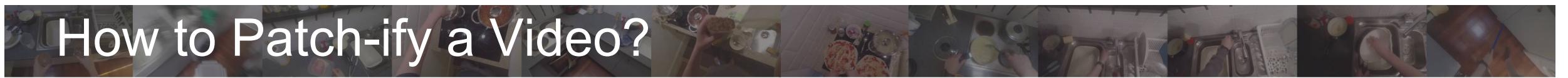


Flatten



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

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

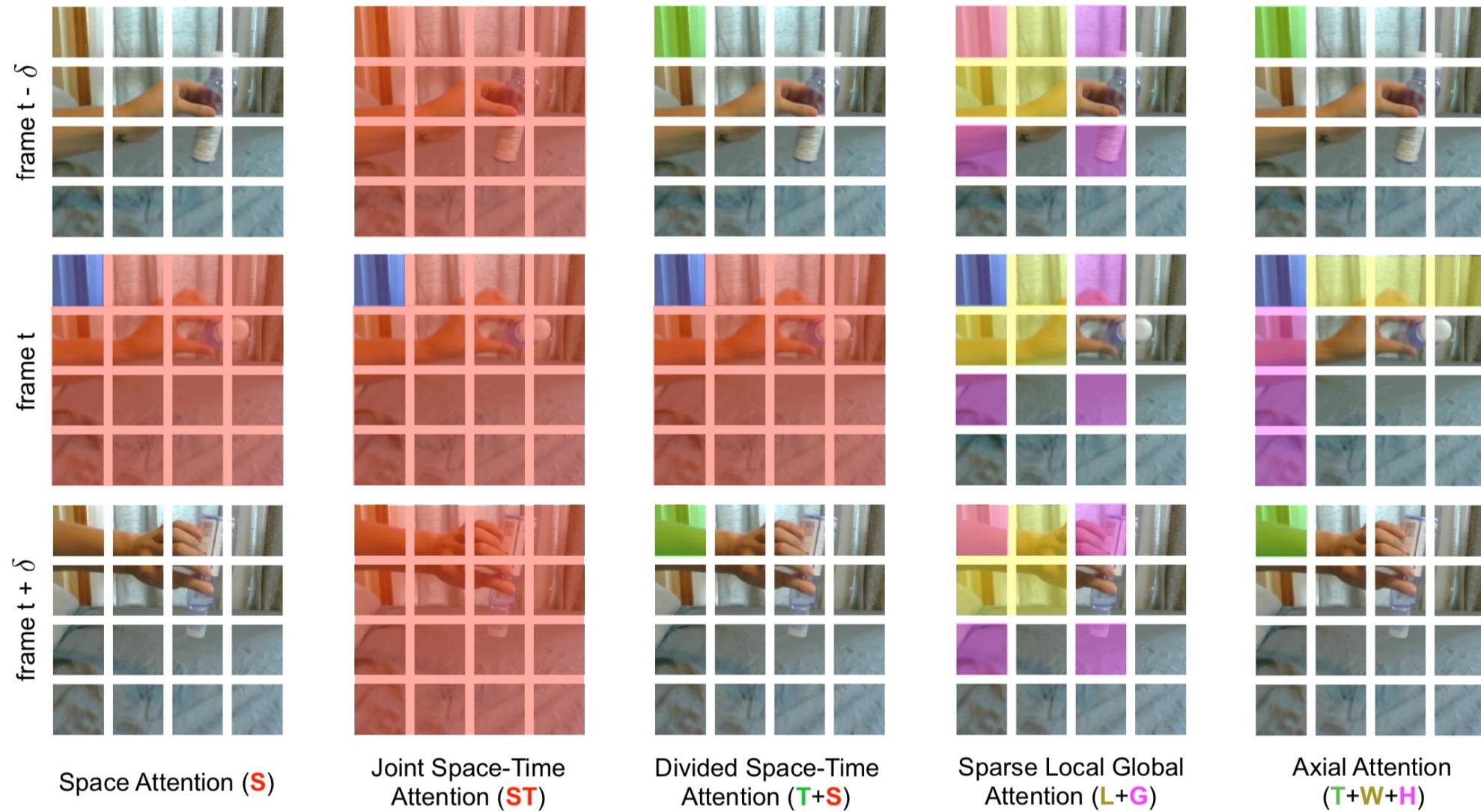


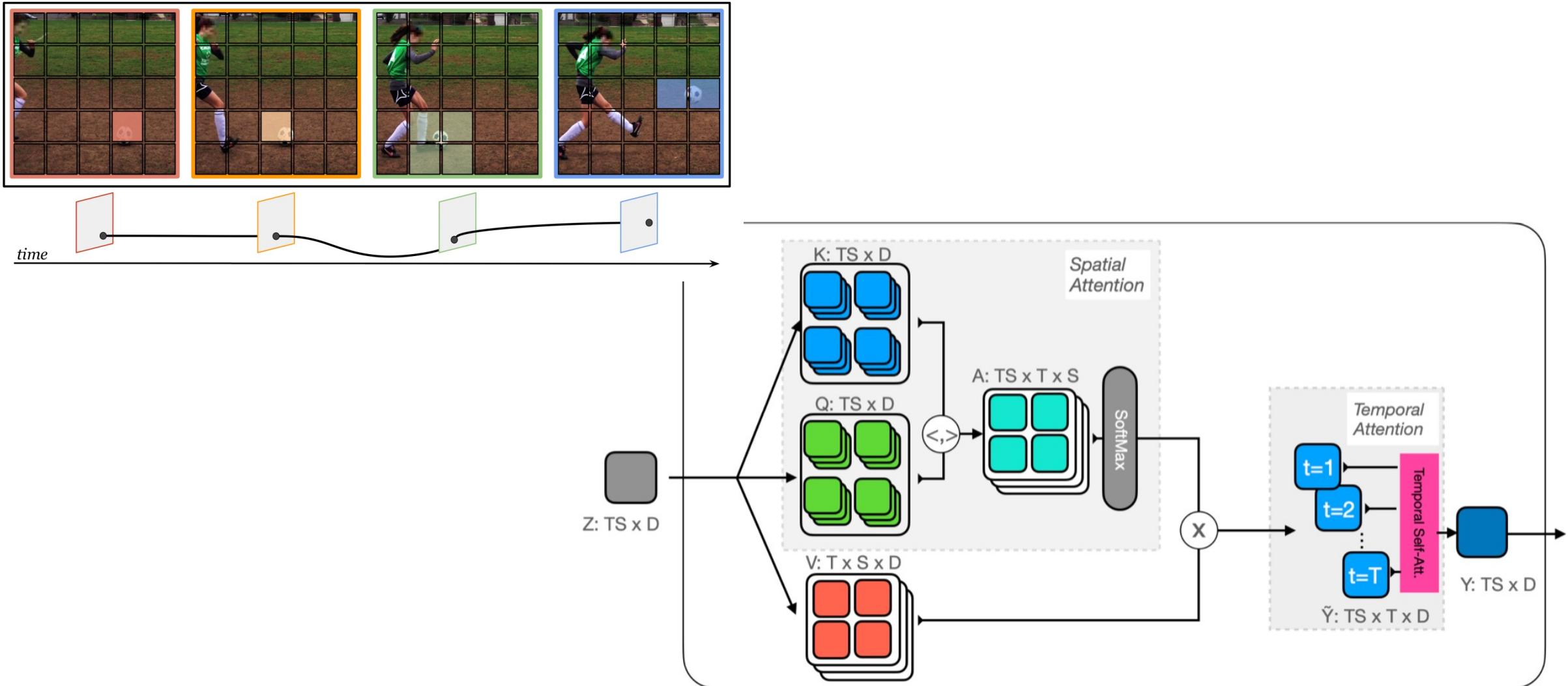
Transformer

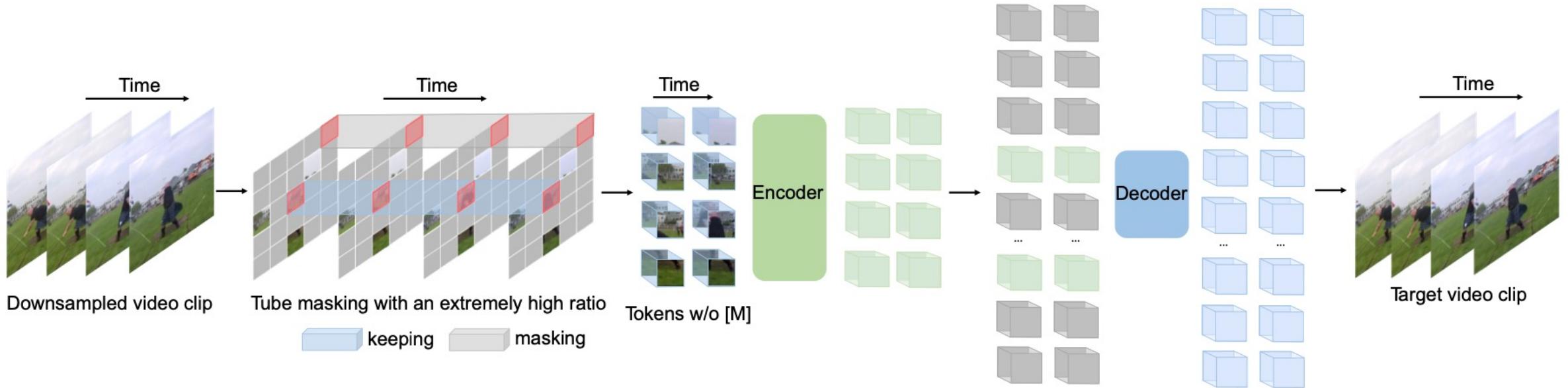


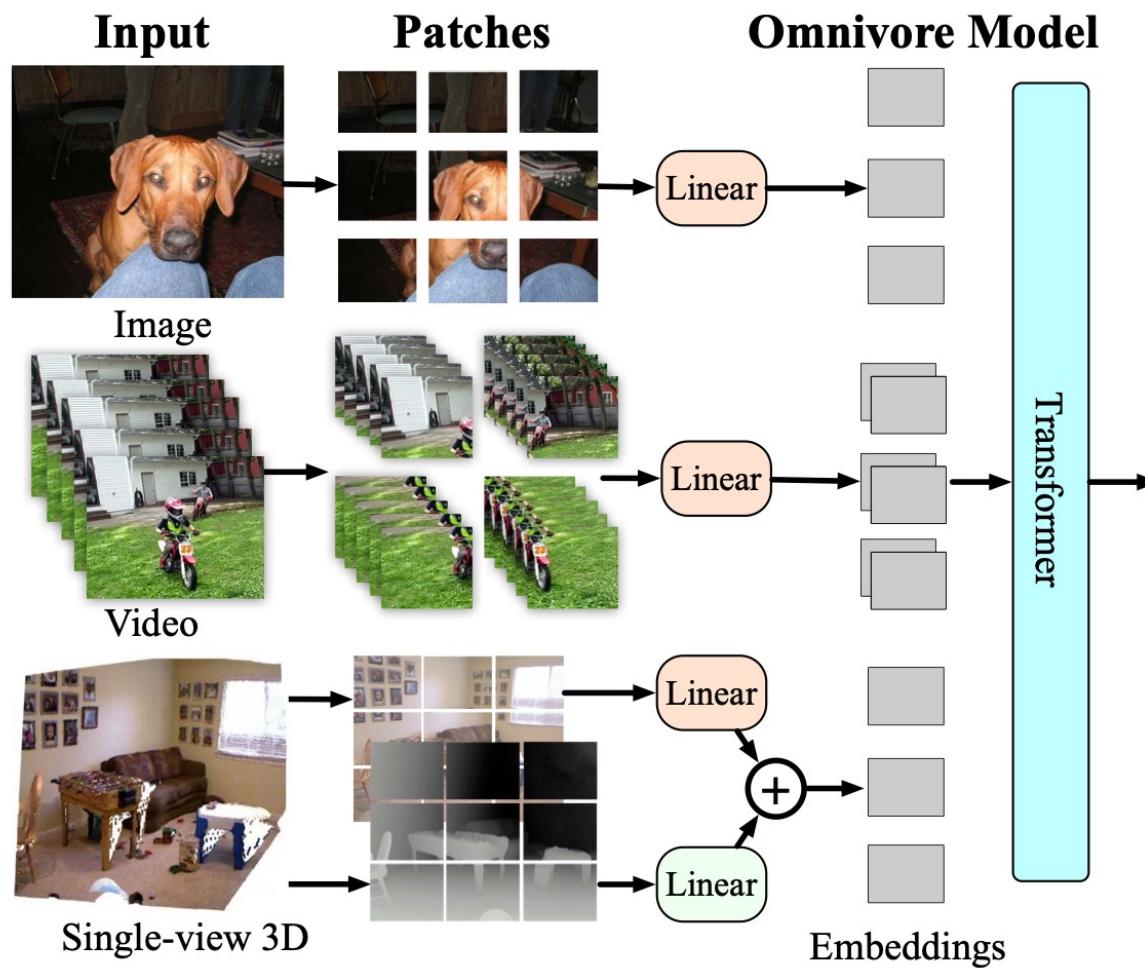
S	1	2	3	4	5	6	7	8	9	1	2	3	4	5	6	7	8	9
T	1	1	1	1	1	1	1	1	1	2	2	2	2	2	2	2	2	2

[TH] x W x C









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



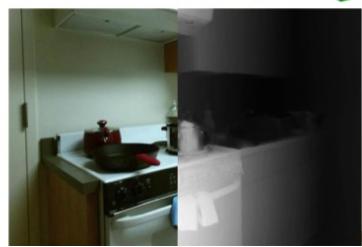
# ImageBind



Web Image-Text



Depth Sensor Data



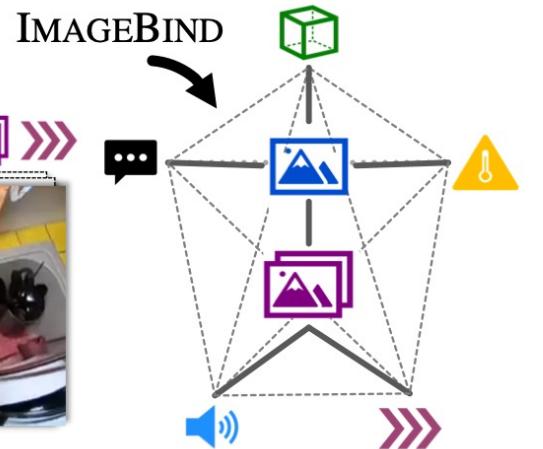
Web Videos



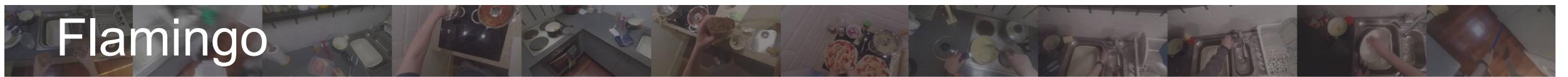
Thermal Data



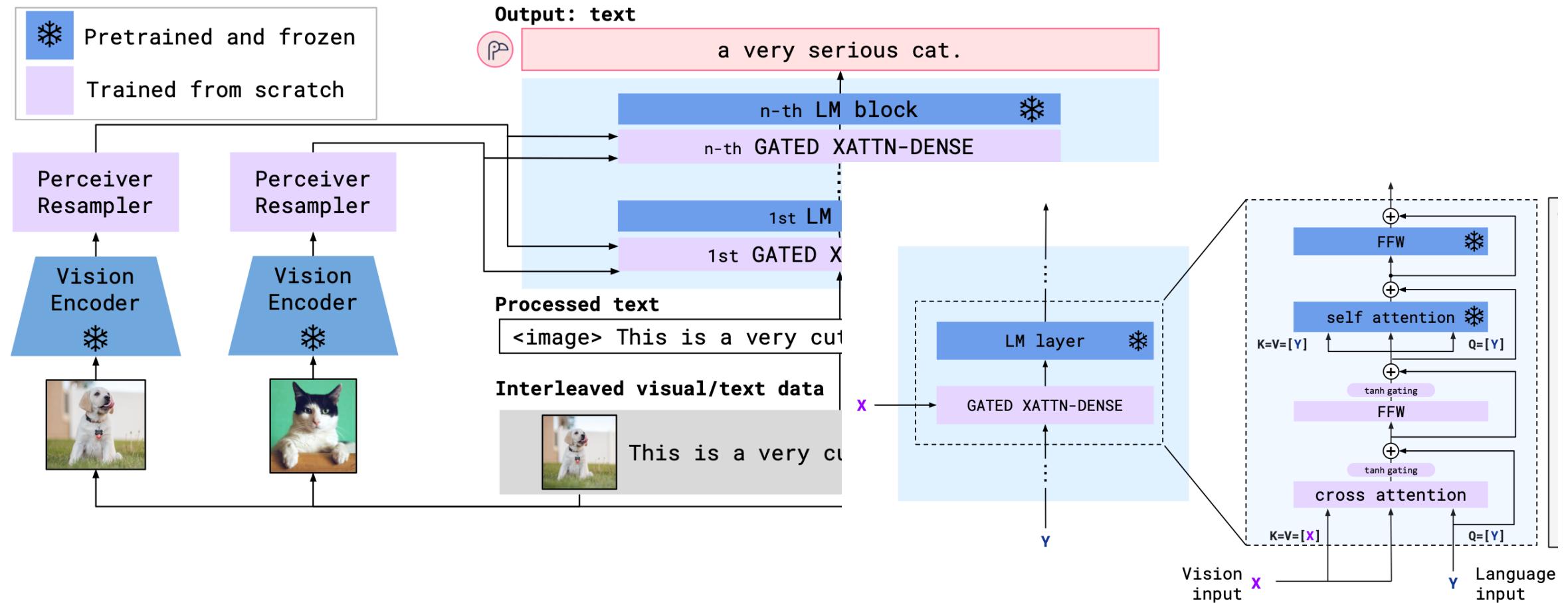
Egocentric Videos



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



# Flamingo



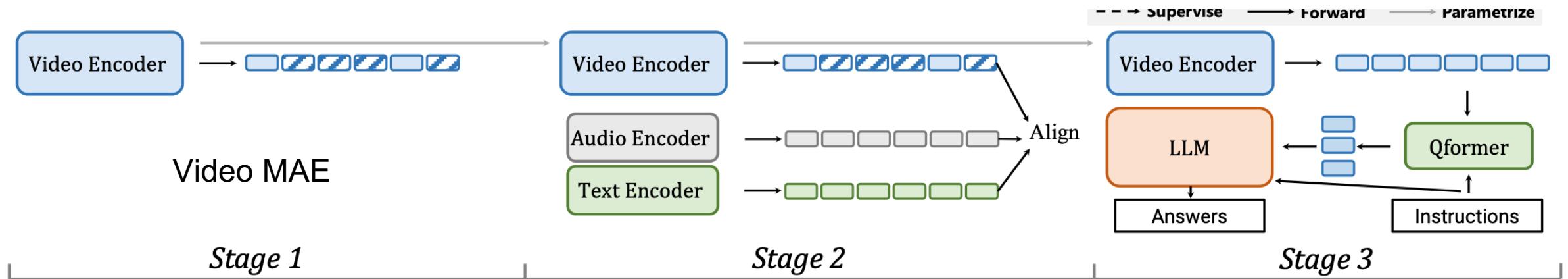


Figure 2: Framework of **InternVideo2**. It consists of three consecutive training phases: unmasked video token reconstruction, multimodal contrastive learning, and next token prediction. In stage 1, the video encoder is trained from scratch, while in stages 2 and 3, it is initialized from the version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



# InternVideo

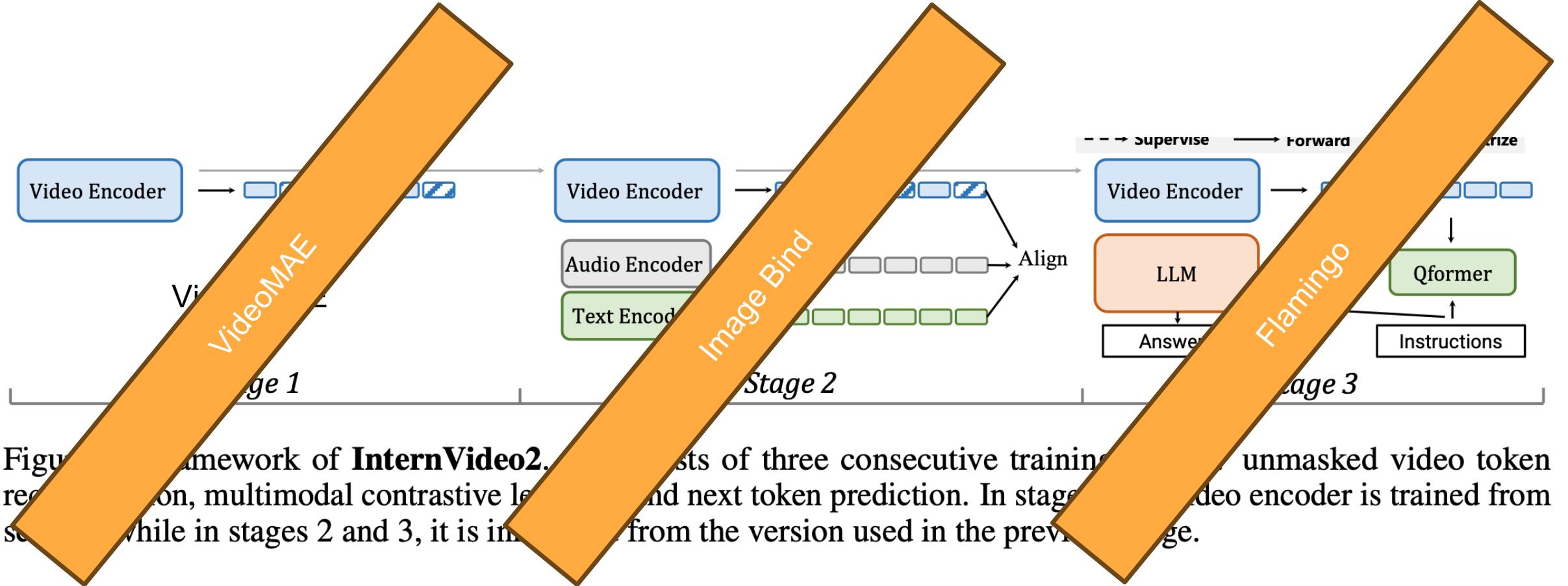


Figure 2: Framework of **InternVideo2**. The framework consists of three consecutive training stages. In stage 1, unmasked video token reconstruction, multimodal contrastive learning, and next token prediction. In stage 2, a video encoder is trained from scratch, while in stages 2 and 3, it is initialized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

# Two types of video understanding tasks

## Video Understanding Tasks

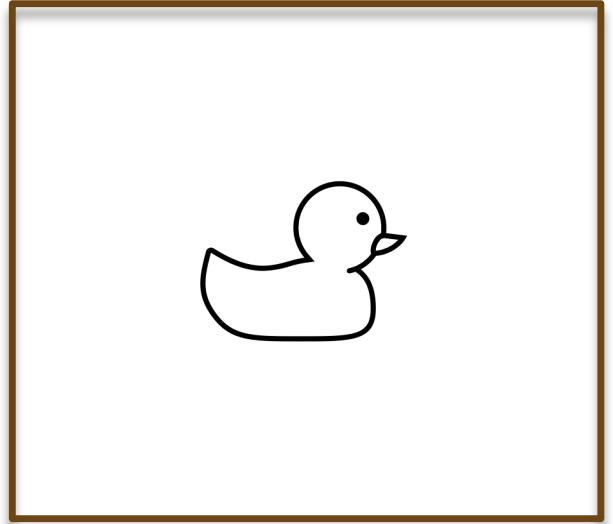
Analogous to Image-based Tasks

Novel Video Tasks



## Image

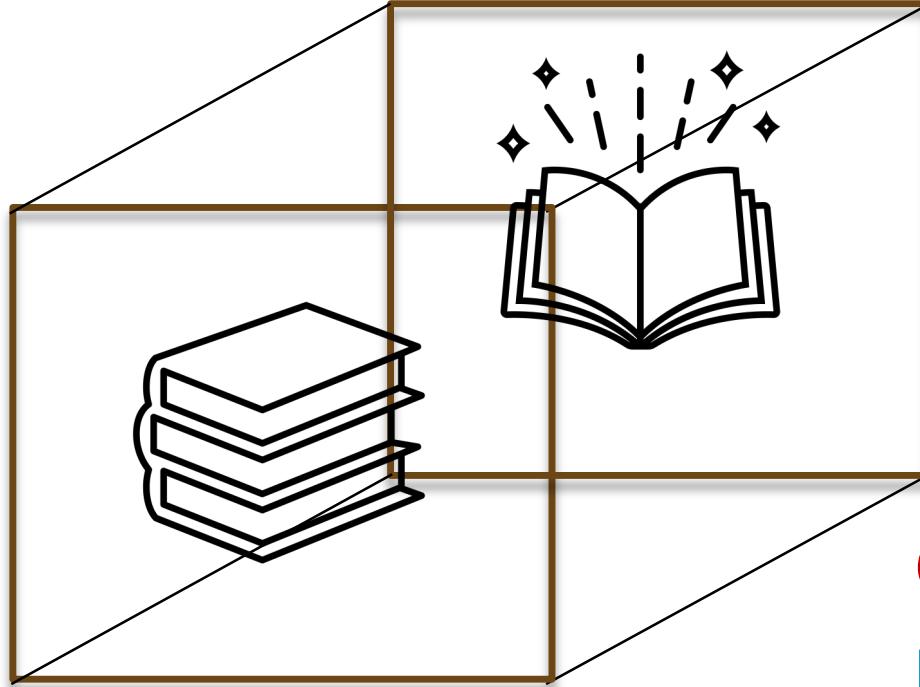
- Object Recognition



Duck

## Video

- Action Recognition

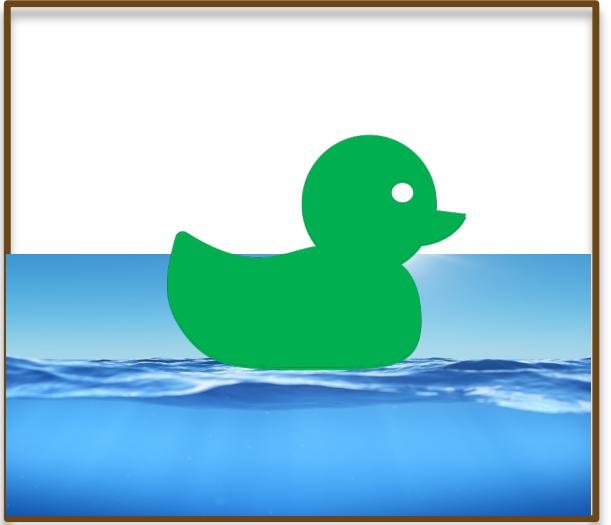


Open  
Book



## Image

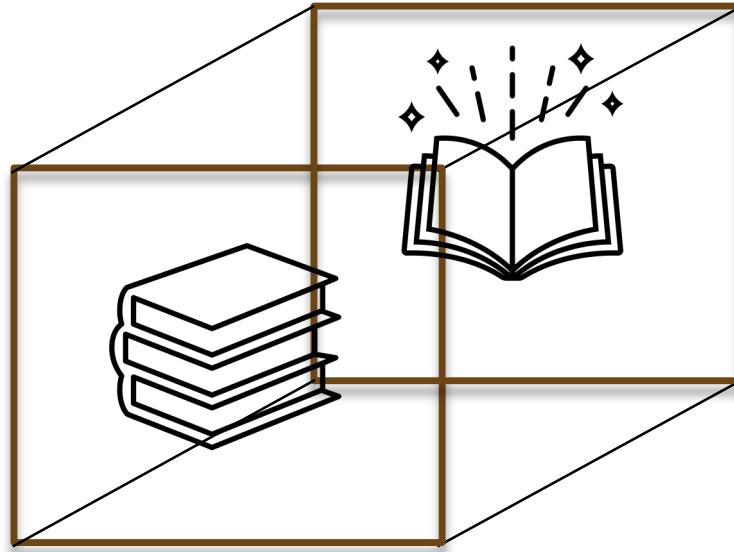
- Image Captioning



A green duck swimming  
In clear water

## Video

- Video Captioning

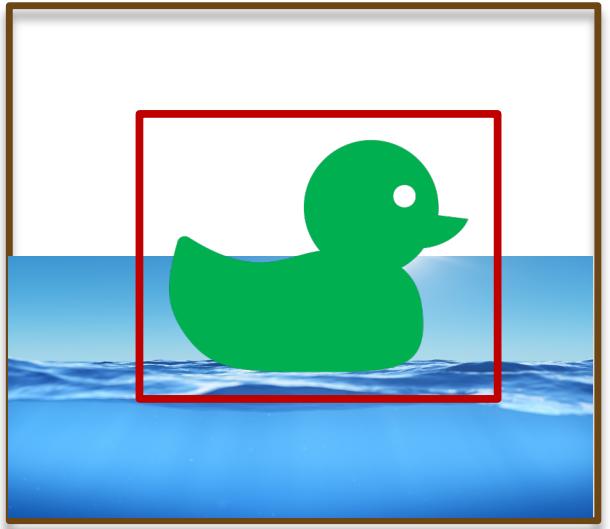


A book picked from top of the pile  
and opened to a page in the middle



## Image

- Object Detection



Duck

## Video

- Action Detection

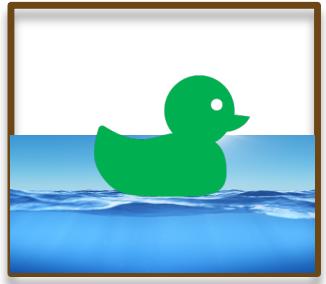


Open Book



## Image

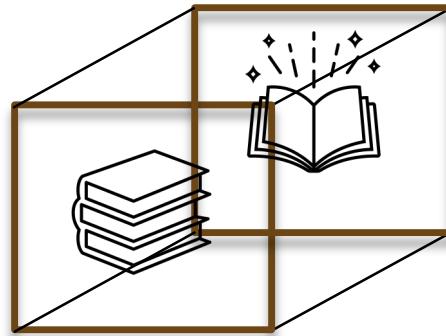
- Image Retrieval



Duck

## Video

- Video Retrieval

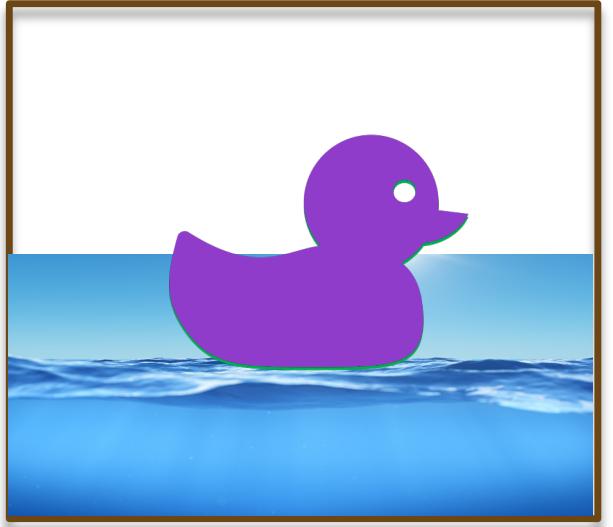


Open Book



## Image

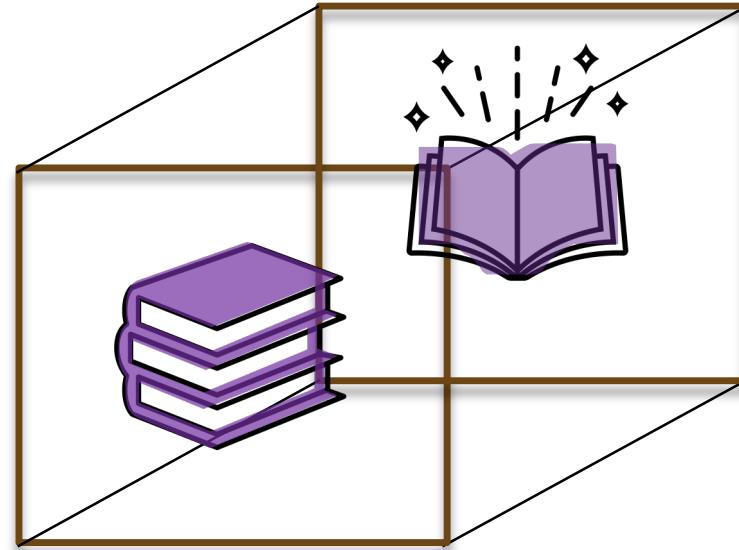
- Object Segmentation



Duck

## Video

- Video Object Segmentation



Book

# EPIC-KITCHENS VISOR

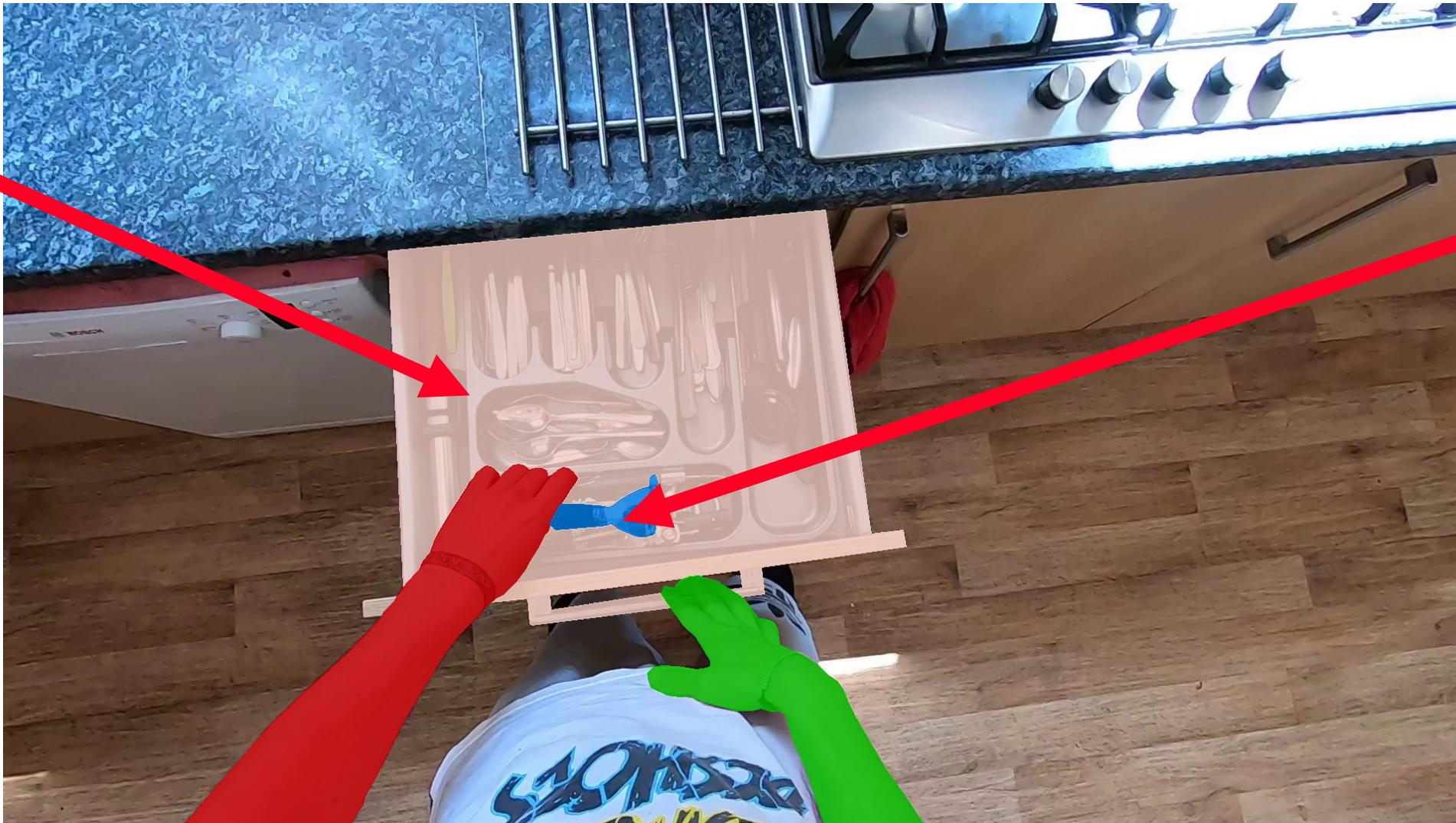
with: Ahmad Darkhalil, Dandan Shan, Bin Zhu, Jian Ma,  
Amlan Kar, Richard Higgins, David Fouhey, Sanja Fidler

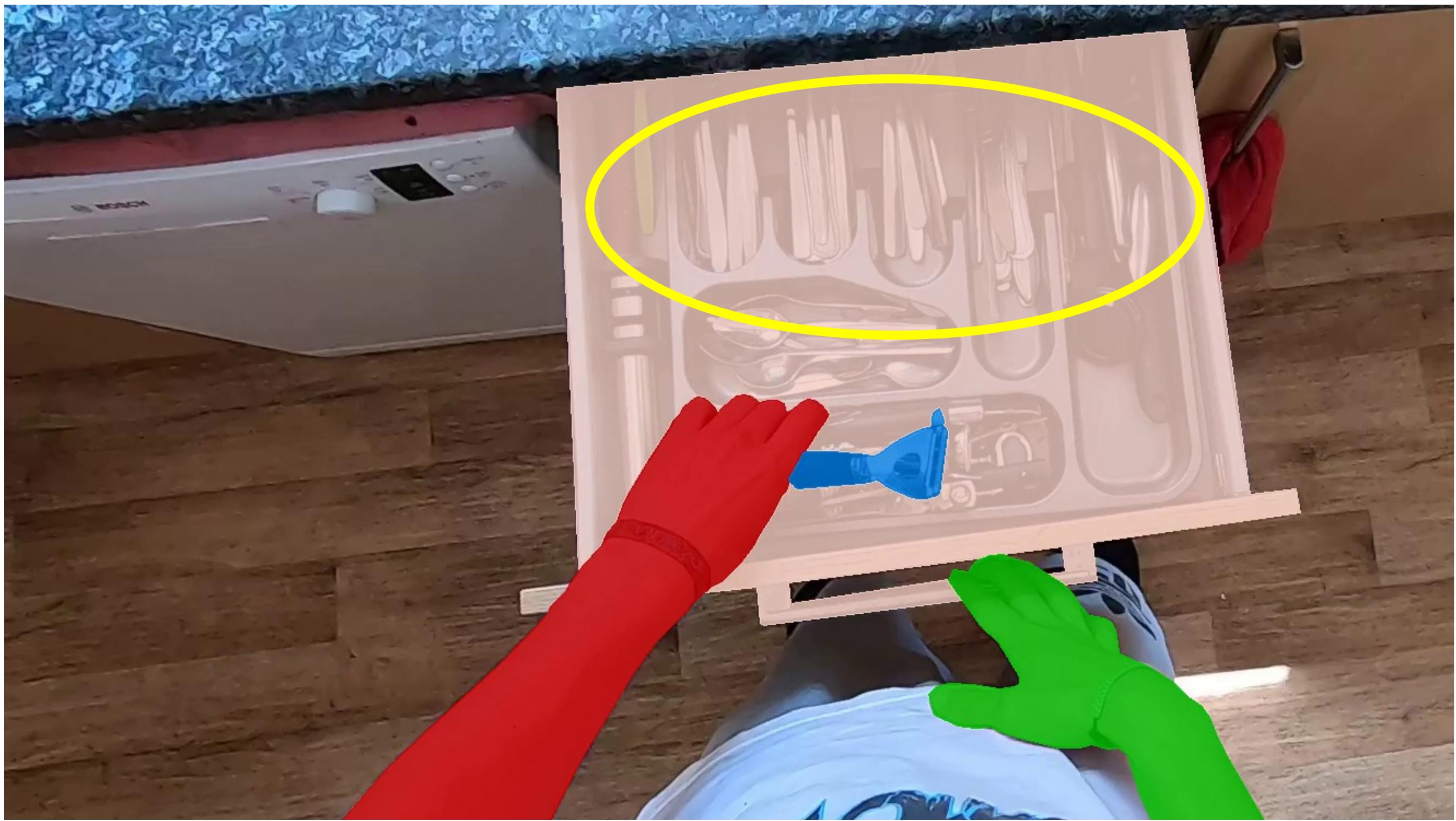


amen  
2025

# 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



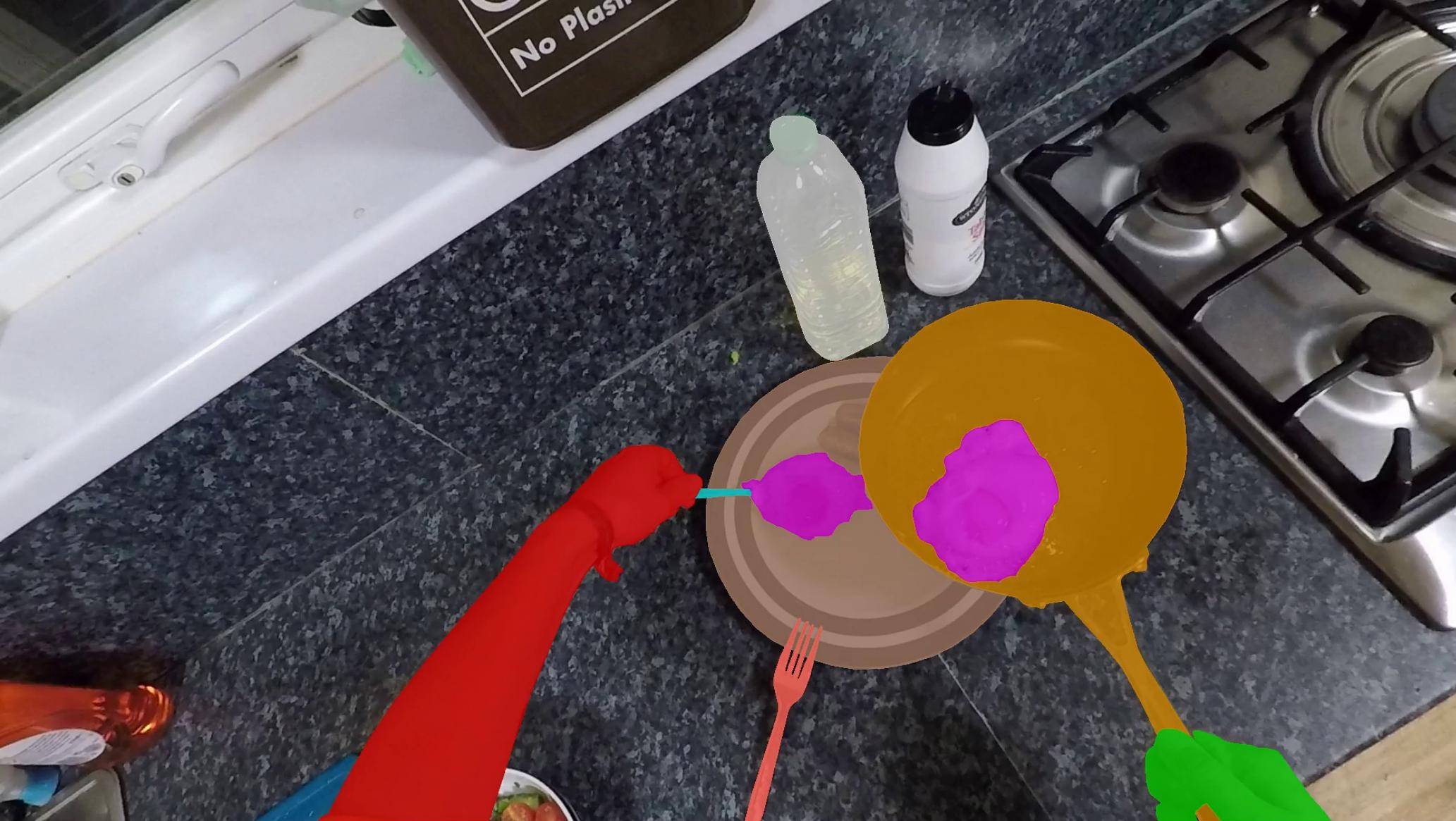
# EPIC-KITCHENS VISOR

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



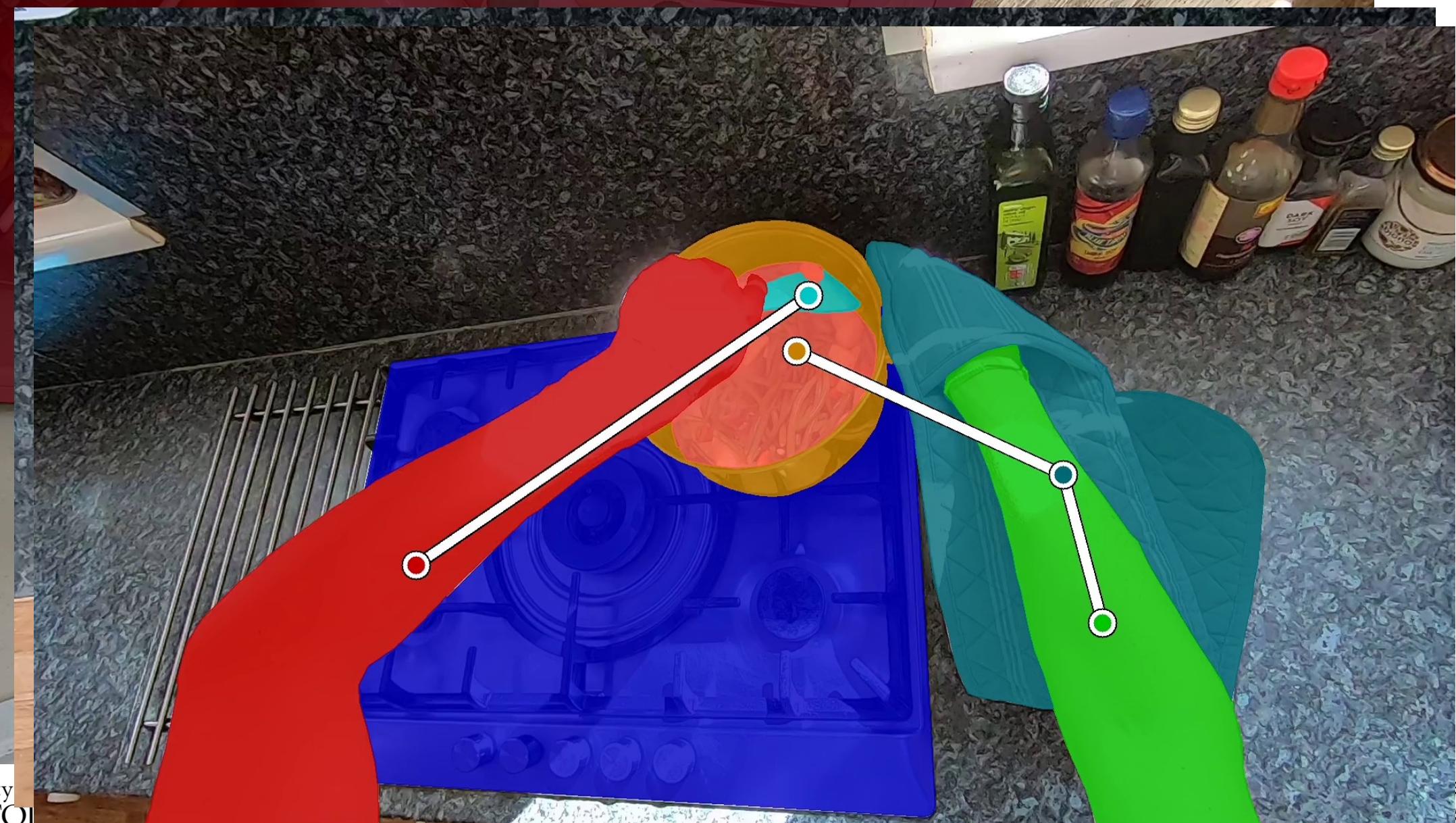
# EPIC-KITCHENS VISOR

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

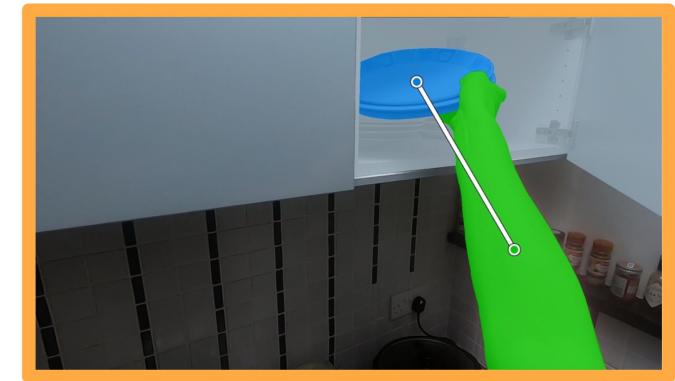
1 Hand, No Contact



2 Hands, No Contact



1 Hand, In Contact



2.7% 41.5%

0.7% 19.4%

27.2% 8.5%



2 Hands, 2 Obj Contacts



2 Hands, Same Contact



2 Hands, 1 In Contact

# EPIC-KITCHENS VISOR

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

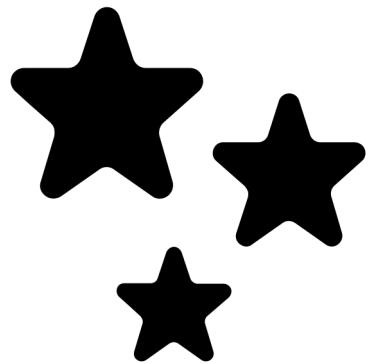




# Analogous Tasks

## Image

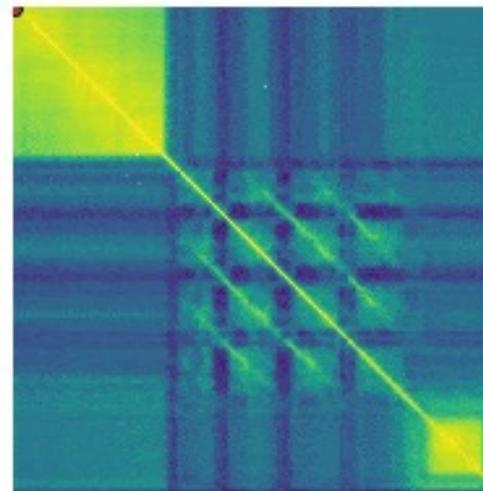
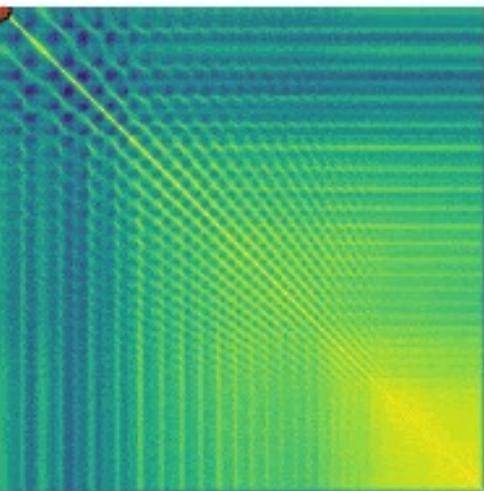
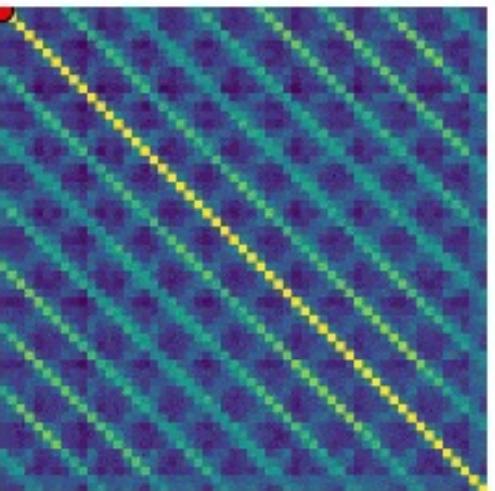
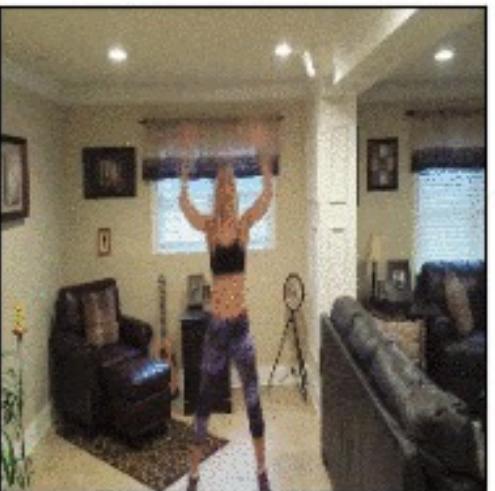
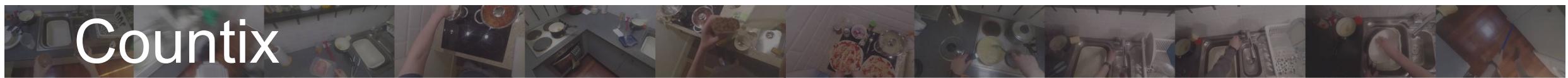
- Object Counting



## Video

- Action Counting





# Every Shot Counts

with: Saptarshi Sinha  
Alexandros Stergiou

RepCount



GT:6

Pred:6

Countix



GT:9

Pred:9

RepCount



GT:32

Pred:32



# Analogous Tasks

## Image

- Text-to-image Generation



Stable Diffusion

## Video

- Text-to-Video Generation



SORA

# 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  
PAISS 2025

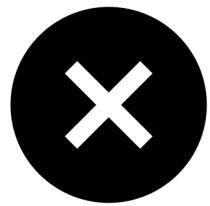


Gemini Veo 3 – May 2025

Dima Damen  
PAISS 2025

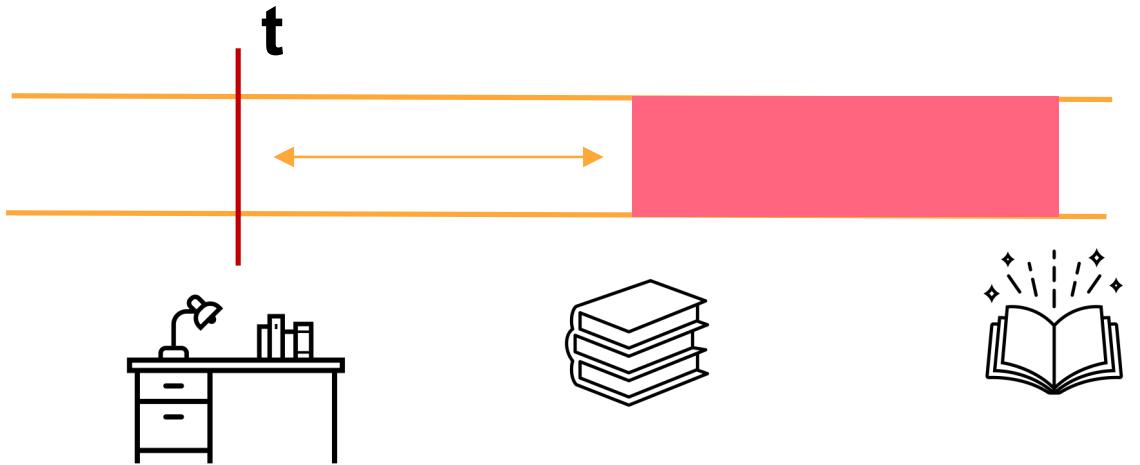


Image



Video

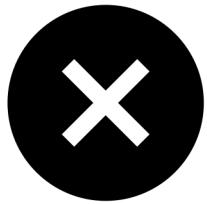
- Action Anticipation  
What will happen after 1 second?



Open Book

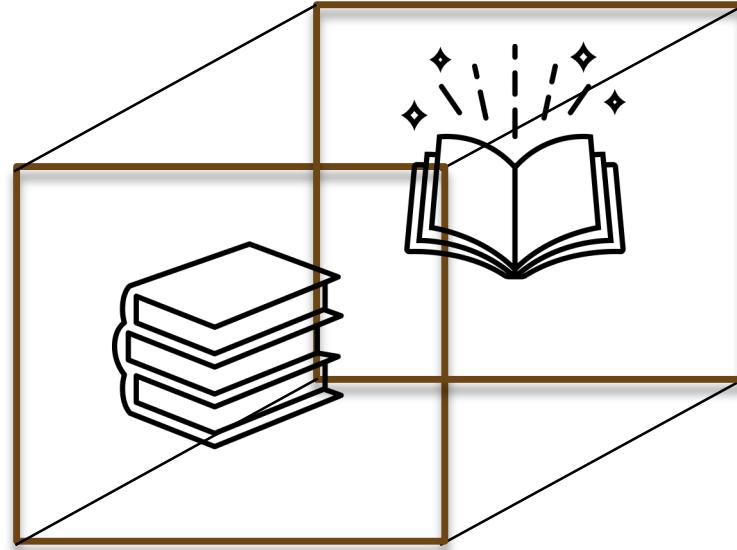
# Non-Analogous Tasks

Image



Video

- Manner of action  
How did you open the book?





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**...

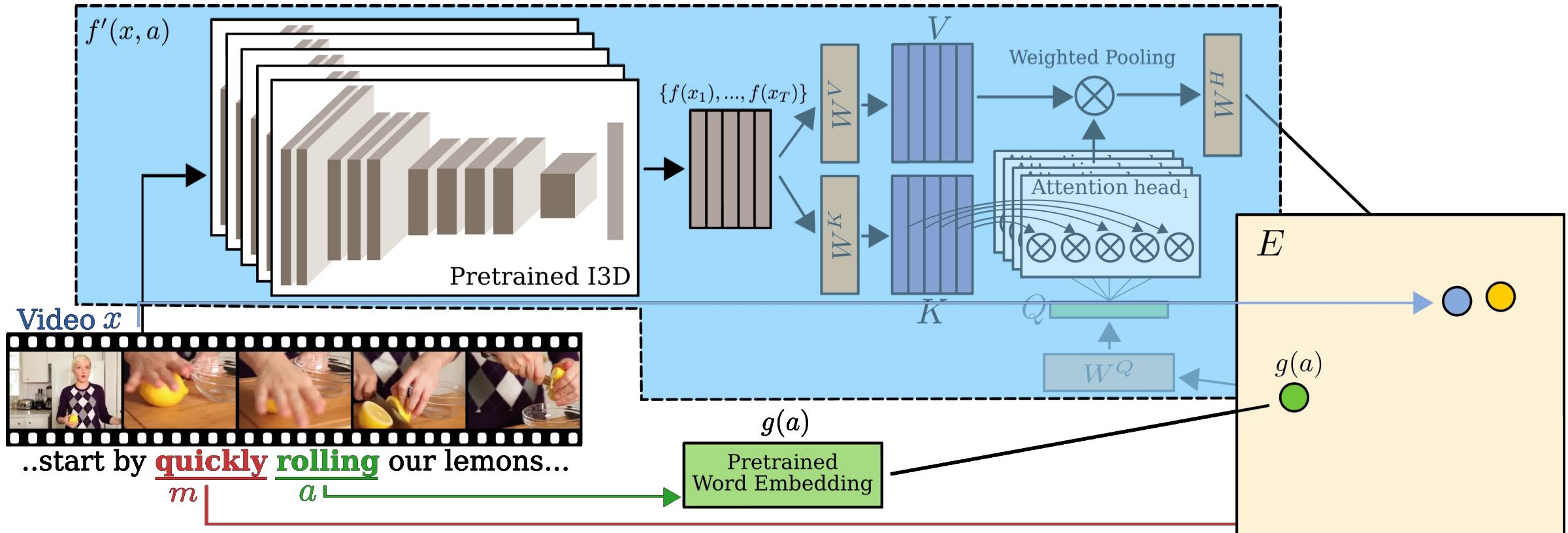
-10 seconds

timestamp

+10 seconds



with: Hazel Doughty  
Ivan Laptev  
Walterio Mayol-Cuevas



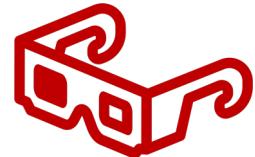


# How?

with: Hazel Doughty  
Ivan Laptev  
Walterio Mayol-Cuevas



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



Motivation and Datasets in  
Egocentric Video Understanding



Video Understanding  
Out of the Frame



Video Understanding:  
Data and Tasks



Teaser: The Wizard of Oz  
& Genie 3



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





with: Vangelis Kazakos  
Arsha Nagrani.  
Andrew Zisserman

Jaesung Huh  
Jacob Chalk

# Multi-modal learning...

- 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

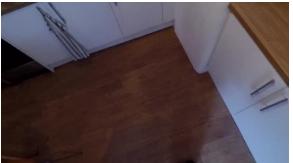
\* : Equal contribution



Dima Damen  
PAISS 2025

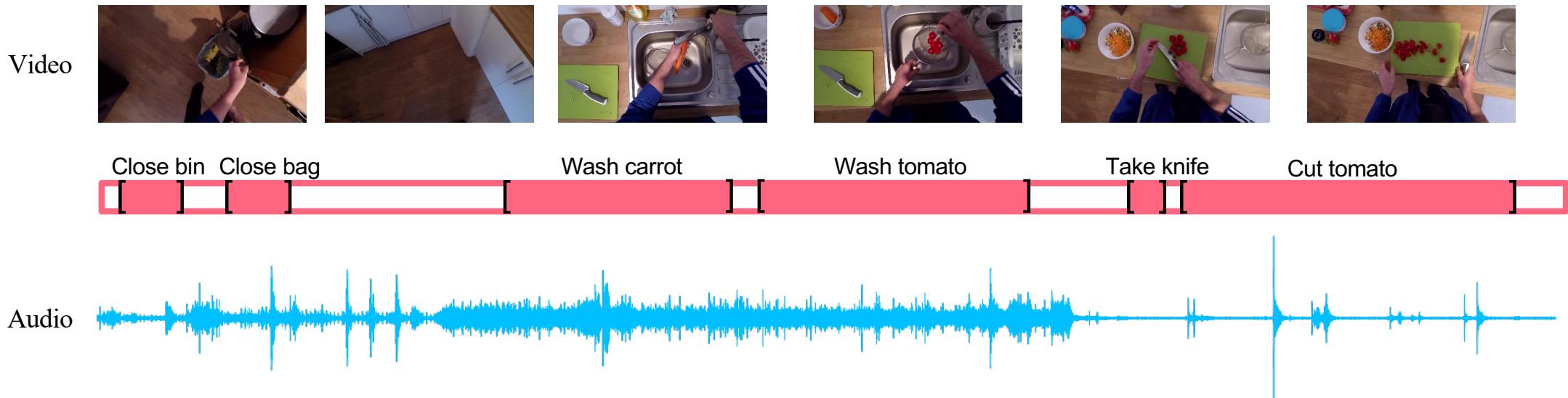


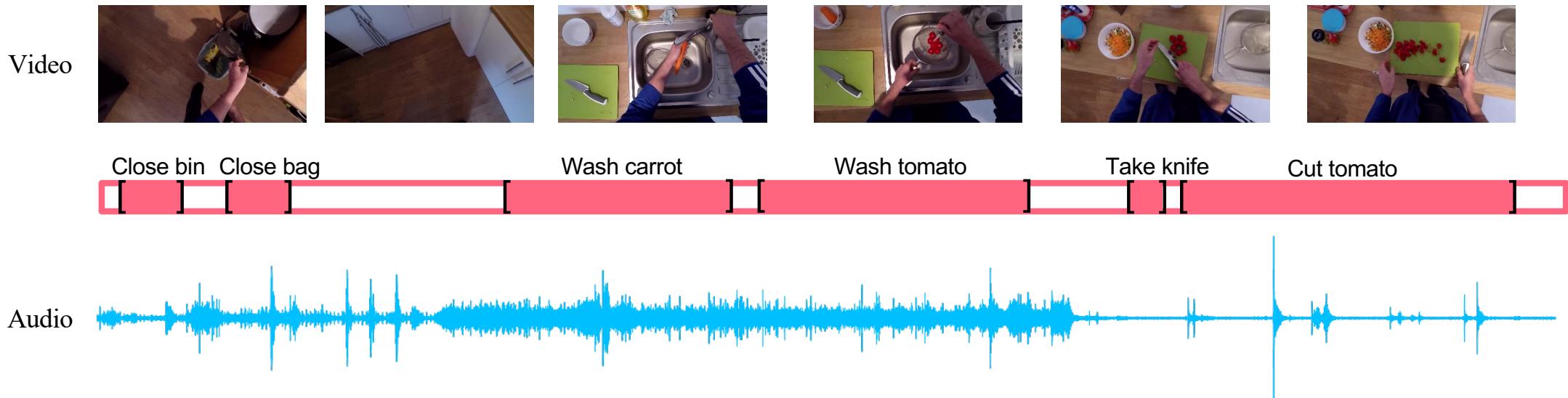
Video

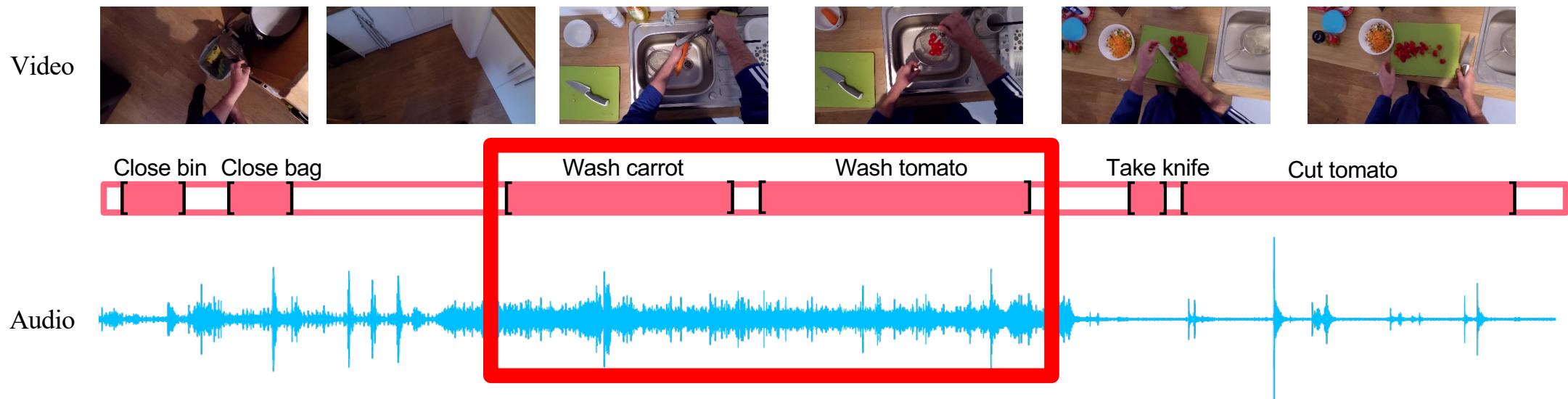


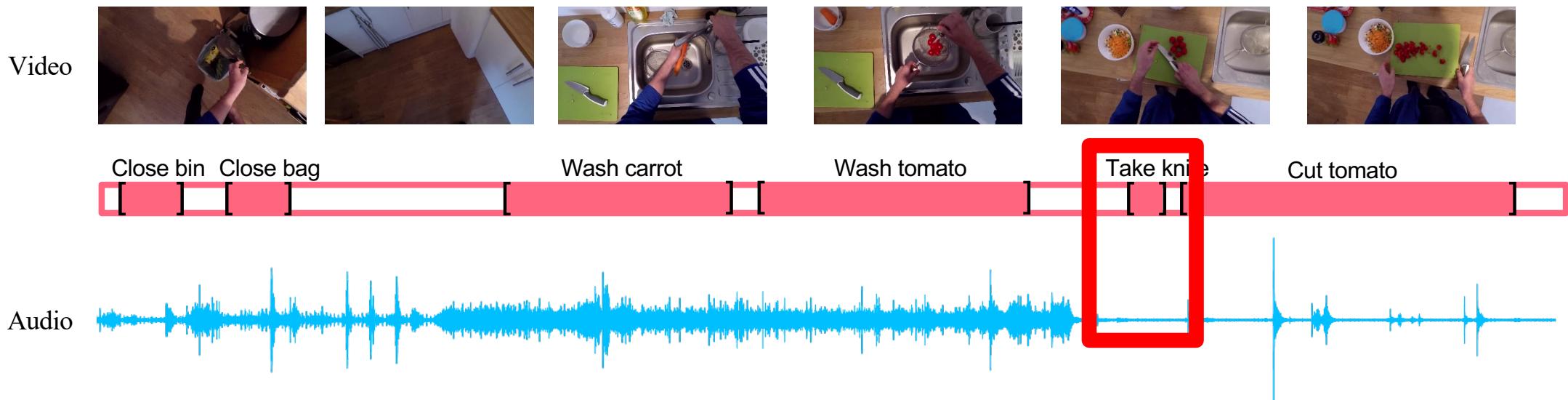
Audio

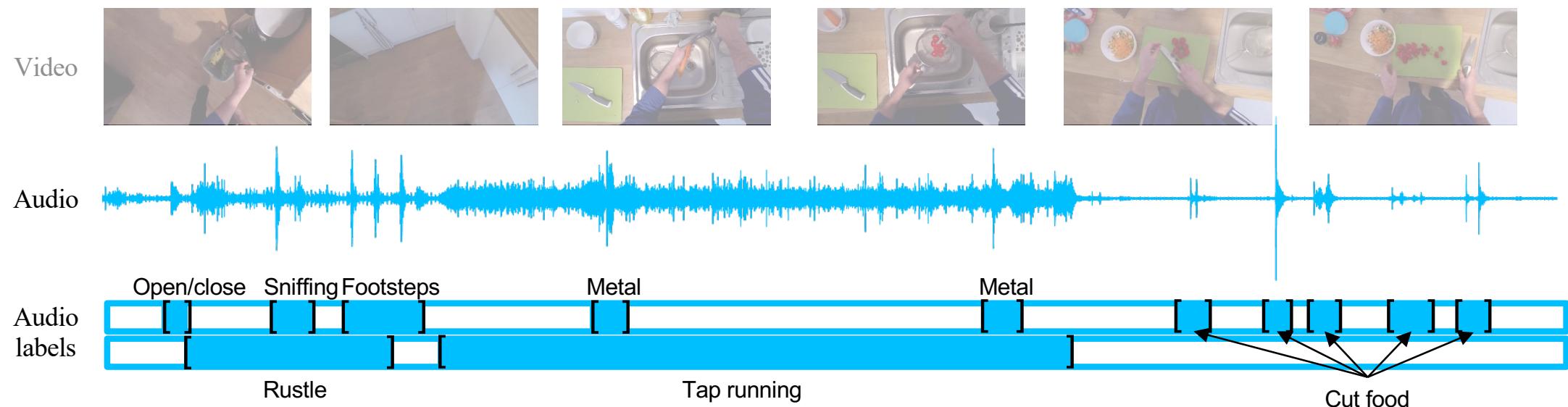








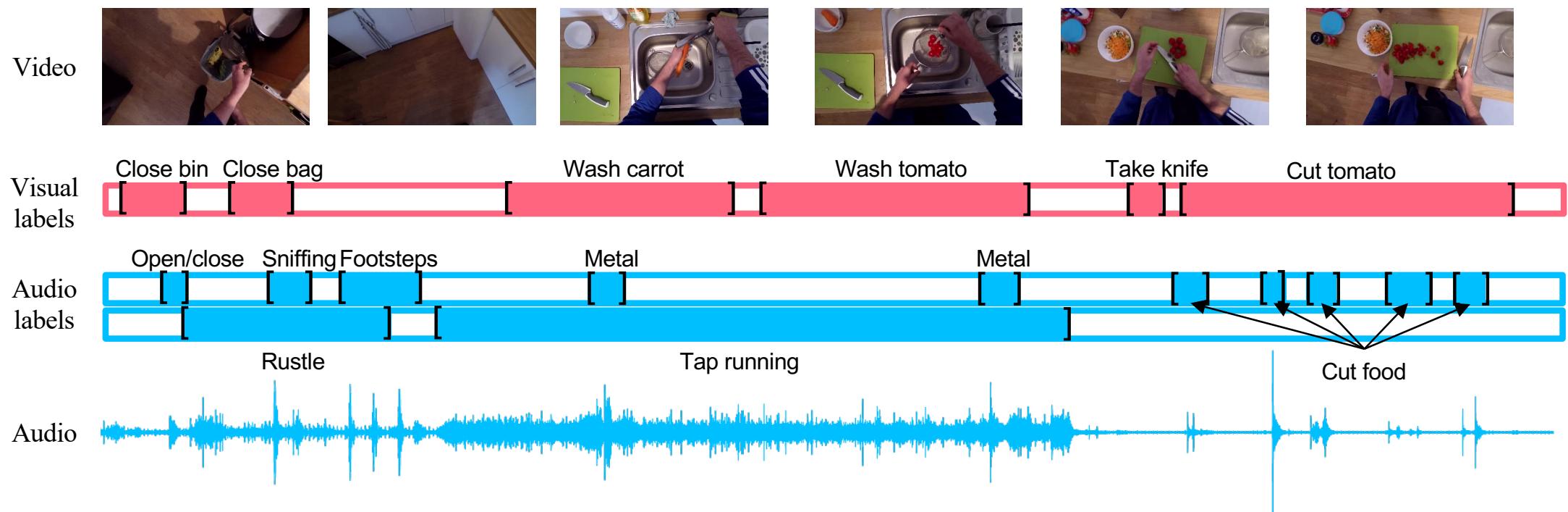






# Motivation

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





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

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

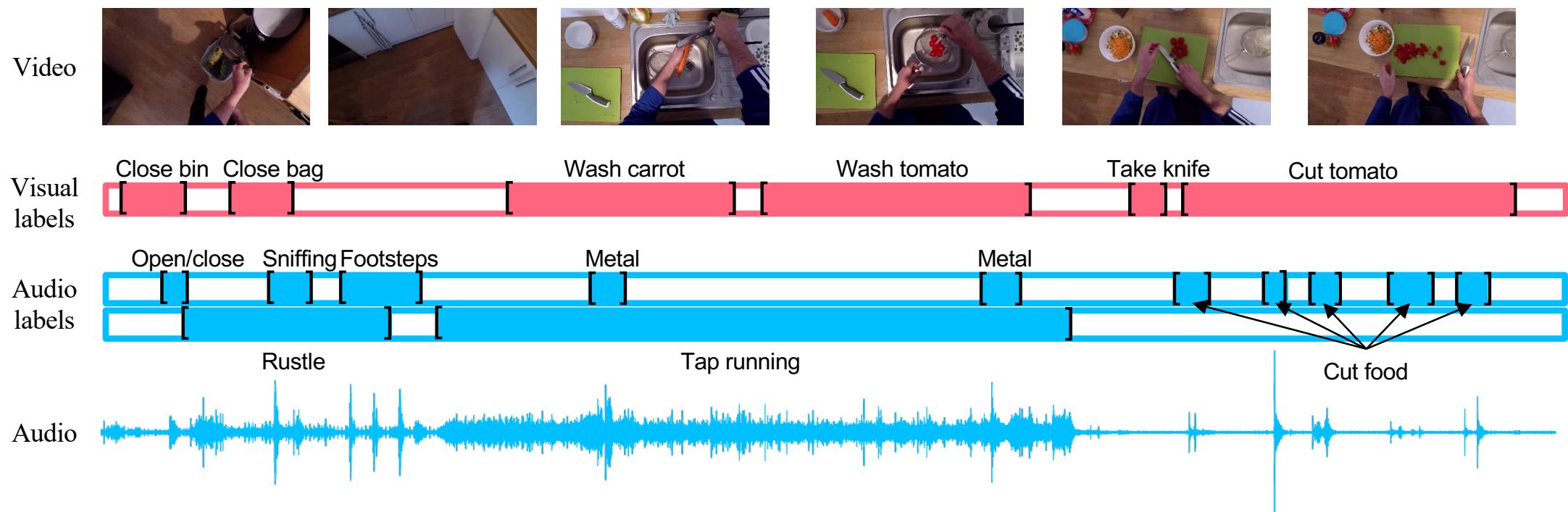
\* : Equal contribution



Dima Damen  
PAISS 2025

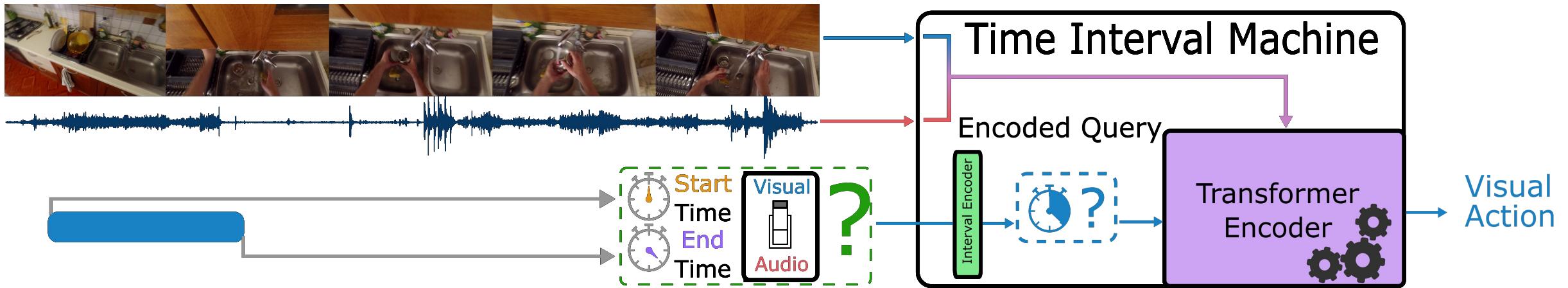
# Multi-Modal Long-Form Dataset

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



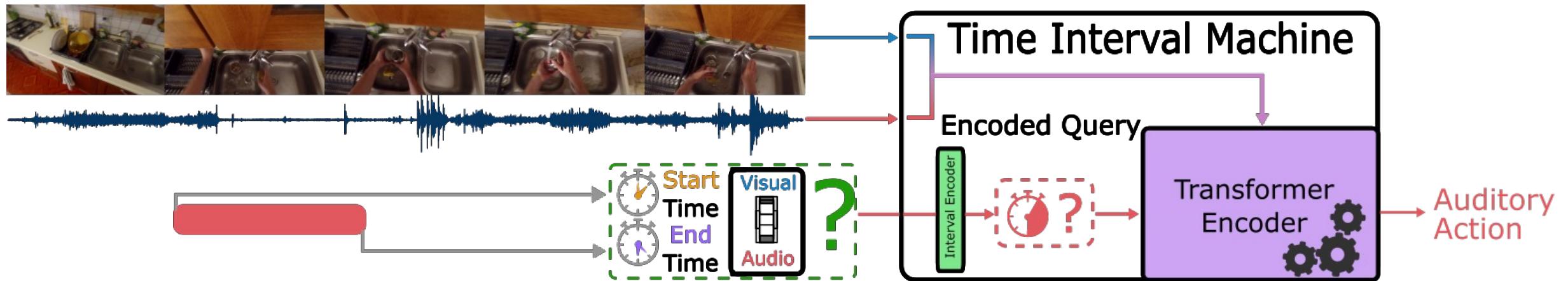
# TIM: A Time-Interval Audio-Visual Machine

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



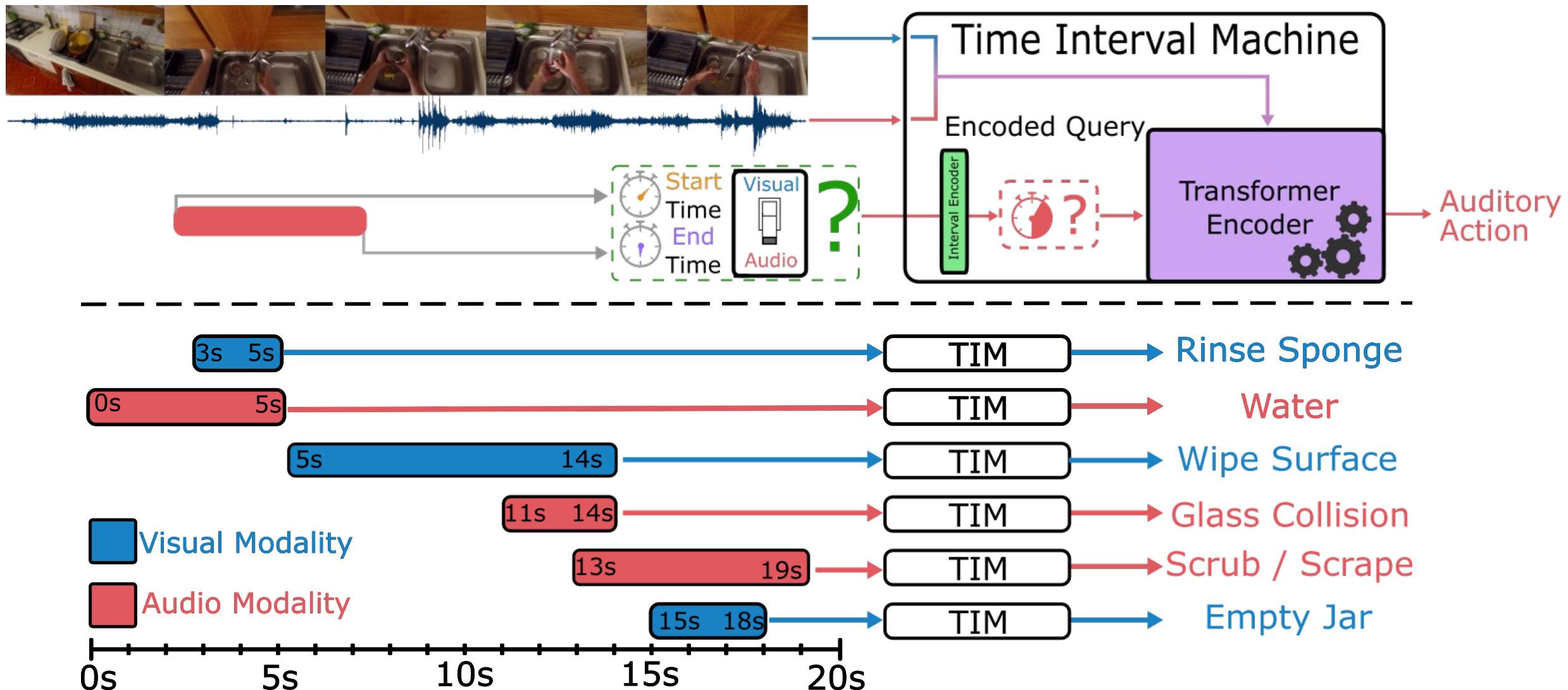
# TIM: A Time-Interval Audio-Visual Machine

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



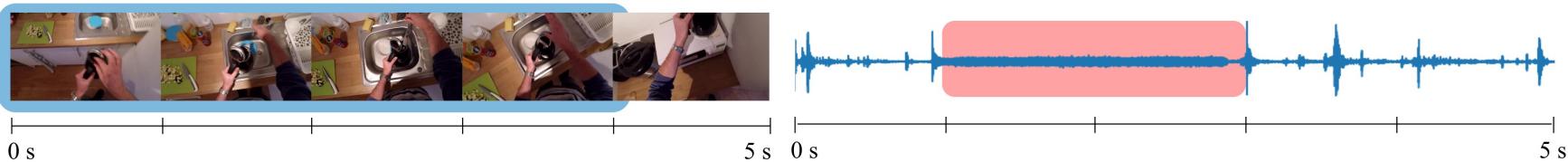
# TIM: A Time-Interval Audio-Visual Machine

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



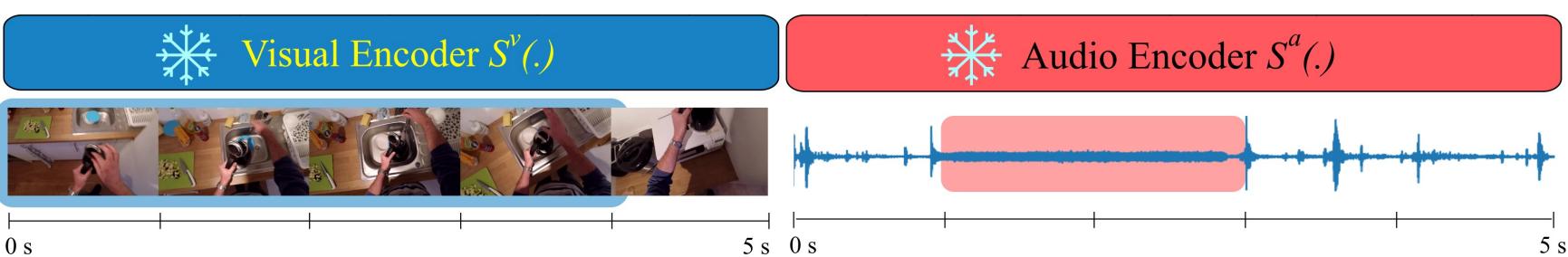
# TIM: A Time-Interval Audio-Visual Machine

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



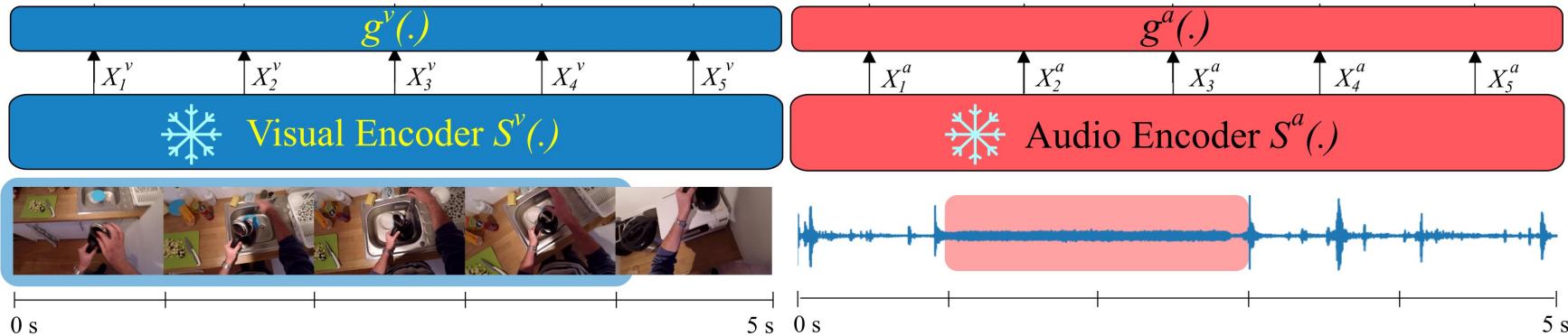
# TIM: A Time-Interval Audio-Visual Machine

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



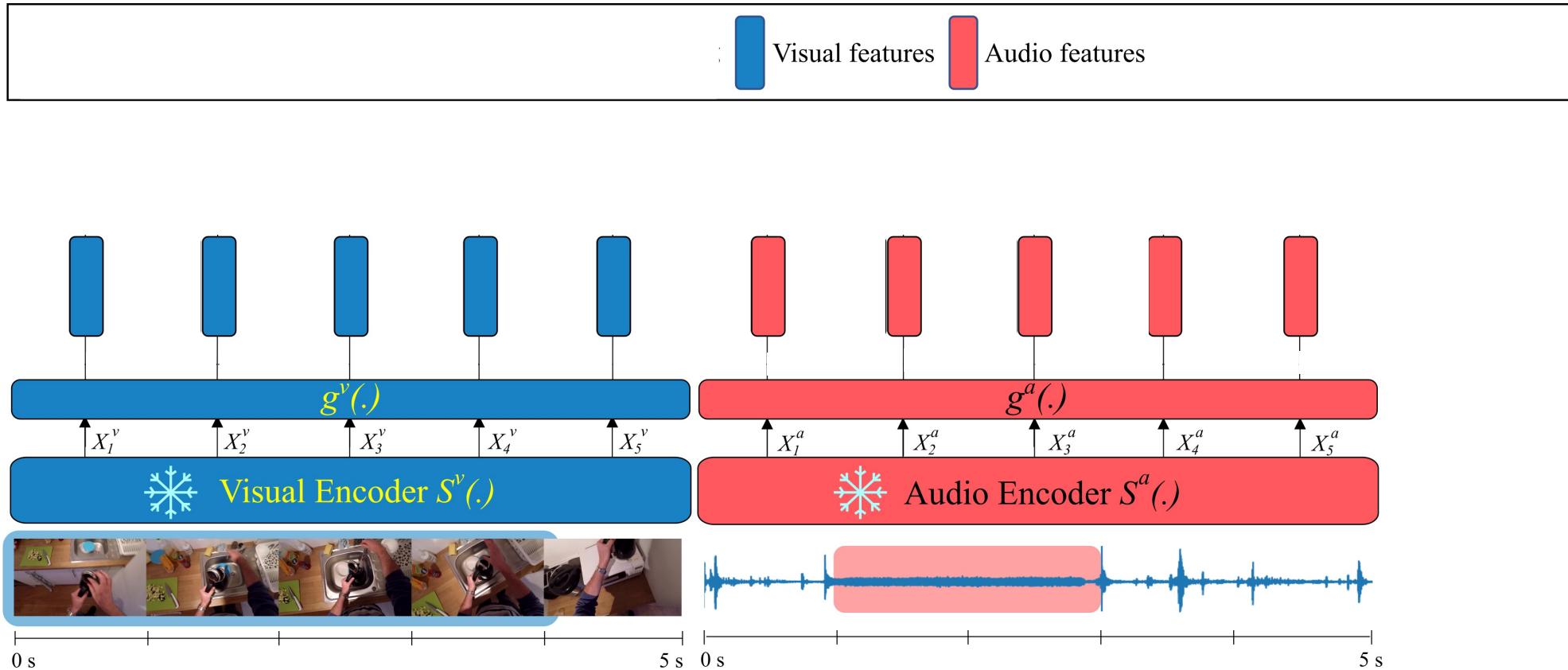
# TIM: A Time-Interval Audio-Visual Machine

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



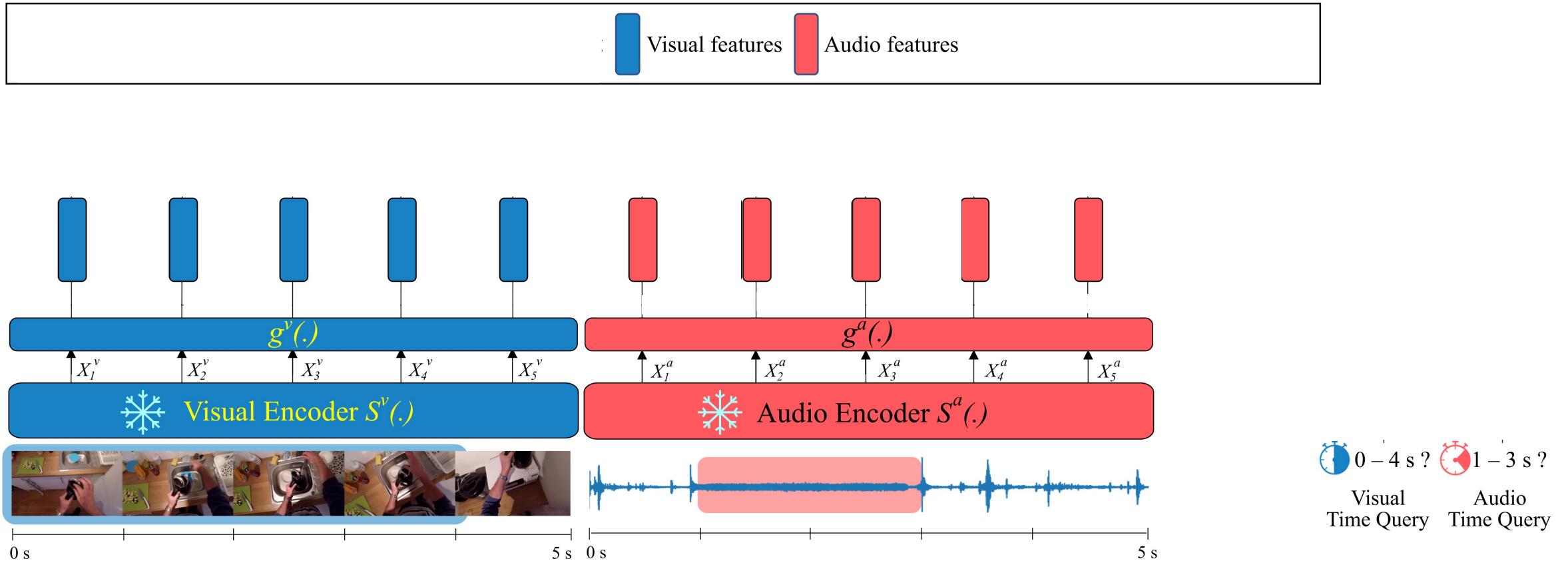
# TIM: A Time-Interval Audio-Visual Machine

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



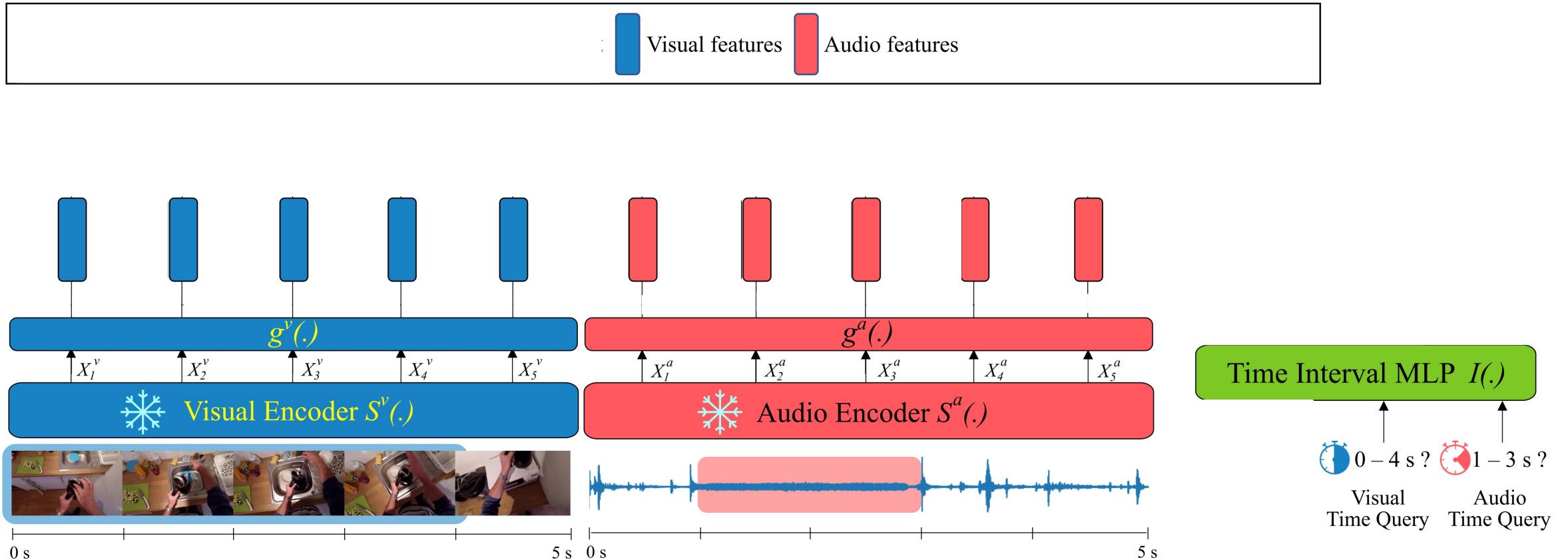
# TIM: A Time-Interval Audio-Visual Machine

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



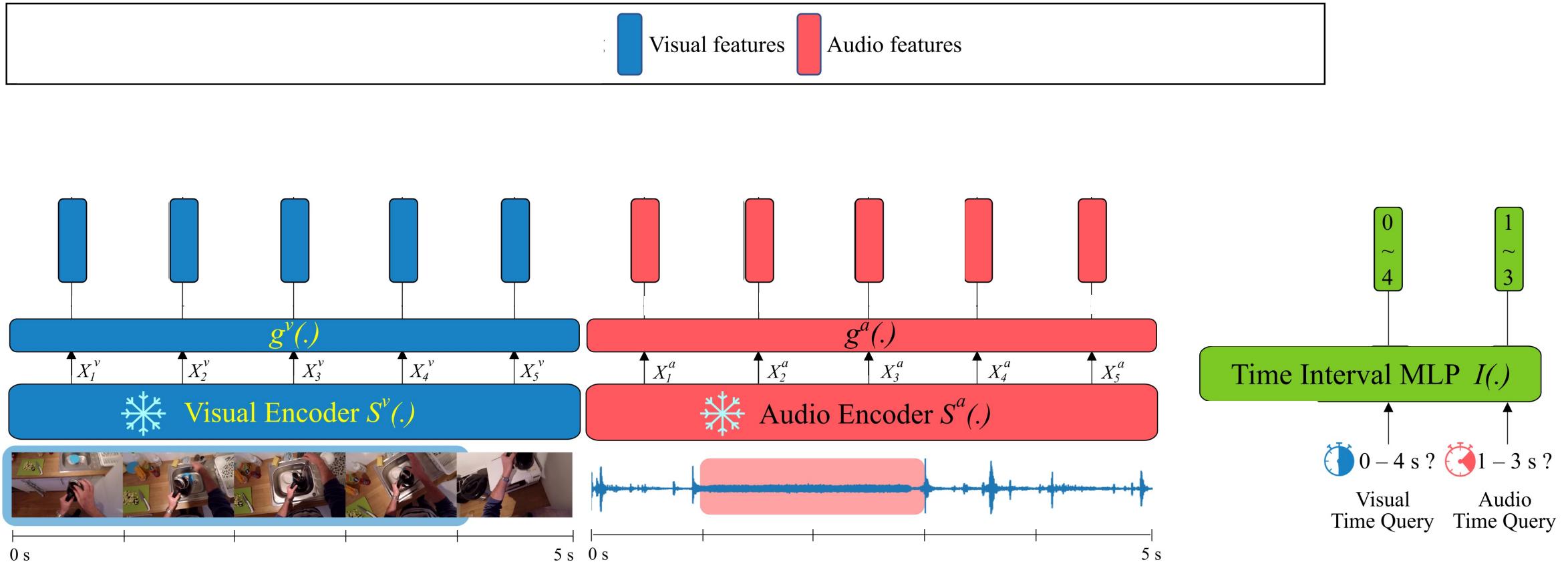
# TIM: A Time-Interval Audio-Visual Machine

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



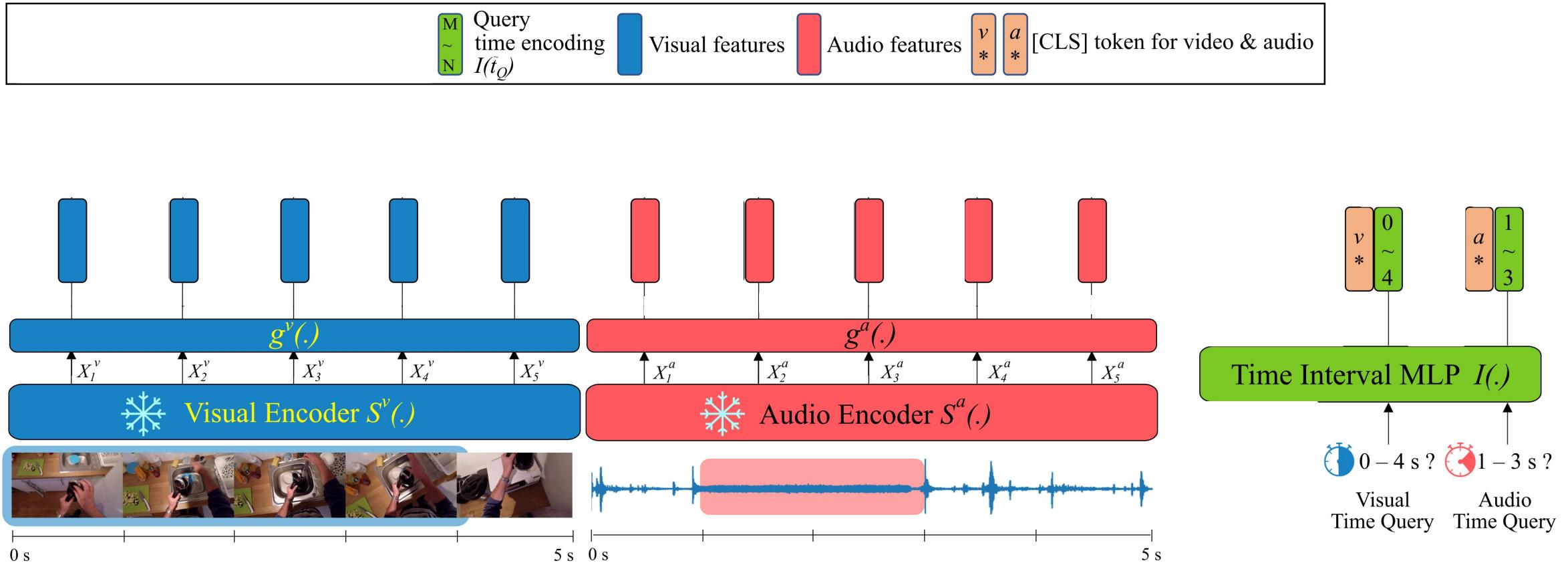
# TIM: A Time-Interval Audio-Visual Machine

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



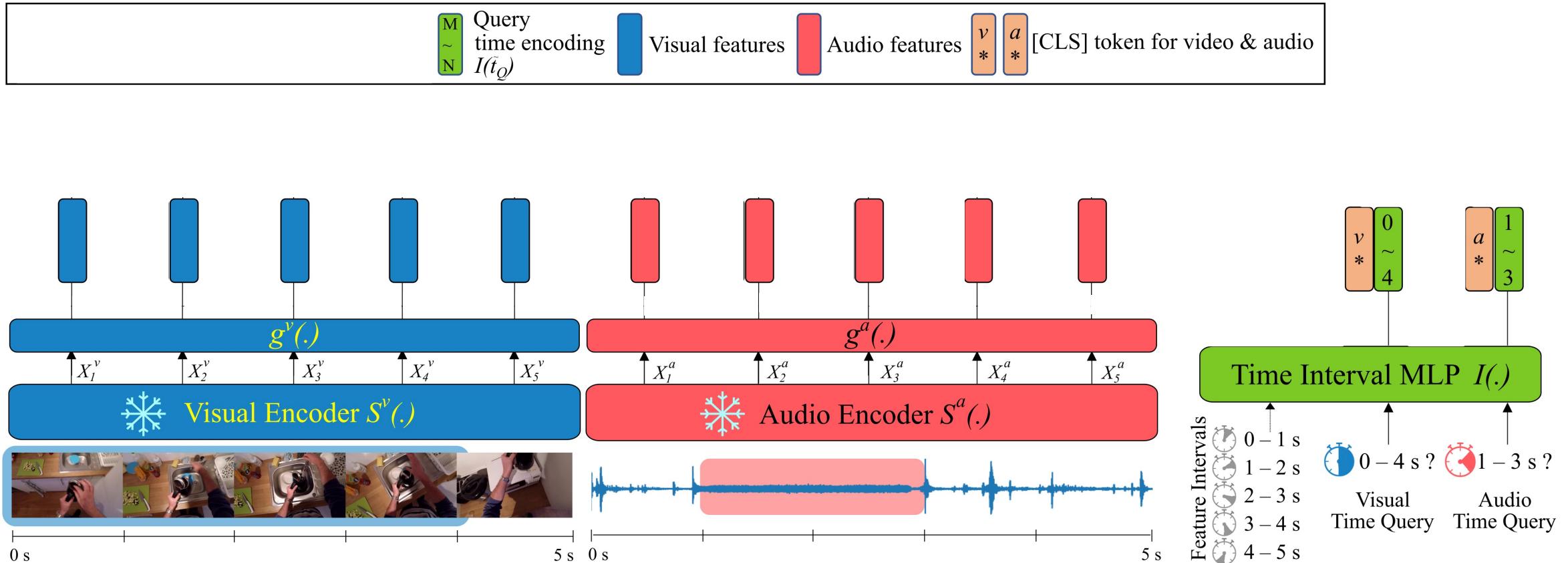
# TIM: A Time-Interval Audio-Visual Machine

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



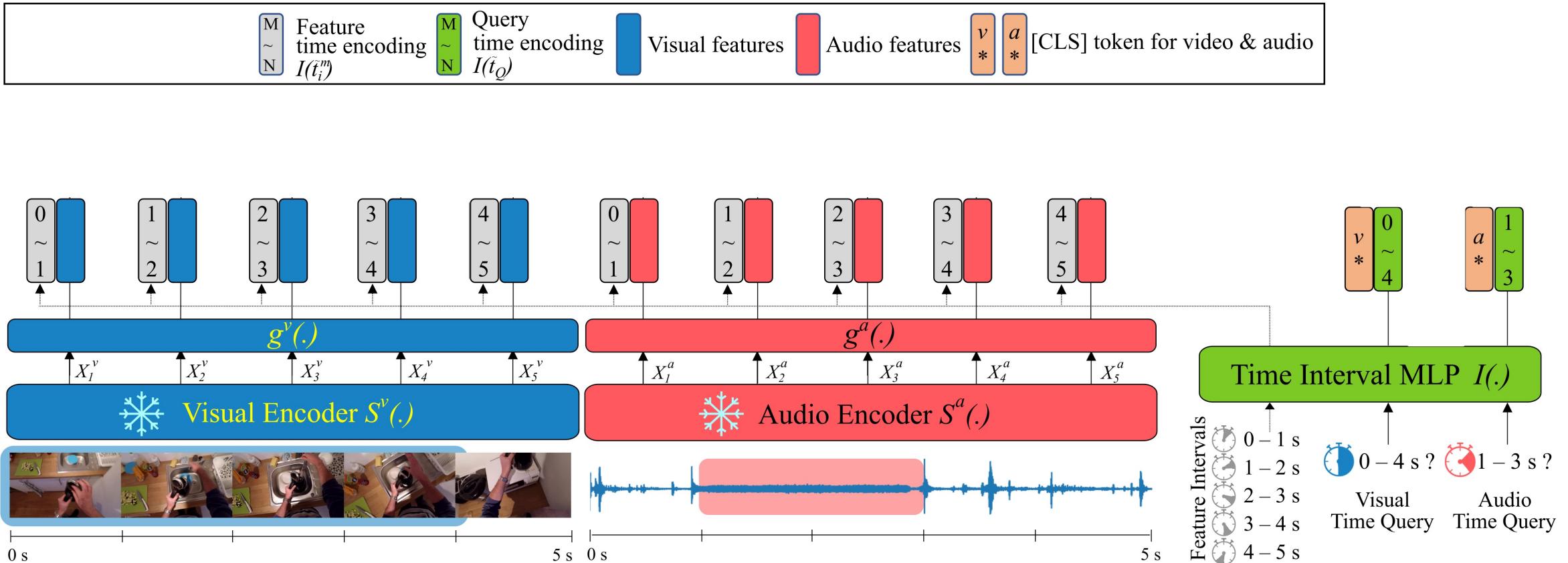
# TIM: A Time-Interval Audio-Visual Machine

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



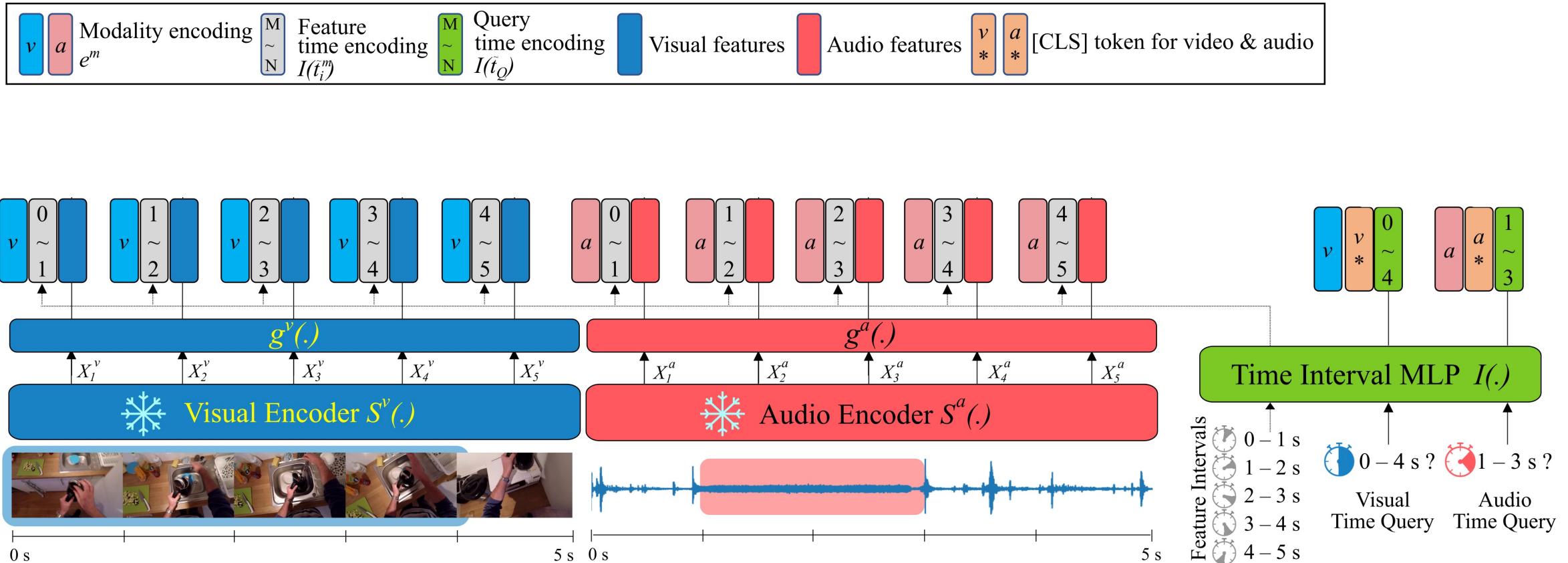
# TIM: A Time-Interval Audio-Visual Machine

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



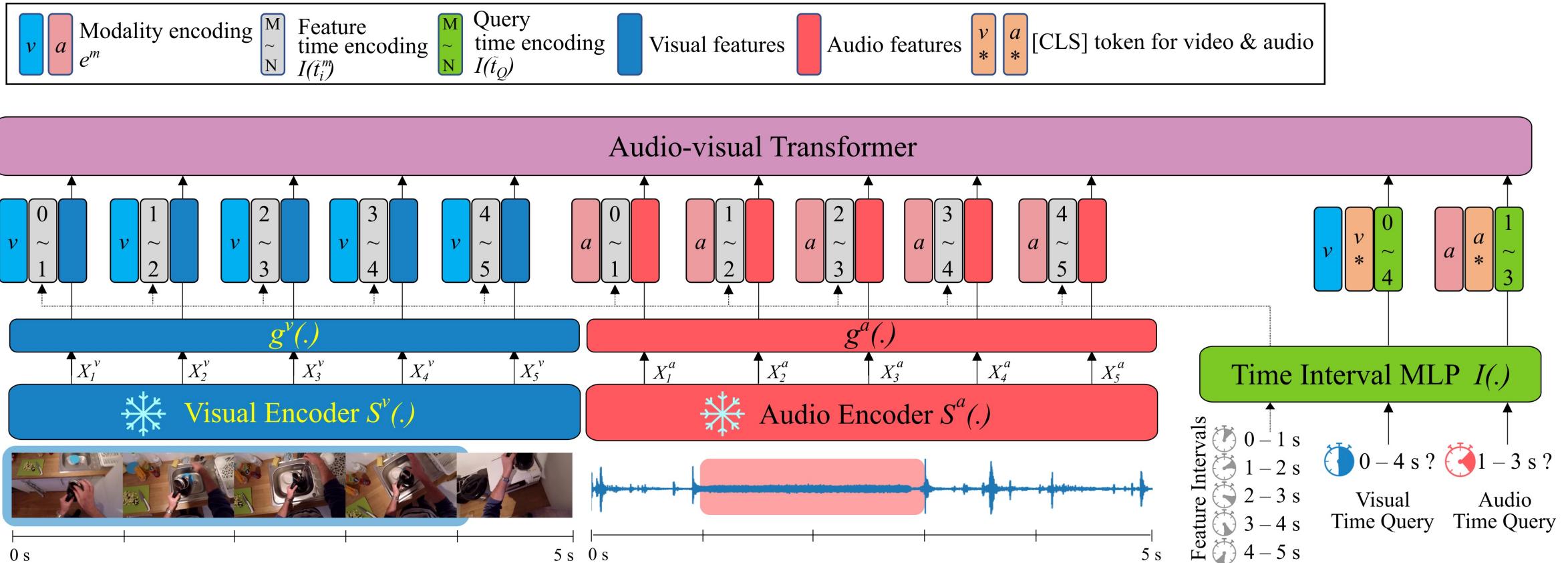
# TIM: A Time-Interval Audio-Visual Machine

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



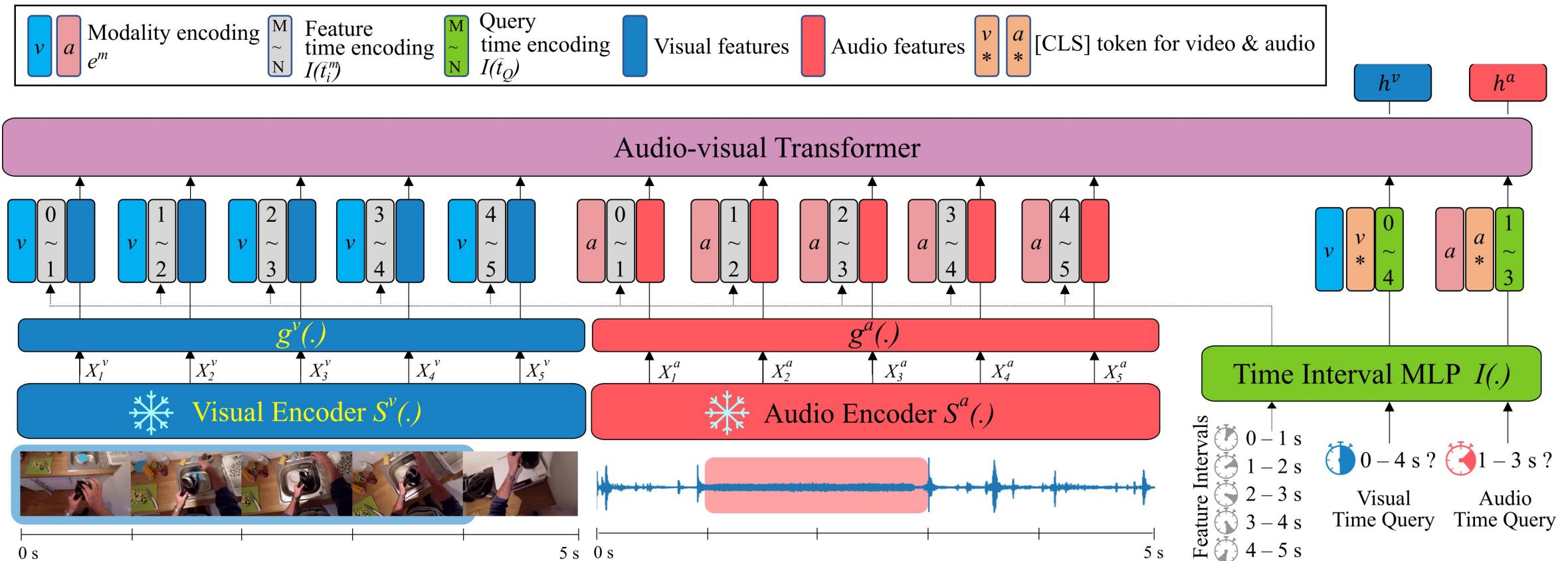
# TIM: A Time-Interval Audio-Visual Machine

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



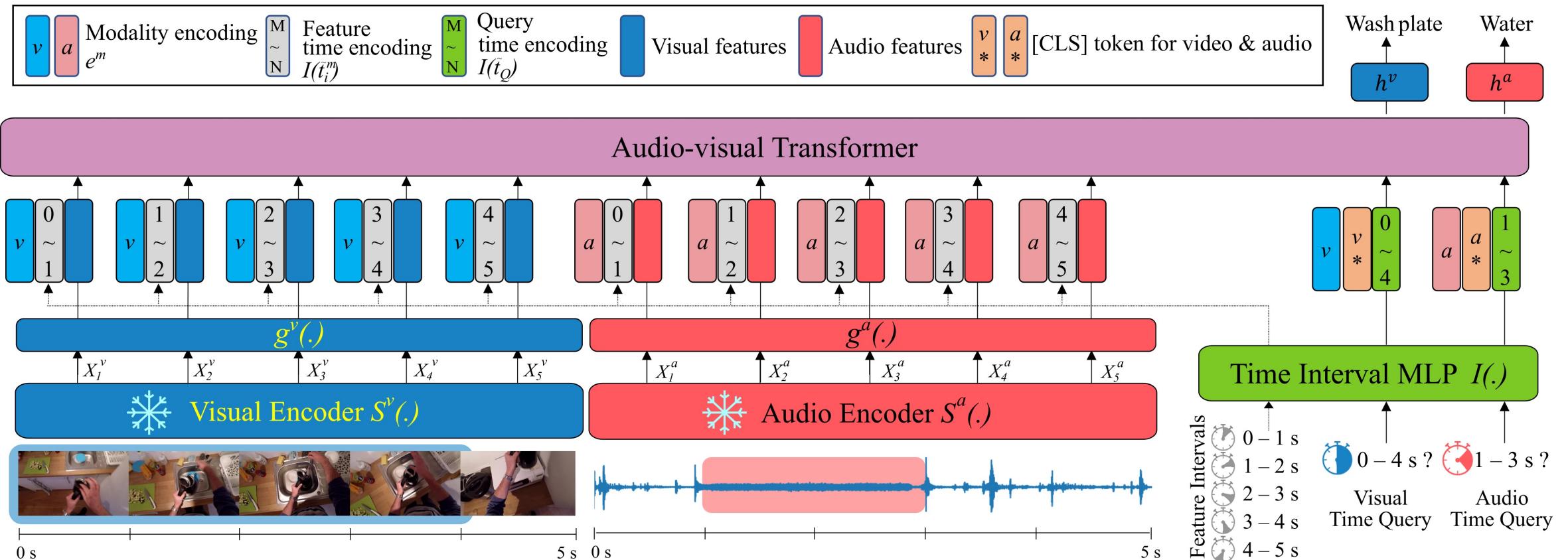
# TIM: A Time-Interval Audio-Visual Machine

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

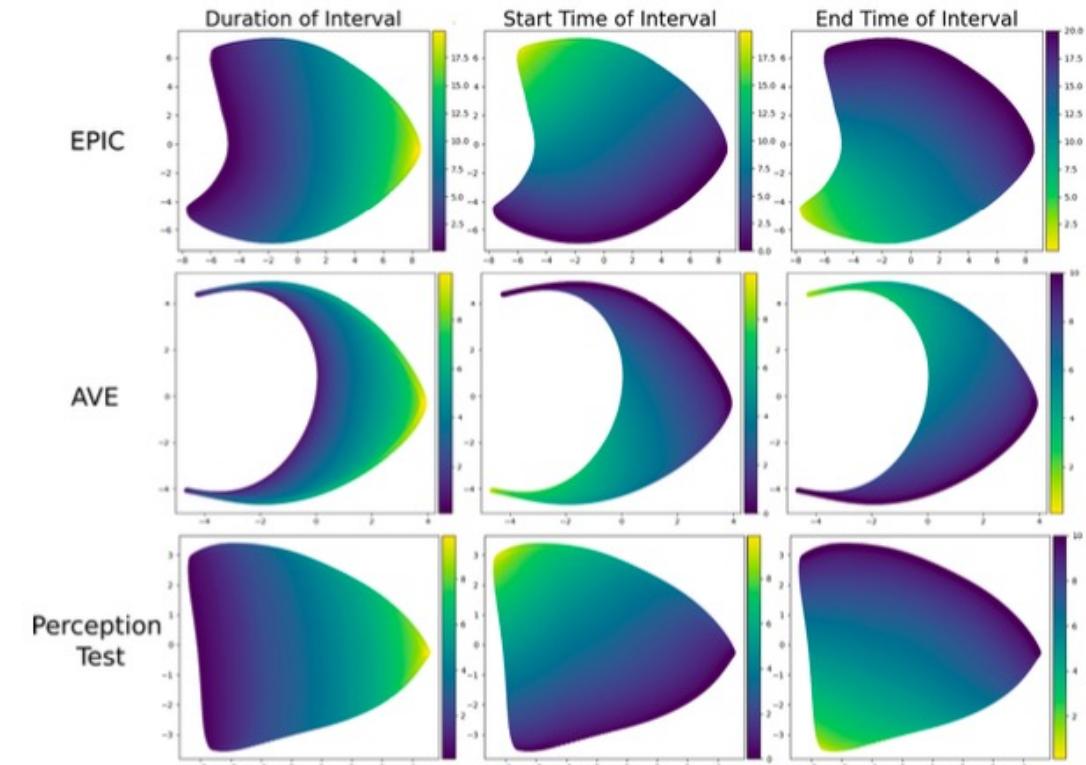
Model	<i>xp</i>	LLM	Verb	Noun	Action
<i>Visual-only models</i>					
MFormer-HR [37]	336p	✗	67.0	58.5	44.5
MoViNet-A6 [27]	320p	✗	72.2	57.3	47.7
MeMViT [55]	224p	✗	71.4	60.3	48.4
Omnivore [14]	224p	✗	69.5	61.7	49.9
MTV [59]	280p	✗	69.9	63.9	50.5
LaViLa (TSF-L) [63]	224p	✓	72.0	62.9	51.0
AVION (ViT-L) [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>

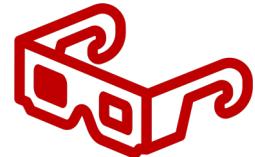
Perception Test Action				
Model	MLP (V)	MTCN [25](A+V)	TIM (V)	TIM (A+V)
<b>Top-1 acc</b>	43.7	51.2	56.1	<b>61.1</b>
Perception Test Sound				
Model	MLP (A)	MTCN [25](A+V)	TIM (A)	TIM (A+V)
<b>Top-1 acc</b>	50.6	52.9	54.8	<b>56.1</b>

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





Motivation and Datasets in  
Egocentric Video Understanding



Video Understanding  
Out of the Frame



Video Understanding:  
Data and Tasks



Teaser: The Wizard of Oz  
& Genie 3



Videos are Multimodal



Outlook into the Future of  
Egocentric Vision



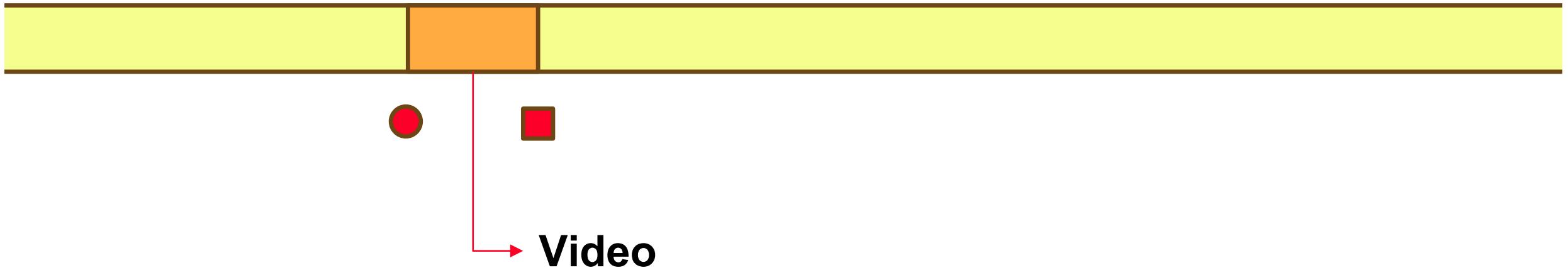
Connected Videos of One's Life



Conclusion

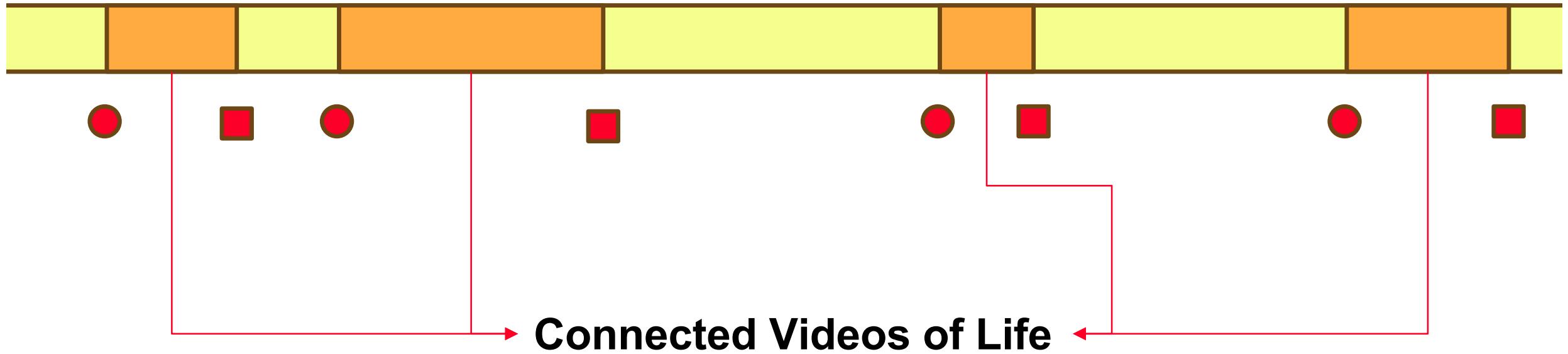


# Current Video Understanding...





# Upcoming Video Understanding...



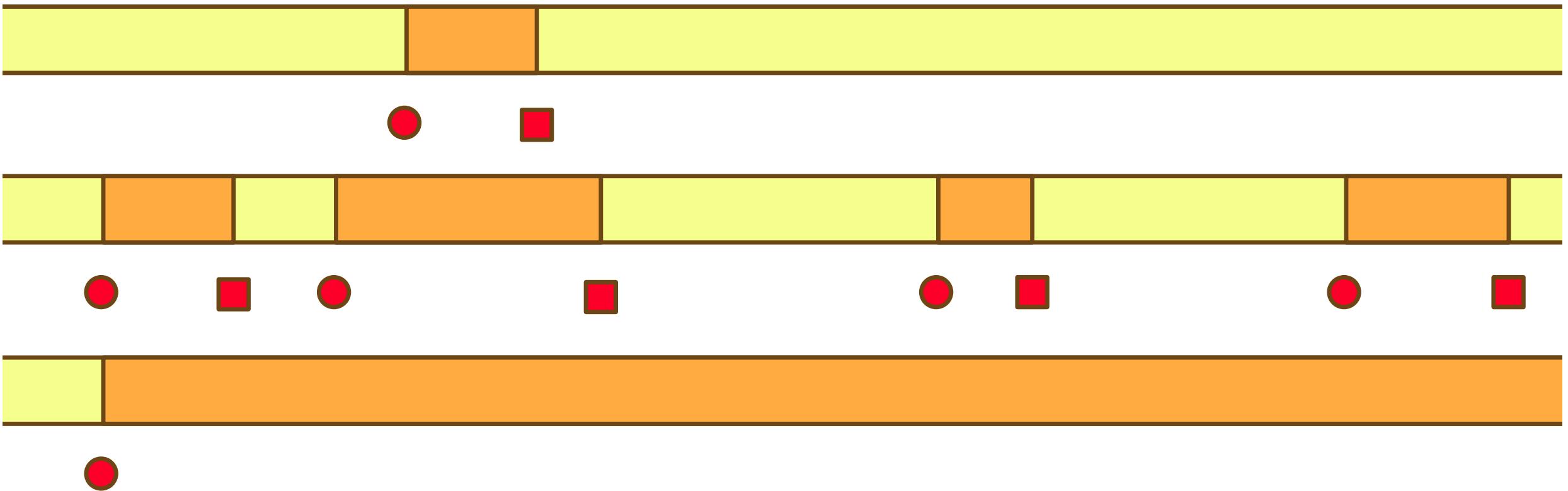


# Eventually...



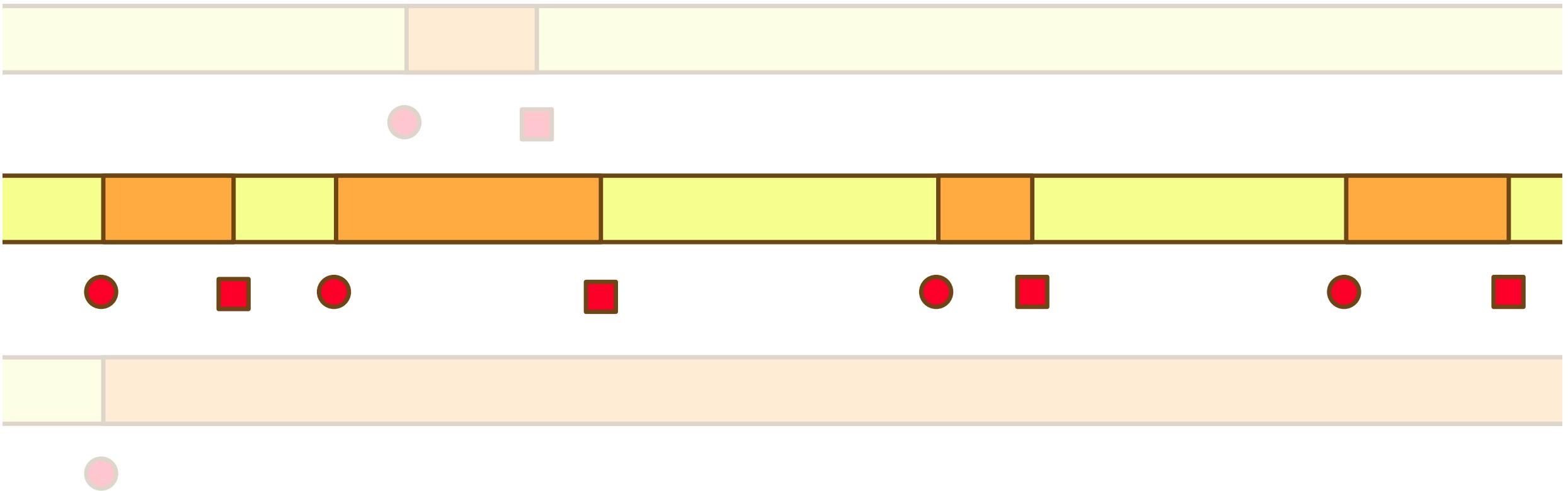


# Egocentric Video Understanding





# 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

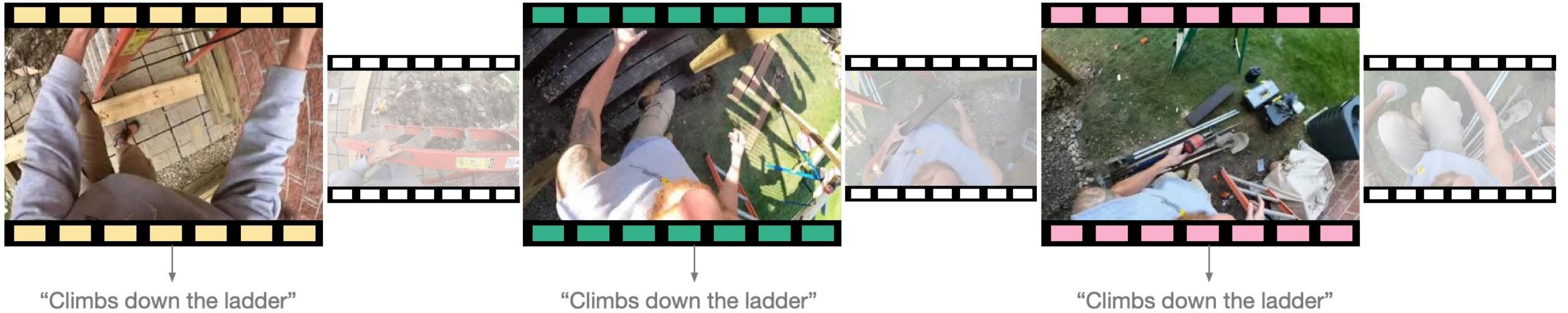
with: Toby Perrett  
Tengda Han  
Andrew Zisserman

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



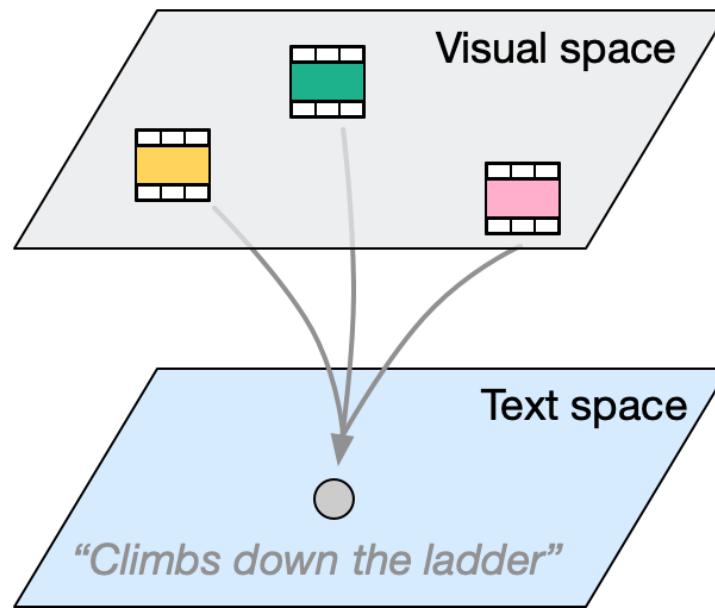
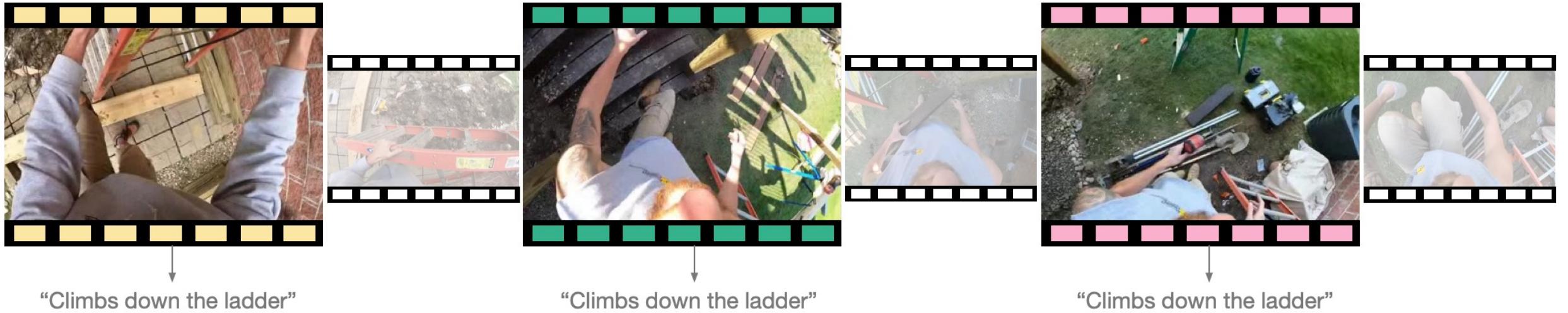
# Unique Video Captioning

with: Toby Perrett  
Tengda Han  
Andrew Zisserman

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

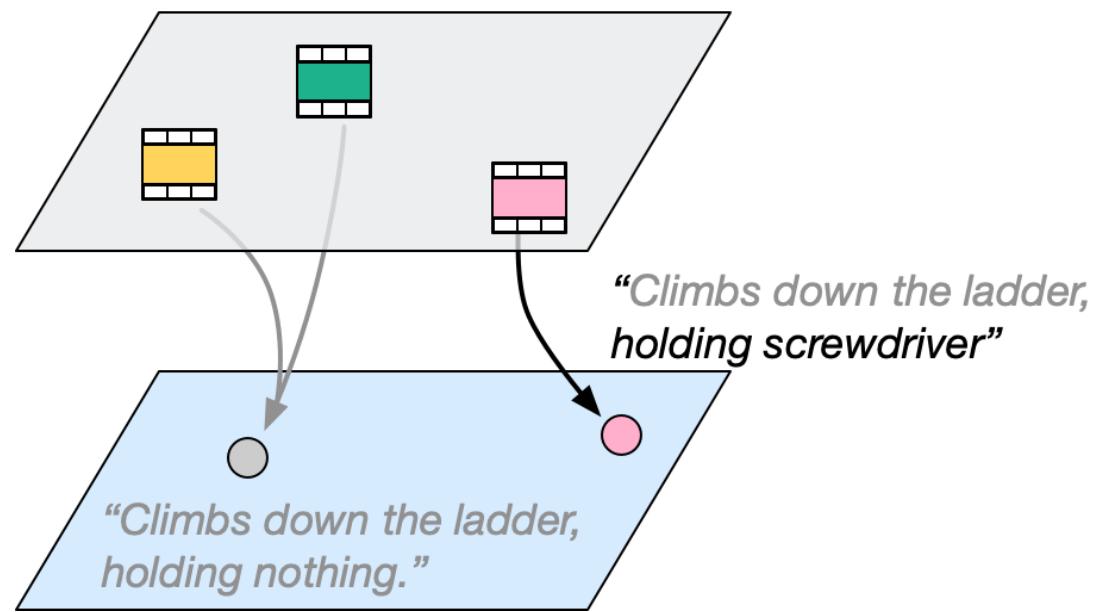
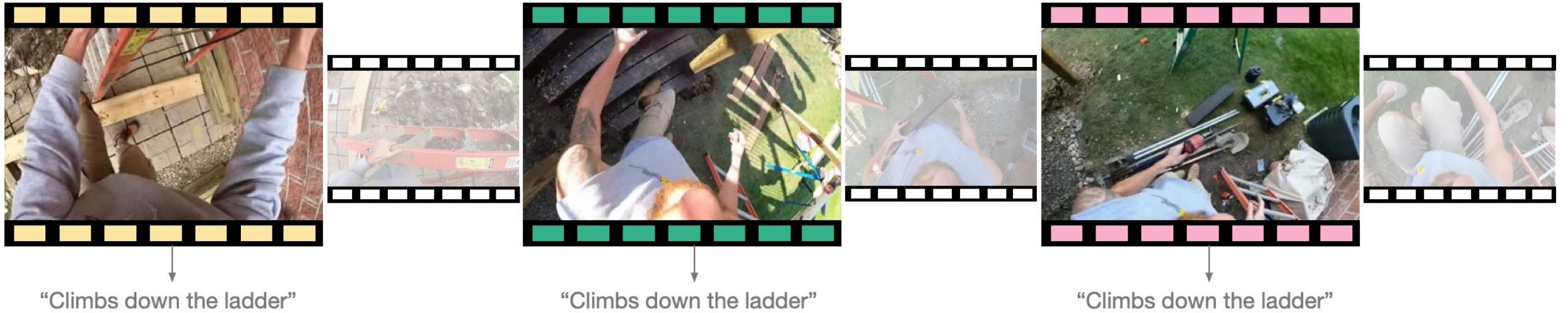
# Unique Video Captioning

with: Toby Perrett  
Tengda Han  
Andrew Zisserman



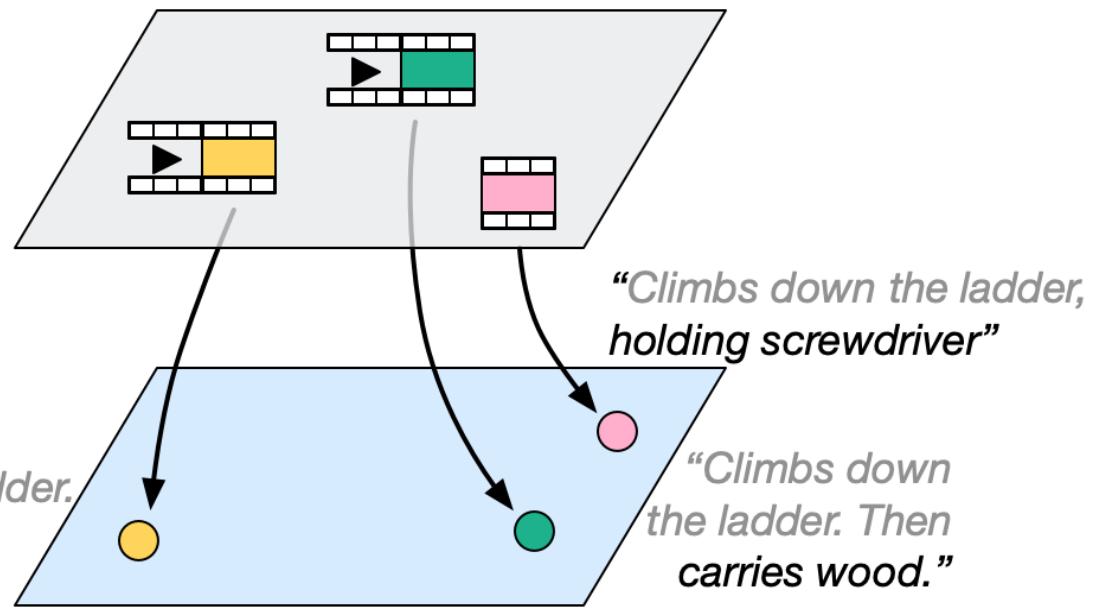
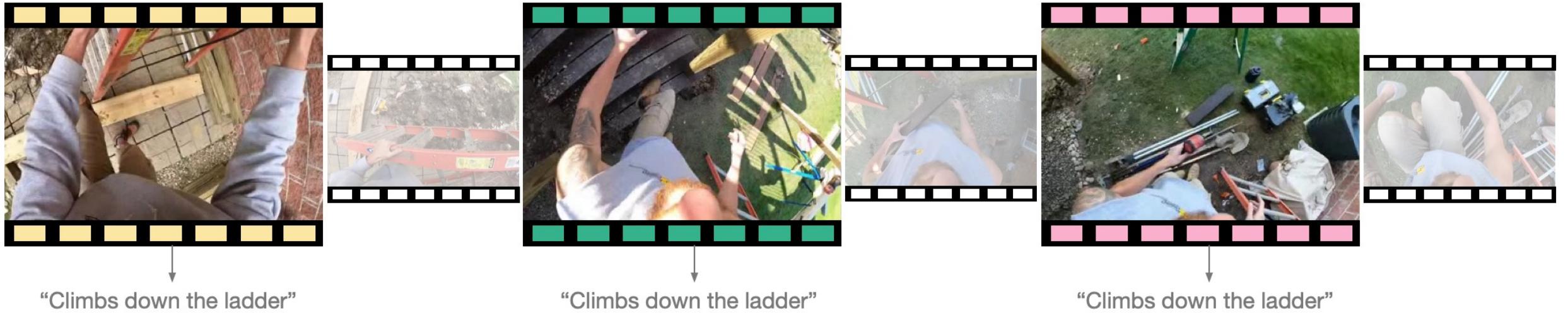
# Unique Video Captioning

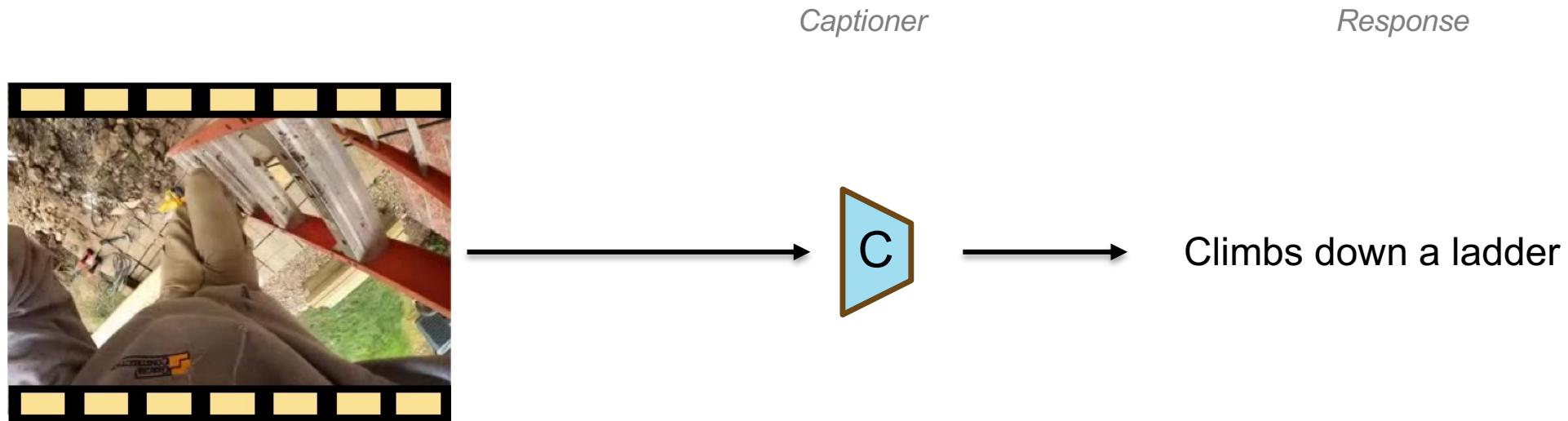
with: Toby Perrett  
Tengda Han  
Andrew Zisserman



# Unique Video Captioning

with: Toby Perrett  
Tengda Han  
Andrew Zisserman







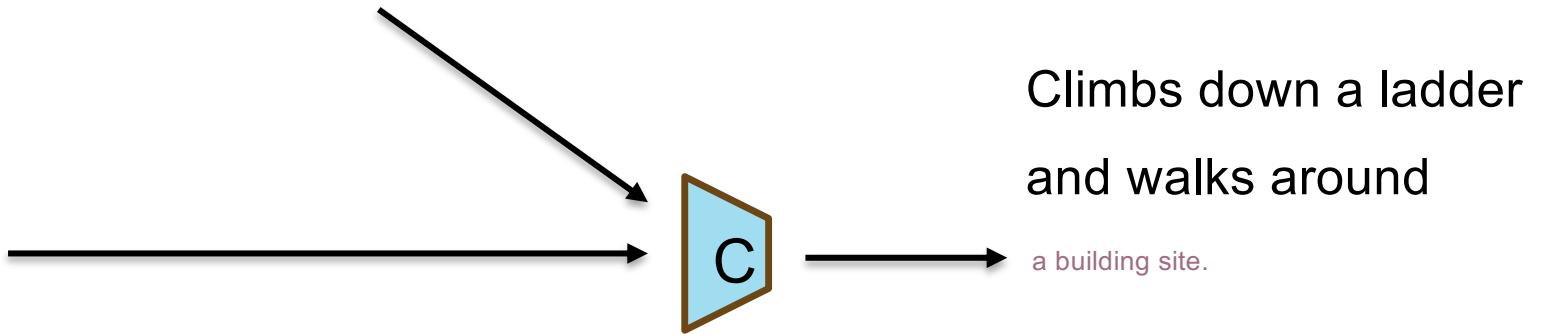
with: Toby Perrett  
Tengda Han  
Andrew Zisserman

*Discriminative prompt*

*Captioner*

*Response*

The person walks around

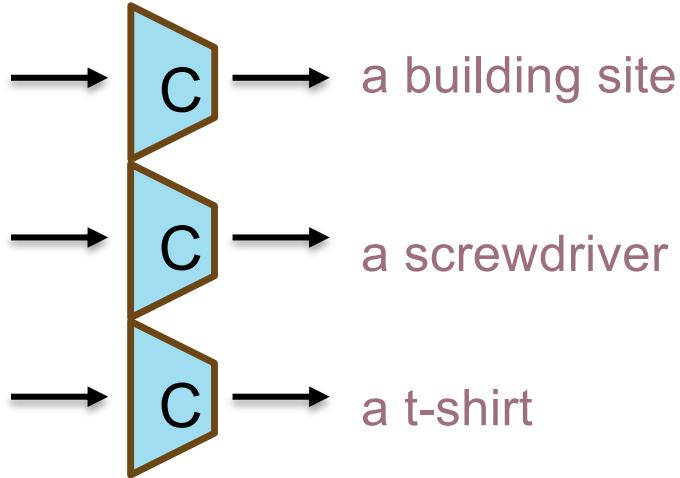


# Captioning by Discriminative Prompting

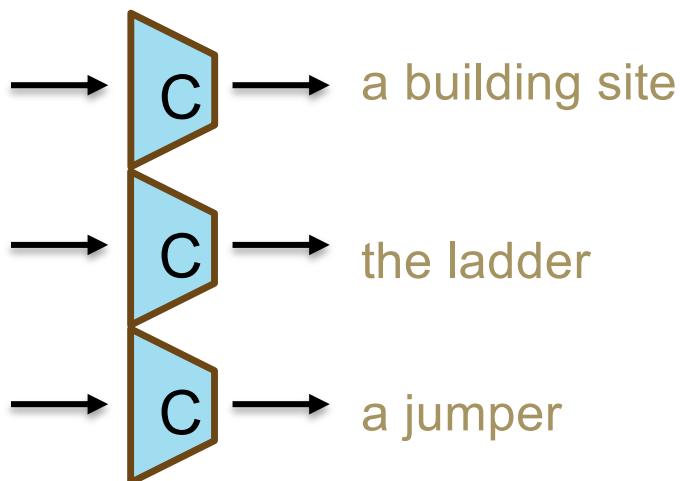
with: Toby Perrett  
Tengda Han  
Andrew Zisserman



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



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



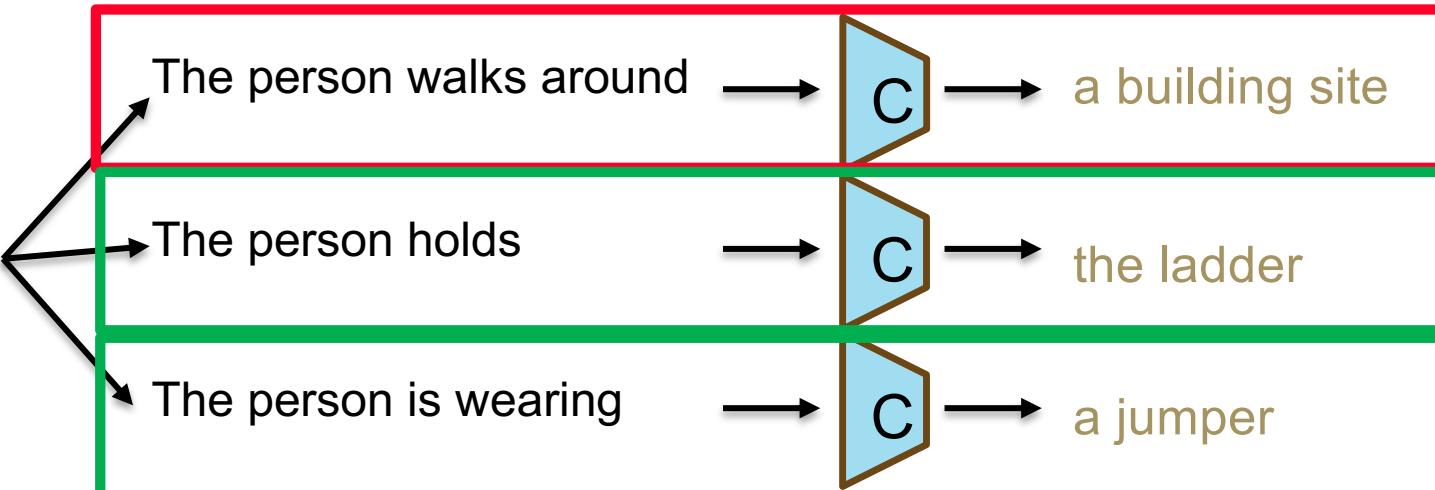
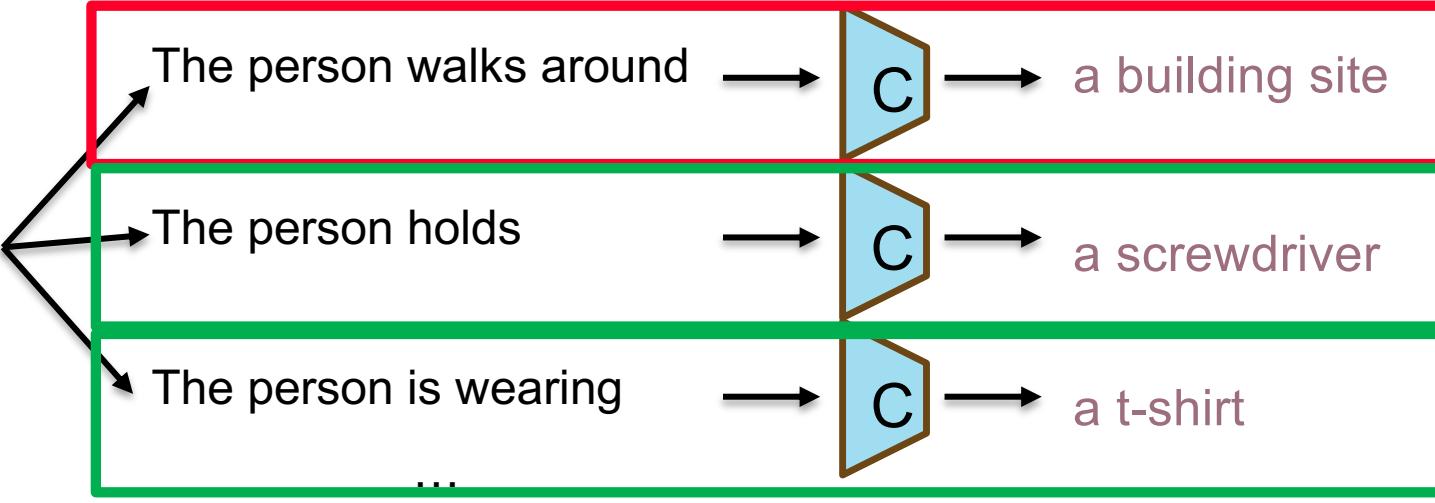
# Captioning by Discriminative Prompting

with: Toby Perrett  
Tengda Han  
Andrew Zisserman



*Discriminative prompts*

*Responses*





# Captioning by Discriminative Prompting

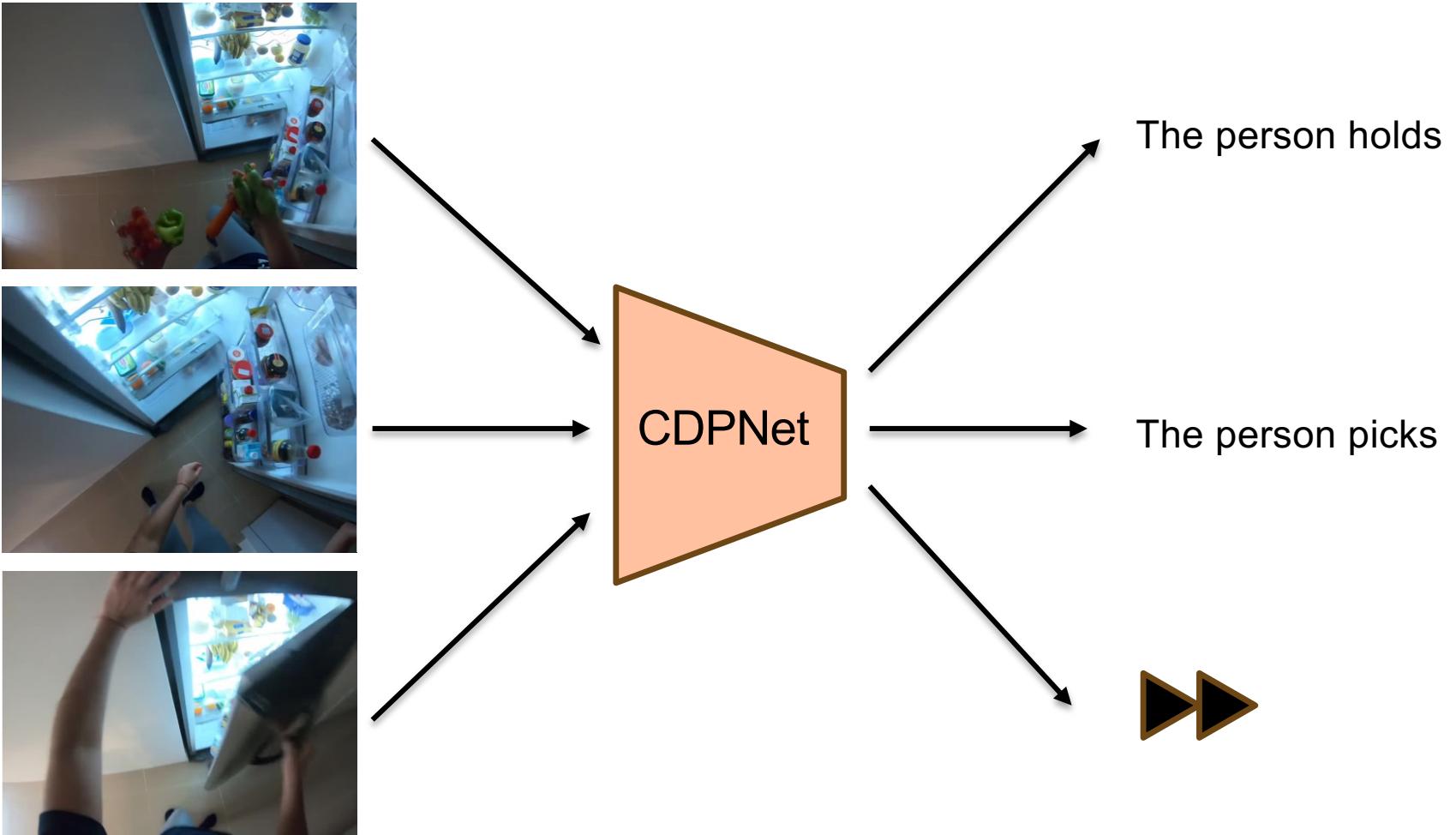
with: Toby Perrett  
Tengda Han  
Andrew Zisserman

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

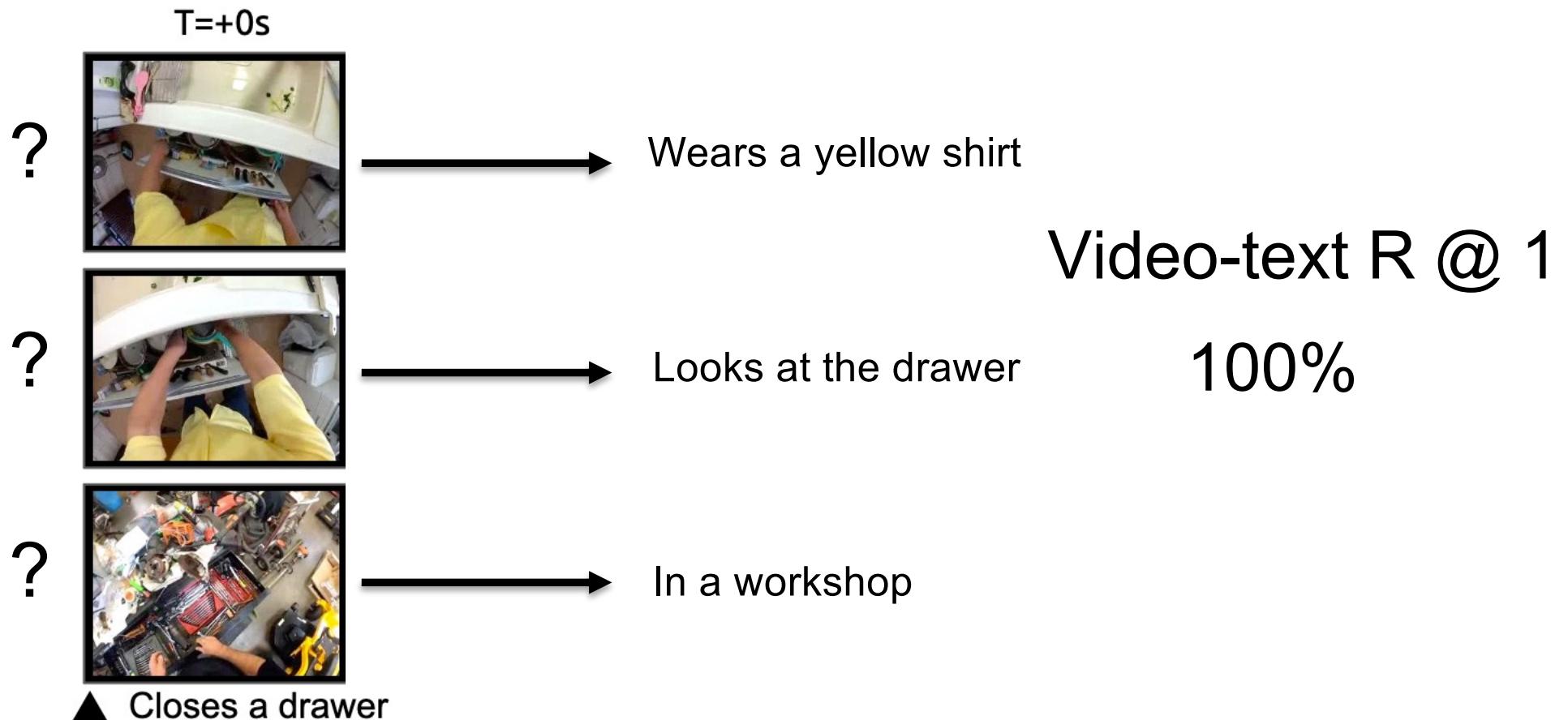
with: Toby Perrett  
Tengda Han  
Andrew Zisserman





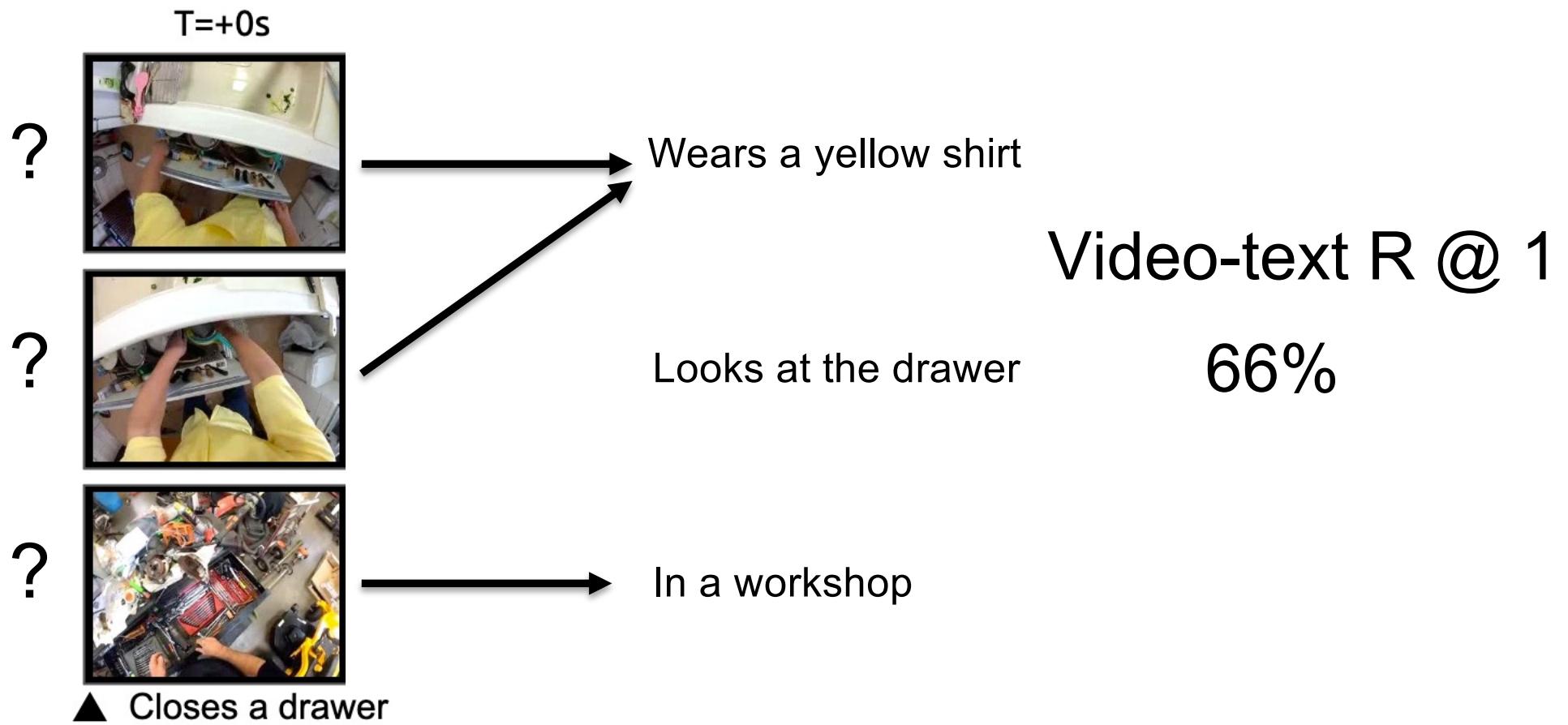
with: Toby Perrett  
Tengda Han  
Andrew Zisserman

Task: Caption every clip, then evaluate retrieval.





Task: Caption every clip, then evaluate retrieval.



# Benchmarks



with: Toby Perrett  
Tengda Han  
Andrew Zisserman

T=+0s



Wears a yellow shirt



Looks at the drawer



In a workshop

▲ Closes a drawer

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



# 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 [tobyperrett.github.io/its-just-another-day/](https://tobyperrett.github.io/its-just-another-day/). The page title is "It's Just Another Day: Unique Video Captioning by Discriminative Prompting". Below the title, it says "ACCV 2024 Oral". The authors listed are Toby Perrett, Tengda Han, Dima Damen, Andrew Zisserman. There are links for "arXiv | Code/Benchmark". The page content includes sections for "Introduction" and "Problem statement", followed by a visual example of captioning errors.

**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.

---

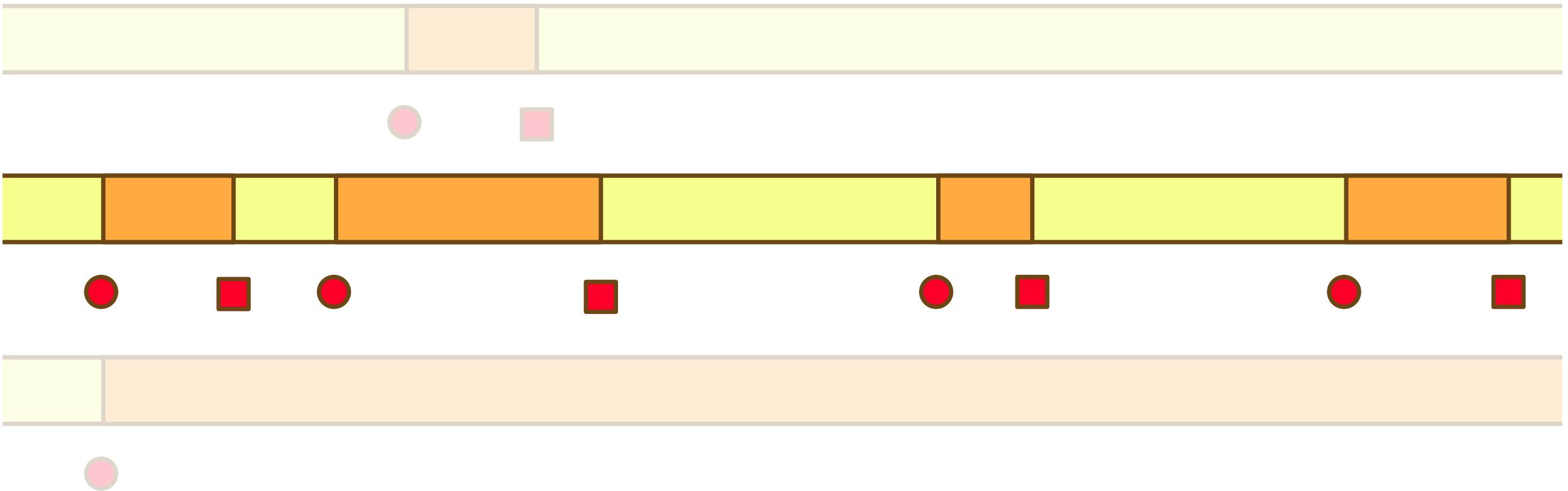
## Problem statement

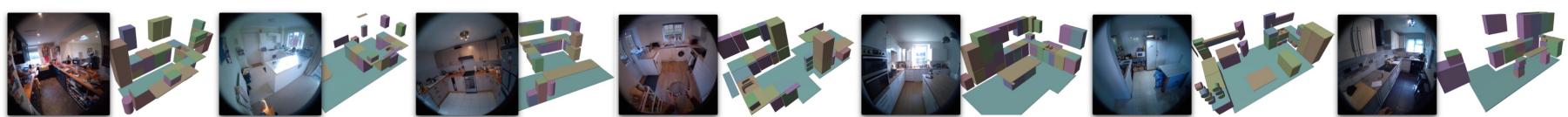
The following figures highlight a shortcoming of current captioning approaches. They caption each clip independently, giving similar captions for similar clips. First, in a timeloop movie:

"A man wakes up"  
"A man sits up"  
"A man wakes up"



# Egocentric Video Understanding





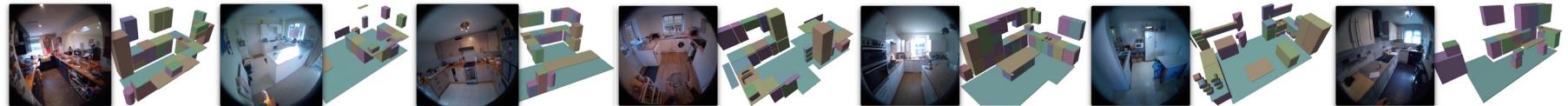
# HD-EPIC: A Highly-Detailed Egocentric Video Dataset



	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		

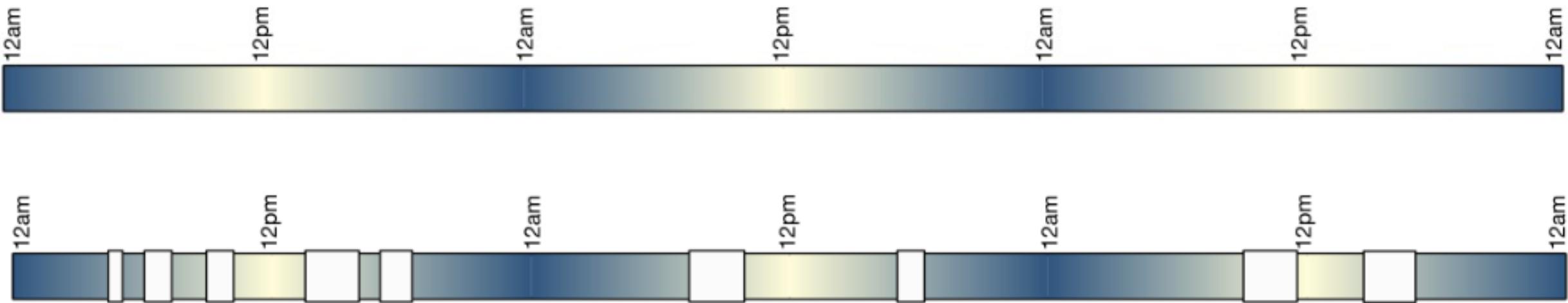
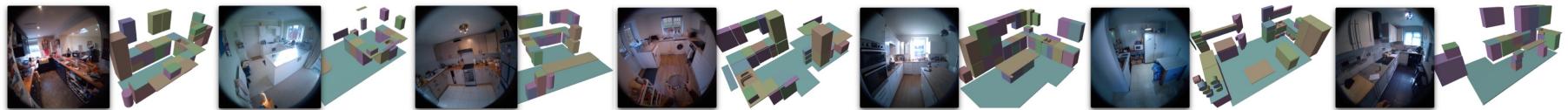


# HD-EPIC



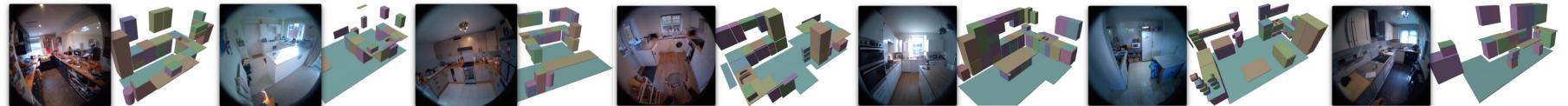


# HD-EPIC





# HD-EPIC



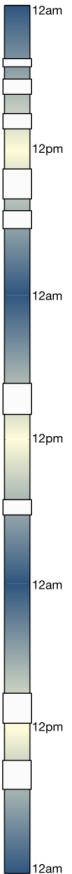
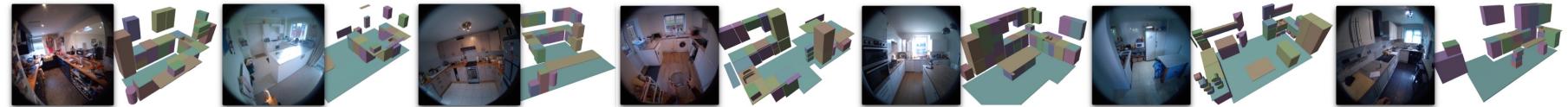
## Recorded over 3 days



Damen  
S 2025

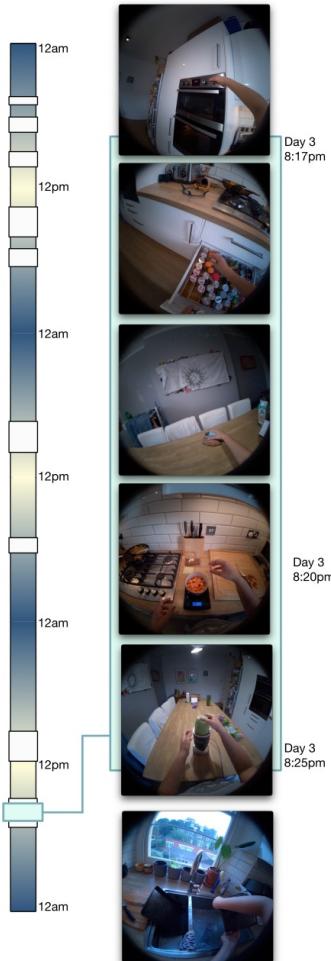
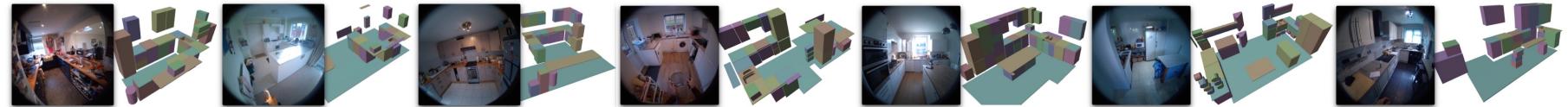


# HD-EPIC



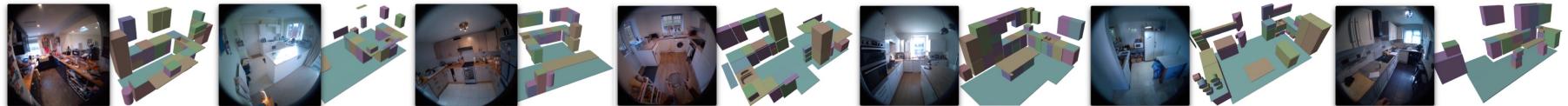


# HD-EPIC





# HD-EPIC



## Recipe: Southwestern Salad

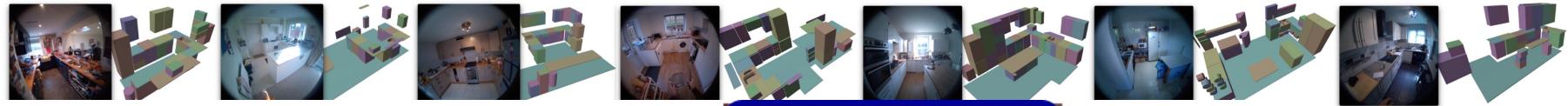
1: Preheat the oven to 400F

Day 3  
1:17pm

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**



## Cacio e Pepe (modified)

Ingredients:

~~200 g~~ penne

~~400g~~ of pasta of your choice

~~1~~ (we recommend bucatini)

~~2~~ tablespoon of black peppercorn

~~30 g~~ <sup>parmigiano</sup>

~~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 1-6 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.

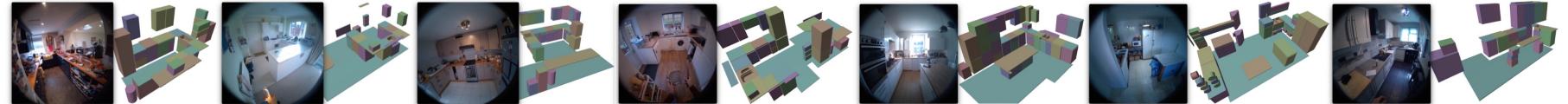


step 3





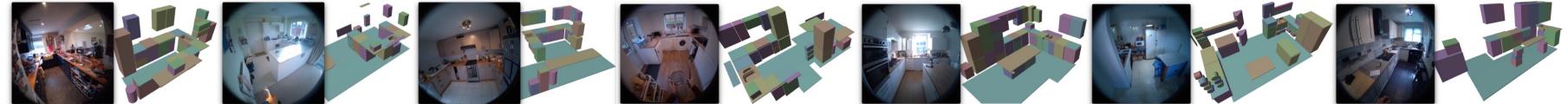
# HD-EPIC



- 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.**



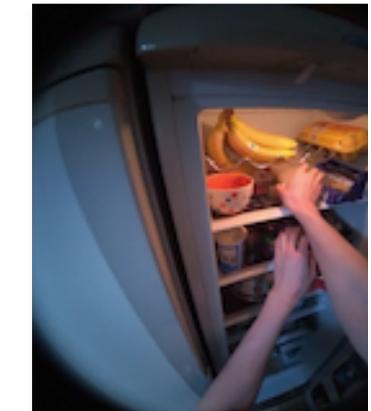
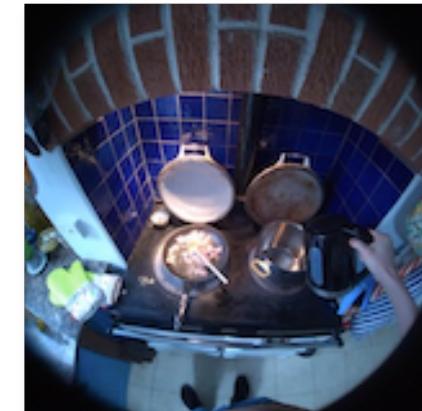
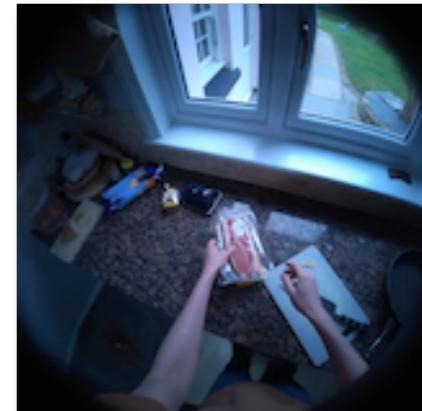
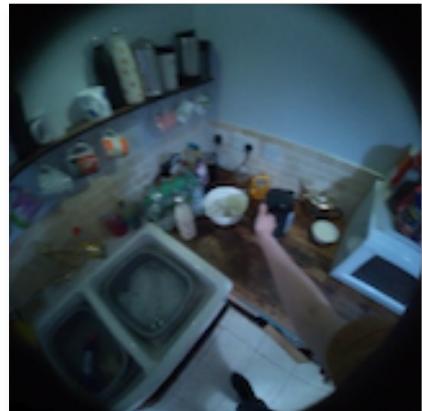
# HD-EPIC



- 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 ...

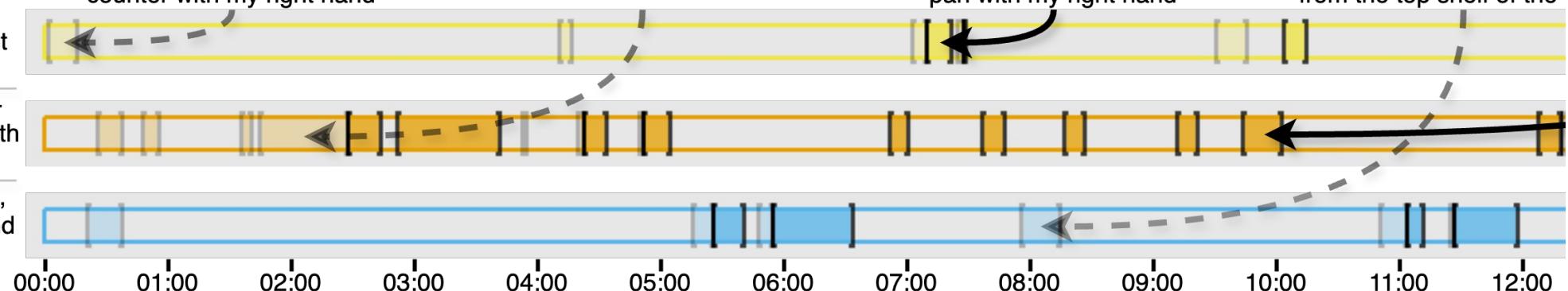
Meanwhile, beat the eggs in a bowl, then finely grate in the Parmesan and mix well.

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

pick up packet of bacon

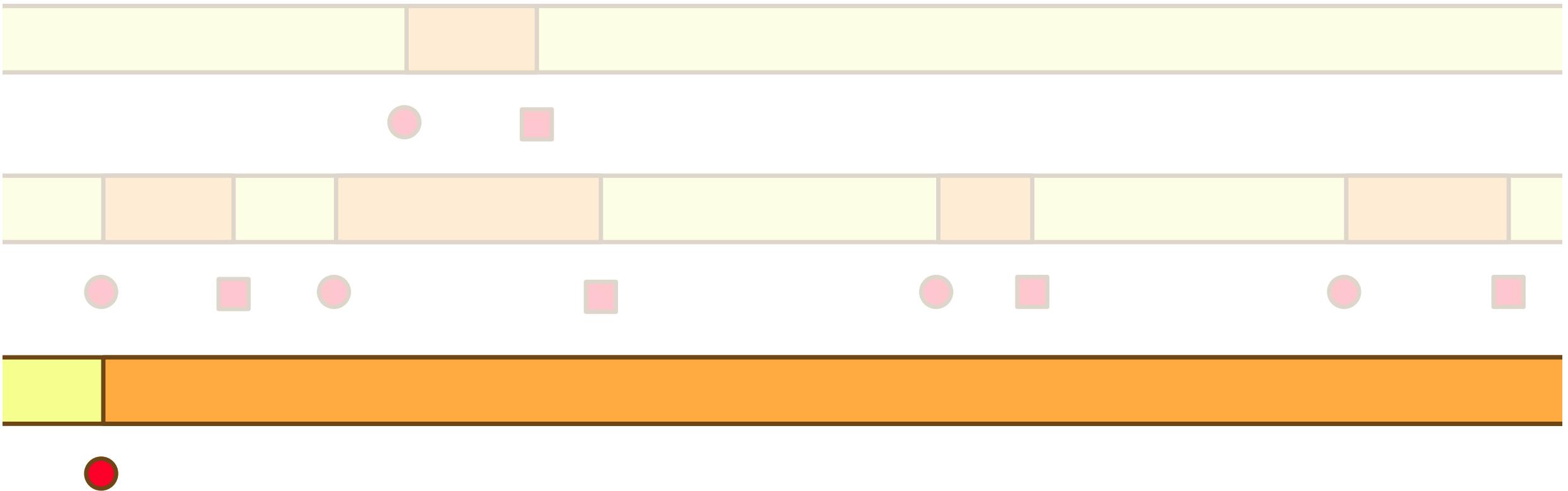
pour water from kettle into the pan with my right hand

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





# 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**

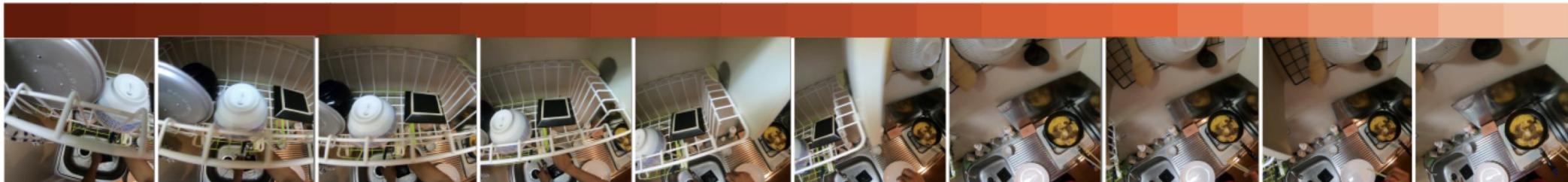
Tengda Han<sup>◊</sup>, Dilara Gokay<sup>◊</sup>, Joseph Heyward<sup>◊</sup>, Chuhan Zhang<sup>◊</sup>  
Daniel Zoran<sup>◊</sup>, Viorica Pătrăucean<sup>◊</sup>, João Carreira<sup>◊</sup>, Dima Damen<sup>◊†</sup>, Andrew Zisserman<sup>◊‡</sup>  
◊Google DeepMind, †University of Bristol, ‡University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



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

Dima Damen  
PAISS 2025

# Learning from Streaming Videos with Orthogonal Gradients

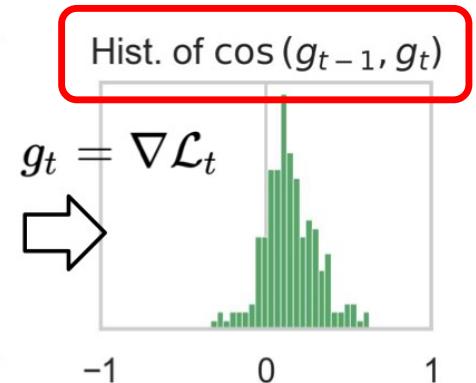


standard IID samples →



# Learning from Streaming Videos with Orthogonal Gradients

⤒  
shuffled  
loading



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

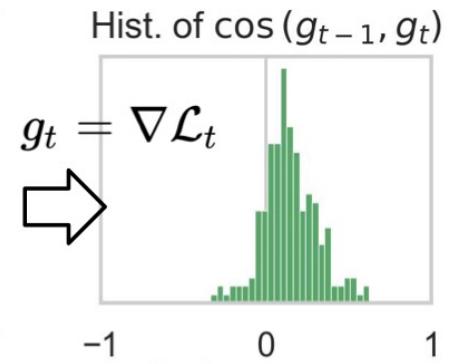


Han et al (2025). Learning from Streaming Video with Orthogonal Gradients. IEEE/CVF Computer Vision and Pattern Recognition (CVPR)

Dima Damen  
PAISS 2025

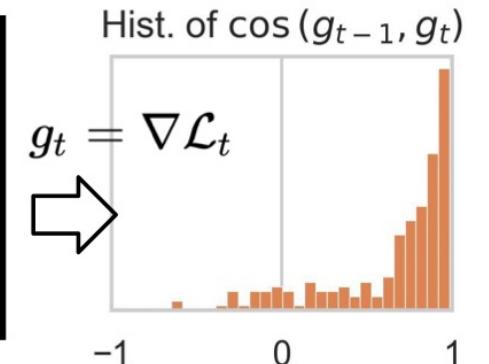
# Learning from Streaming Videos with Orthogonal Gradients

⤒  
shuffled  
loading



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

⌚  
sequential  
loading



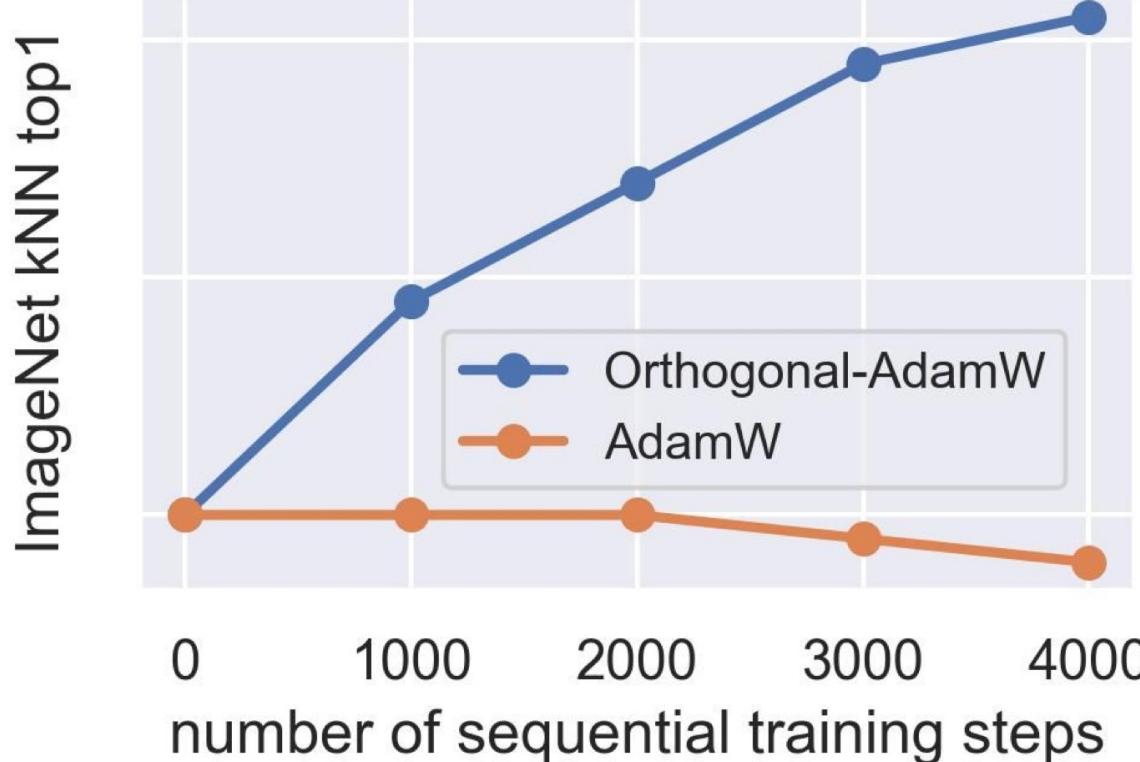
gradients are **highly correlated** over training steps



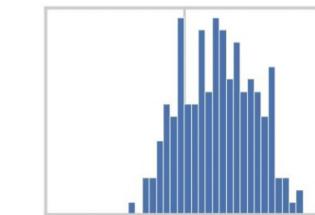
Han et al (2025). Learning from Streaming Video with Orthogonal Gradients. IEEE/CVF Computer Vision and Pattern Recognition (CVPR)

Dima Damen  
PAISSL 2025

# Learning from Streaming Videos with Orthogonal Gradients

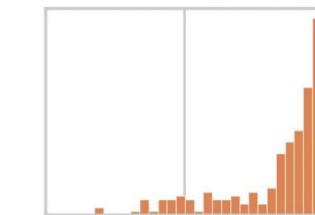


Hist. of  $\cos(g_{t-1}, g_t)$



Orthogonal Optimizer

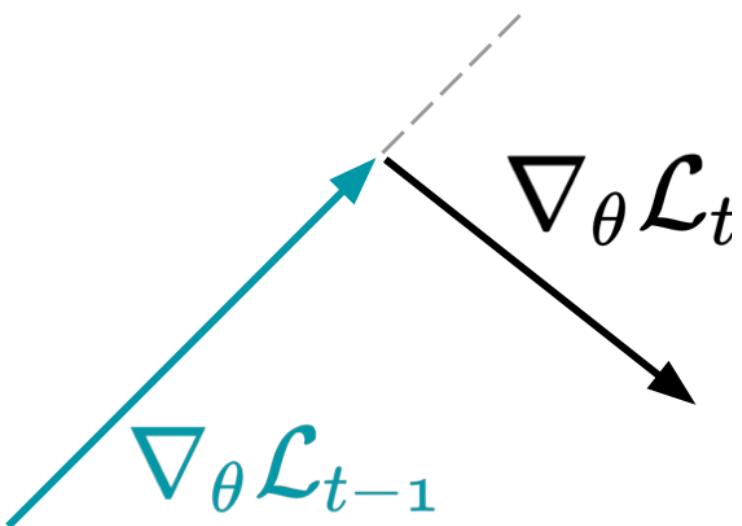
Hist. of  $\cos(g_{t-1}, g_t)$



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

Dima Damen  
PAISS 2025

# 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  
PAISS 2025

# 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:  $g_t = \nabla_{\theta} \mathcal{L}(\theta_{t-1}; \mathcal{B}_t)$

8:   Update biased first moment estimate:  $m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$

9:   Update biased second moment estimate:  $v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$

10:   Compute bias-corrected first moment:  $\hat{m}_t = \frac{m_t}{1 - \beta_1^t}$

11:   Compute bias-corrected second moment:  $\hat{v}_t = \frac{v_t}{1 - \beta_2^t}$

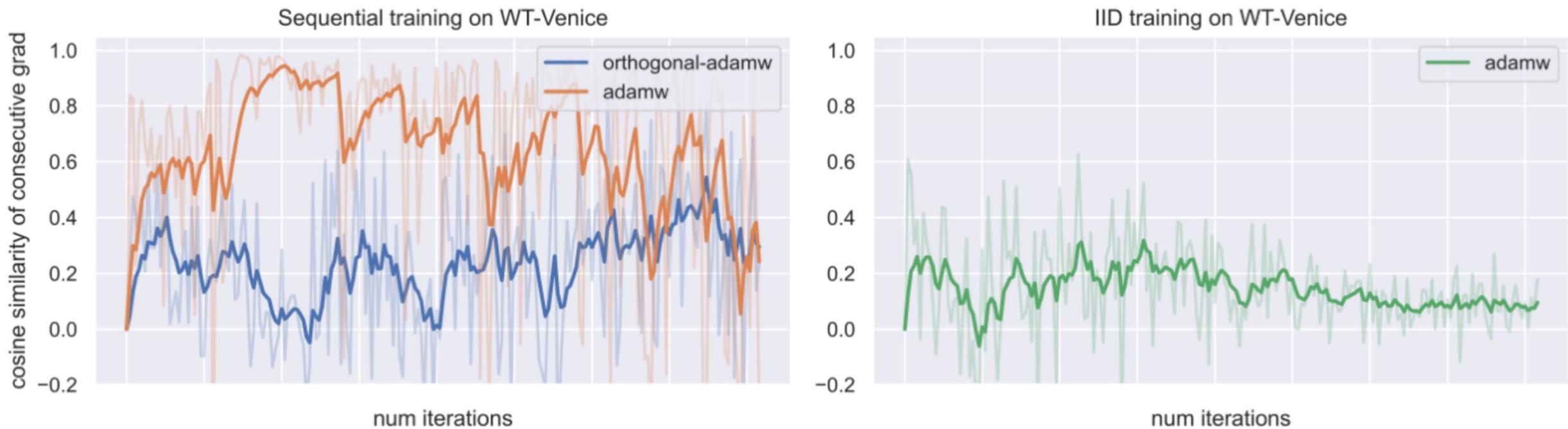
12:   Apply weight decay:  $\theta_{t-1} = \theta_{t-1} - \eta \lambda \theta_{t-1}$

13:   Update parameters:  $\theta_t = \theta_{t-1} - \eta \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon}$

14: **end for**



# Learning from Streaming Videos with Orthogonal Gradients

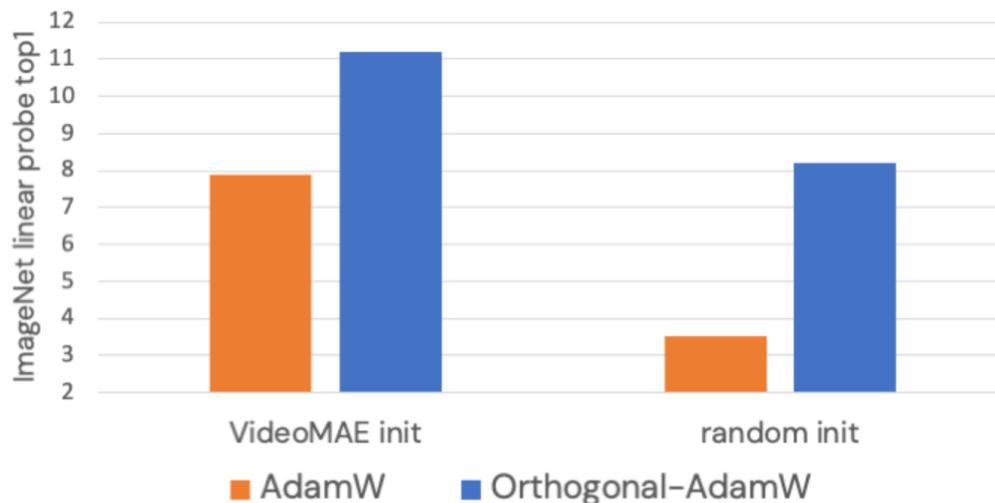


Han et al (2025). Learning from Streaming Video with Orthogonal Gradients. IEEE/CVF Computer Vision and Pattern Recognition (CVPR)

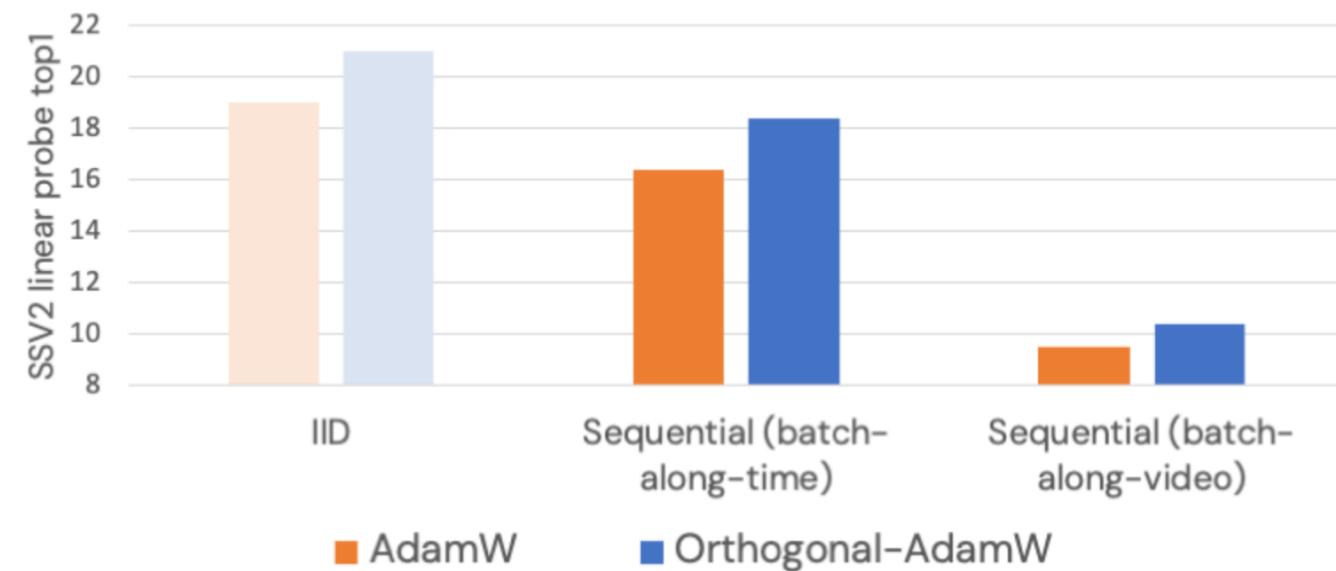
Dima Damen  
PAISS 2025

# Learning from Streaming Videos with Orthogonal Gradients

DoRA sequential pretraining on  
WalkingTour-Venice

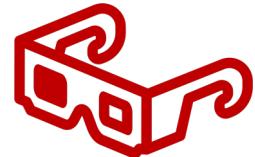


VideoMAE pretraining on SSV2



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

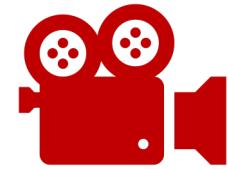
Dima Damen  
PAISS 2025



Motivation and Datasets in  
Egocentric Video Understanding



Video Understanding  
Out of the Frame



Video Understanding:  
Data and Tasks



Teaser: The Wizard of Oz  
& Genie 3



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**



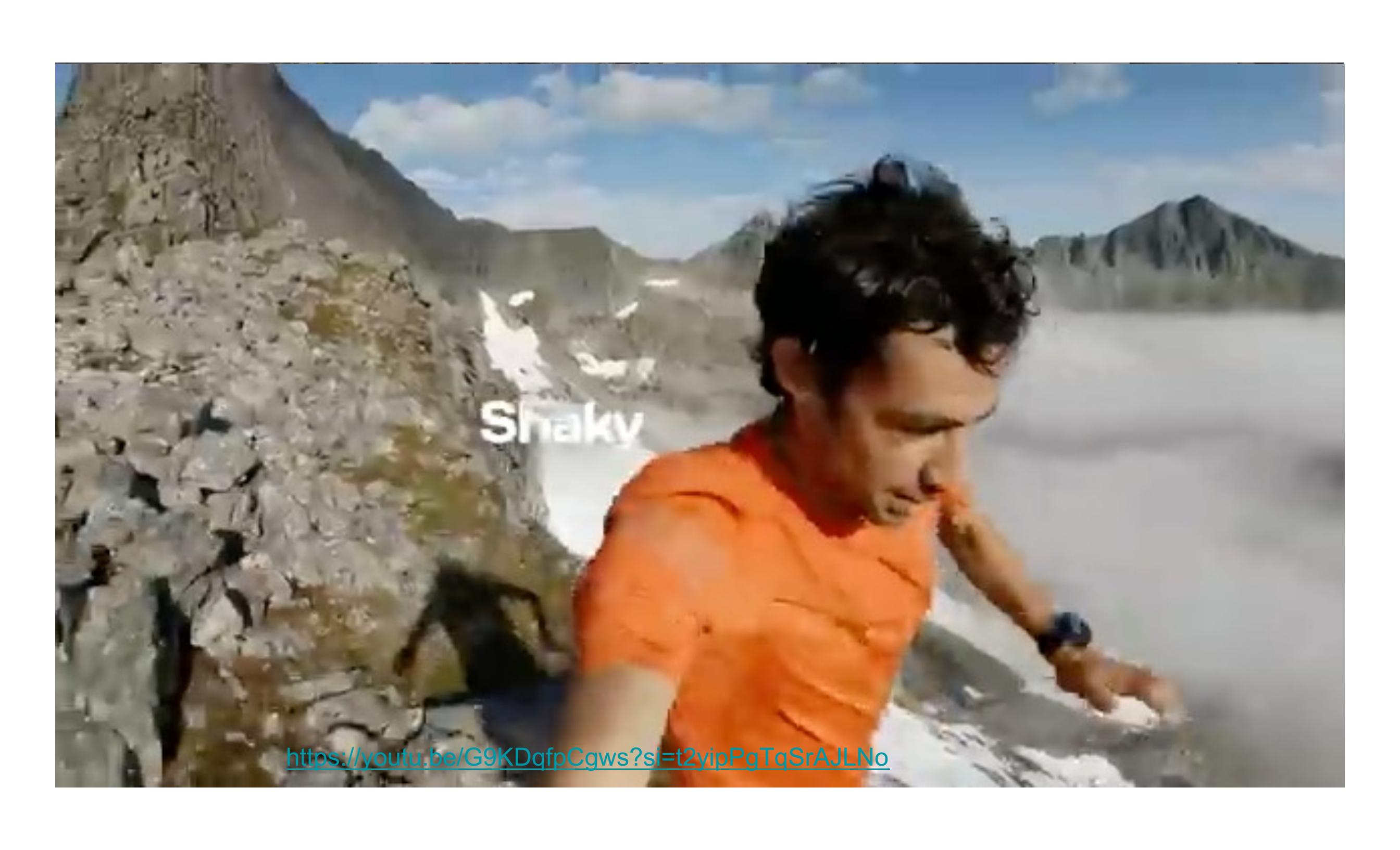
# 4 years later (2018)

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](#).

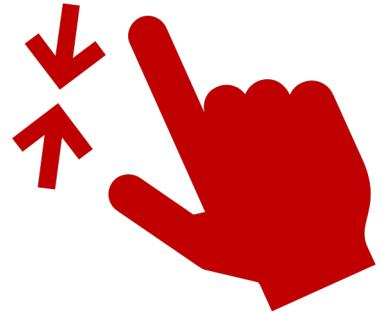
The screenshot shows the GoPro website homepage. At the top, there are navigation links for News, Awards, Support, and Gift Cards. The main headline is "NEW HERO13 Black in Forest Green". Below the headline, there's a message: "Looks like you're in United Kingdom (GB). Choose the country you want to shop in." The menu includes Cameras, Apps, Accessories, Lifestyle Gear, GoPro Subscription, and Shop by Activity. Below the menu, there are links for ALL NEWS, PROS, IN THE WILD, TOOLS OF THE TRADE, HEROES, COMMUNITY, and EVENTS. The main content area features a large image of a roller coaster track with the text "SHAKY VIDEO IS DEAD: HERO7 BLACK IS HERE". Below this, there's a video player showing a GoPro video titled "GoPro: Introducing HERO7 Black in 4K - Shaky Video is Dead". The video player includes controls for "More videos", "HERO 7 16 NEW THINGS TO KNOW", "HyperSmooth", and "Previous GoPro". The video duration is 0:29 / 2:00. The GoPro logo is visible in the top left corner of the video player.





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

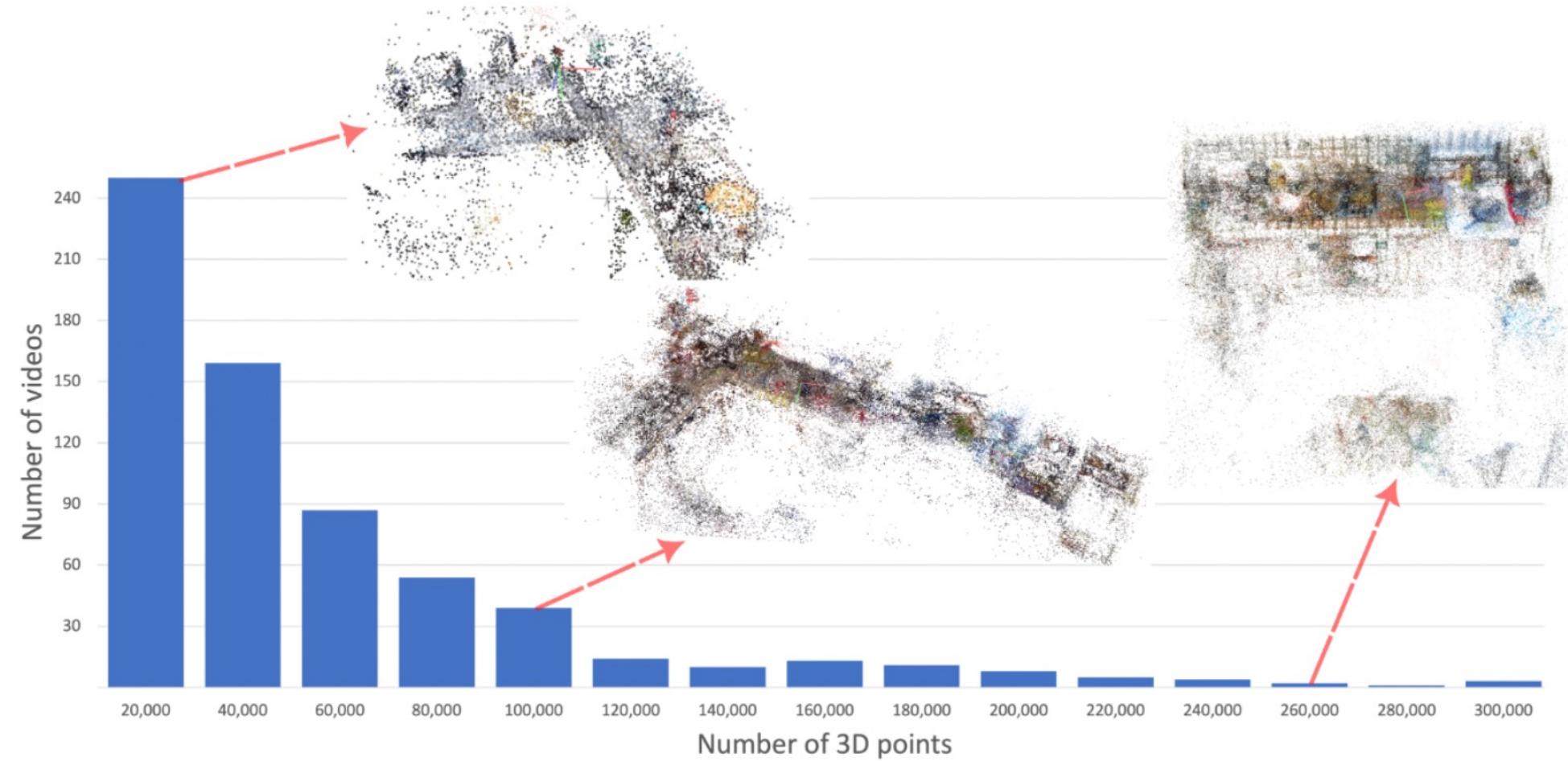


**EPIC-KITCHENS**

# 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



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





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



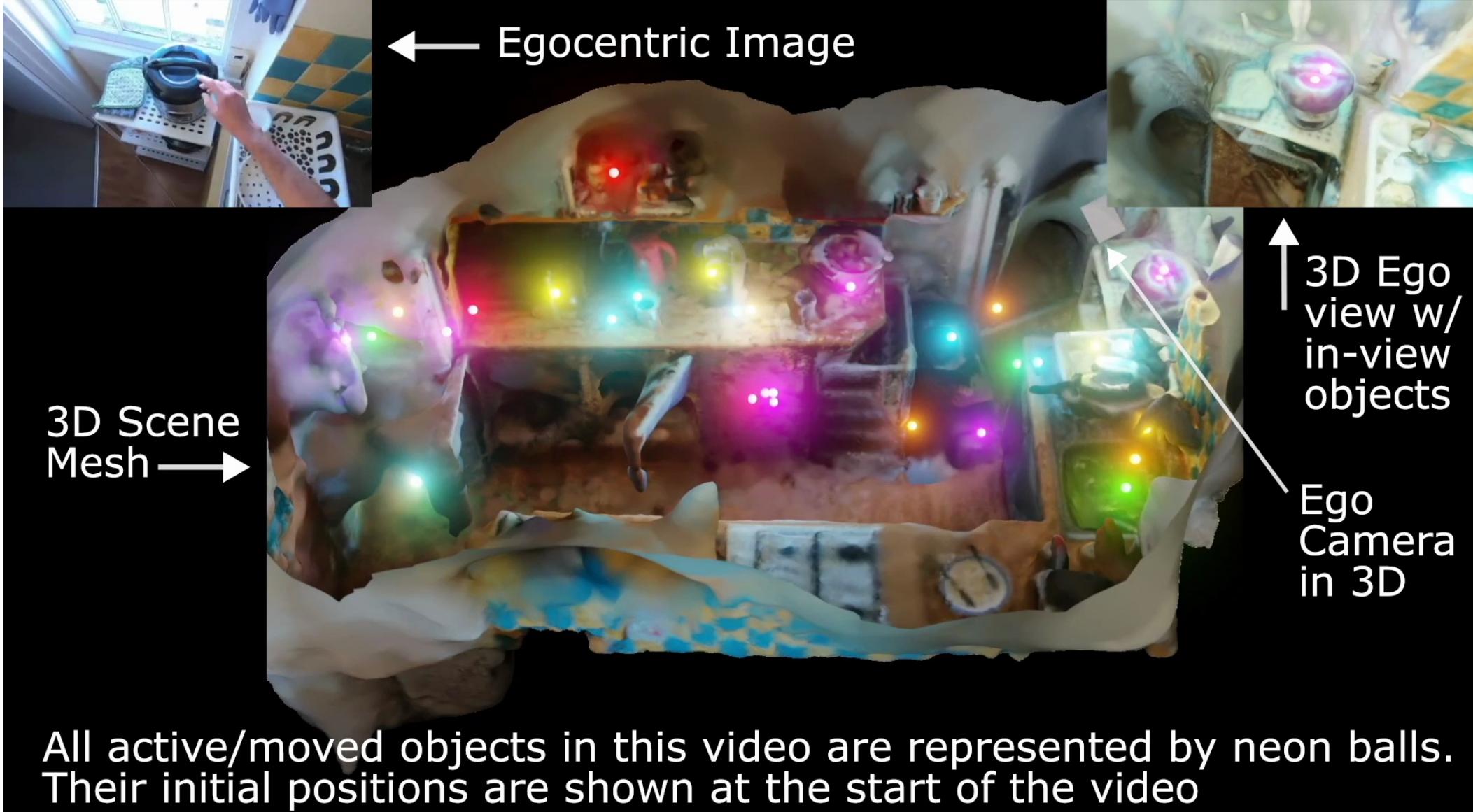
Plizzari et al (2025). Spatial Cognition from Egocentric Video: Out of Sight, Not Out of Mind. 3DV

...na Damen  
PAISS 2025

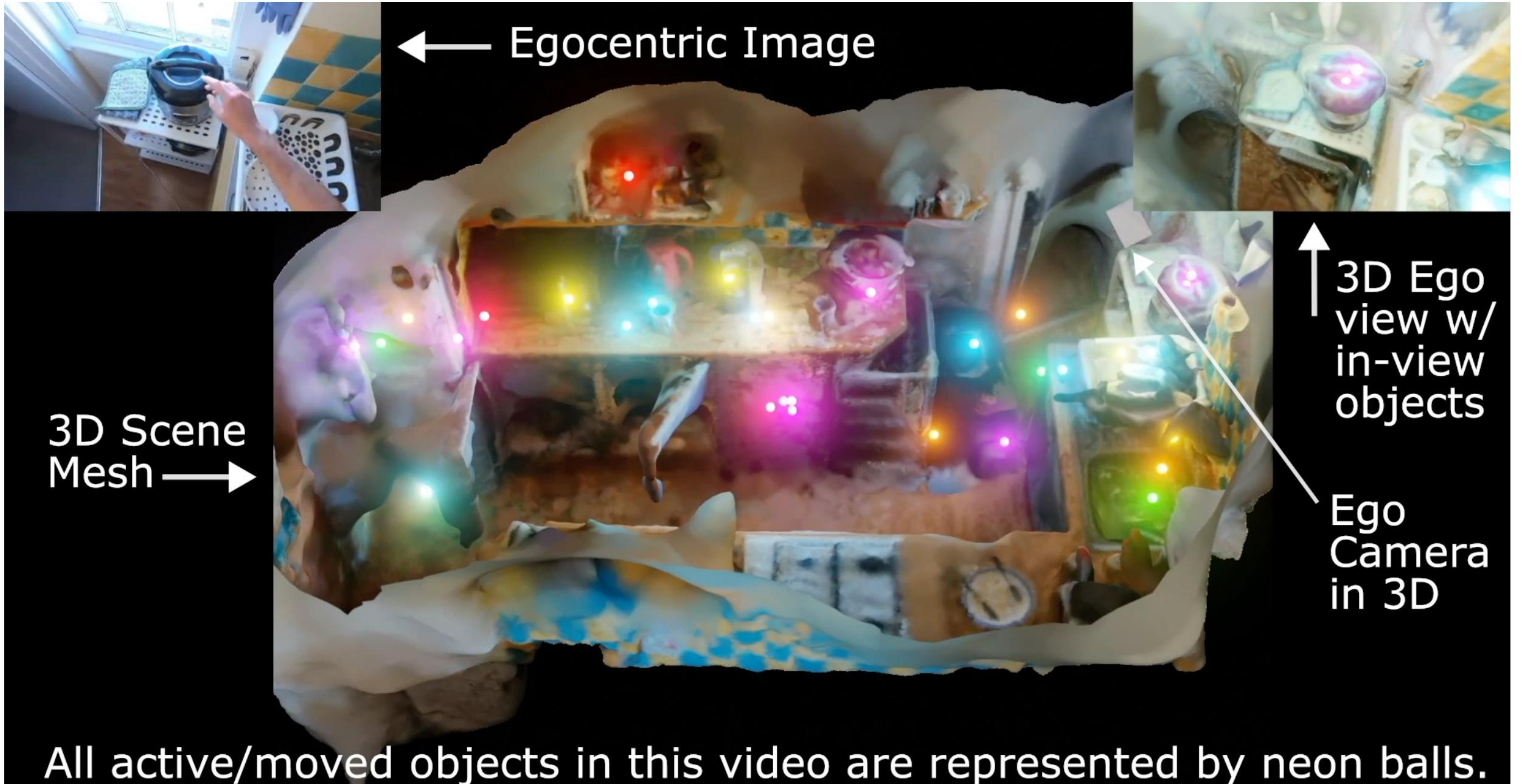


with: Chiara Plizzari  
Toby Perrett

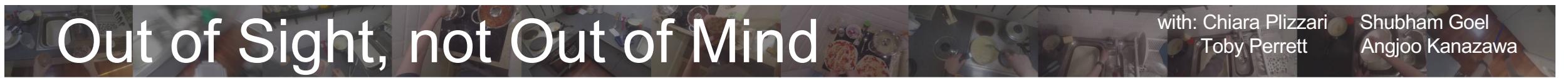
Shubham Goel  
Angjoo Kanazawa



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

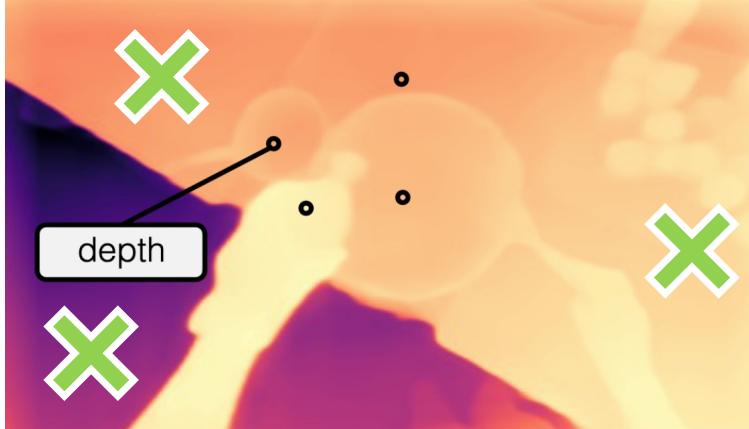
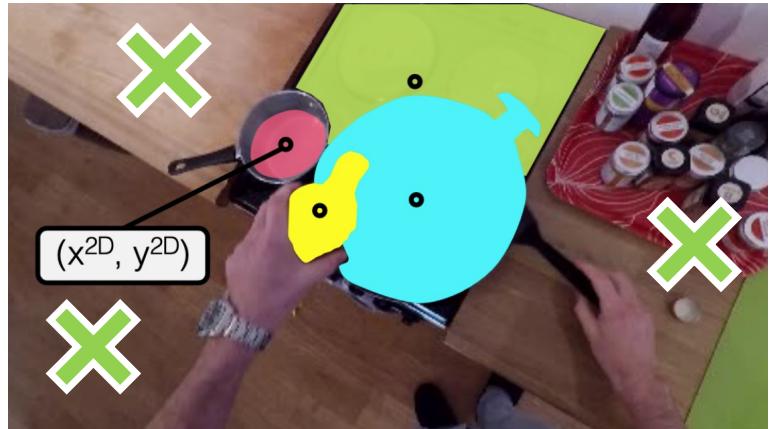
with: Chiara Plizzari  
Toby Perrett

Shubham Goel  
Angjoo Kanazawa

Lift

Match

Keep



0.0 ... 1.0

0.3m ... 1.8m

# Out of Sight, not Out of Mind

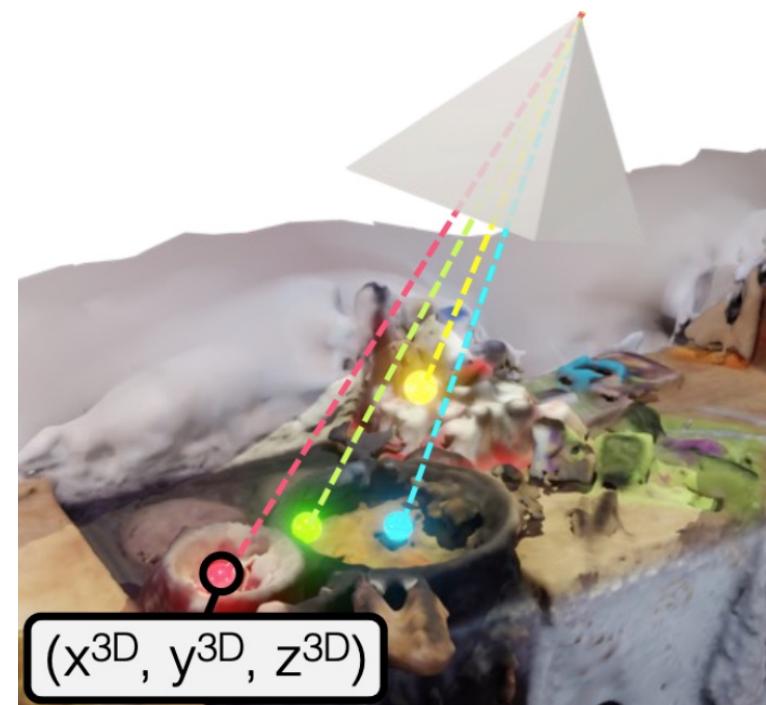
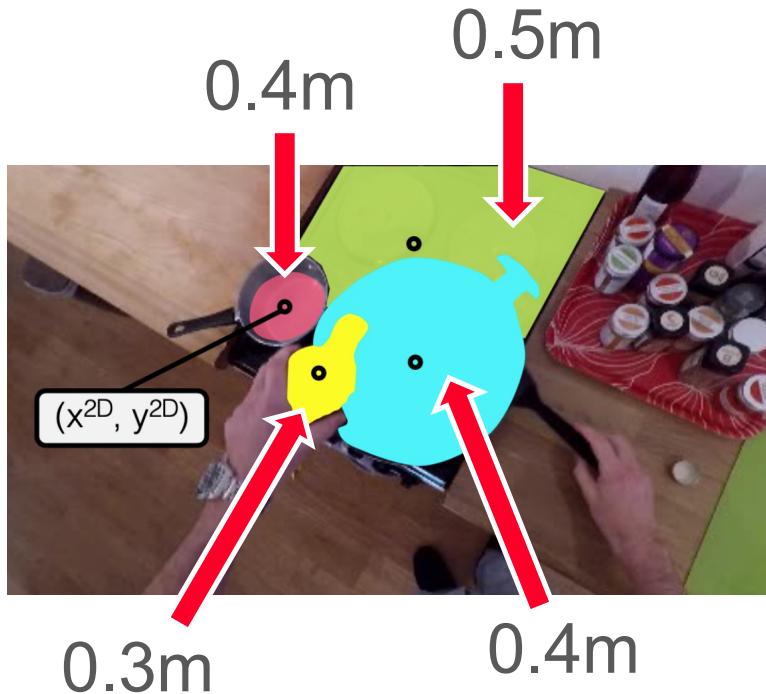
with: Chiara Plizzari  
Toby Perrett

Shubham Goel  
Angjoo Kanazawa

Lift

Match

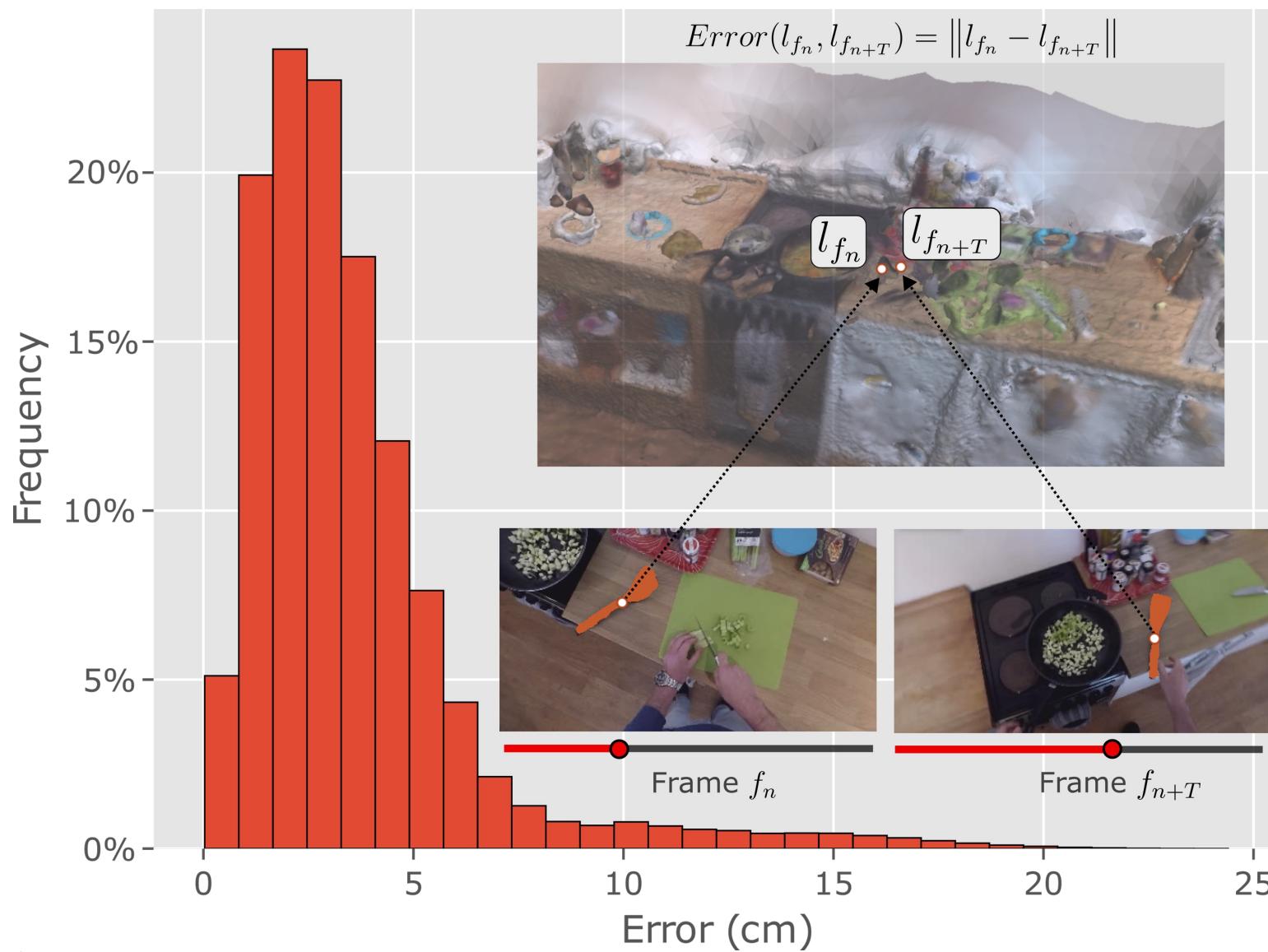
Keep



# 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

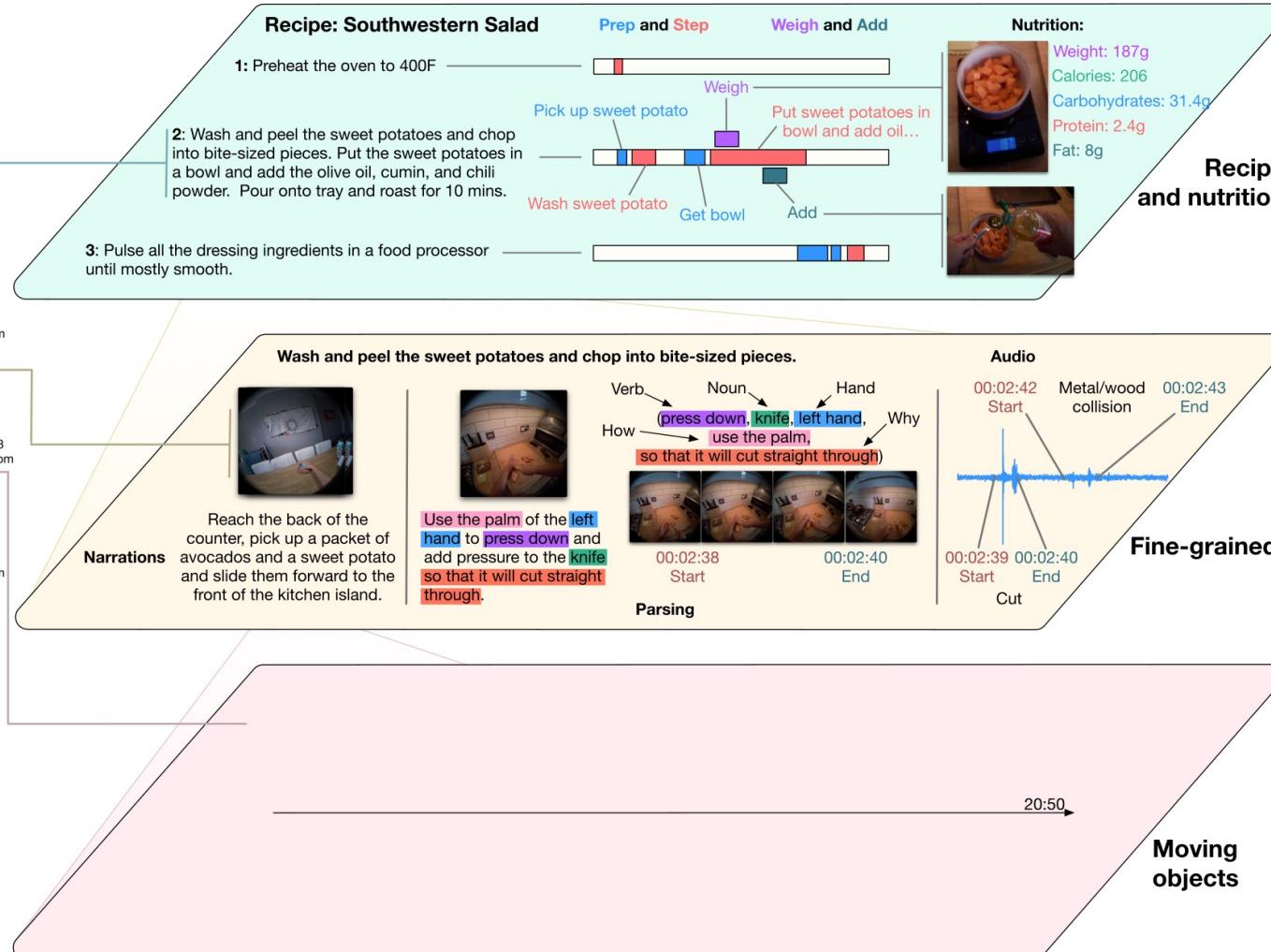
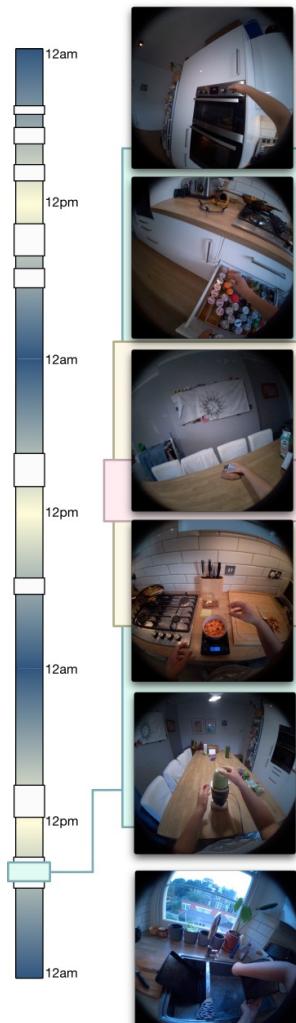
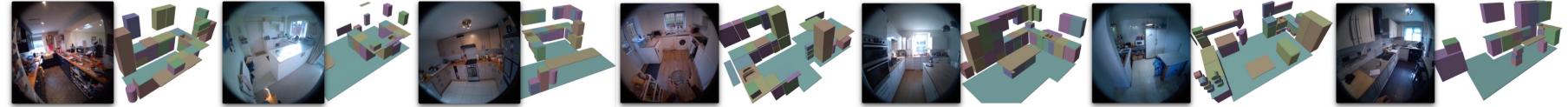


Plizzari et al (2025). Spatial Cognition from Egocentric Video: Out of Sight, Not Out of Mind. 3DV

...na Damen  
PAISS 2025

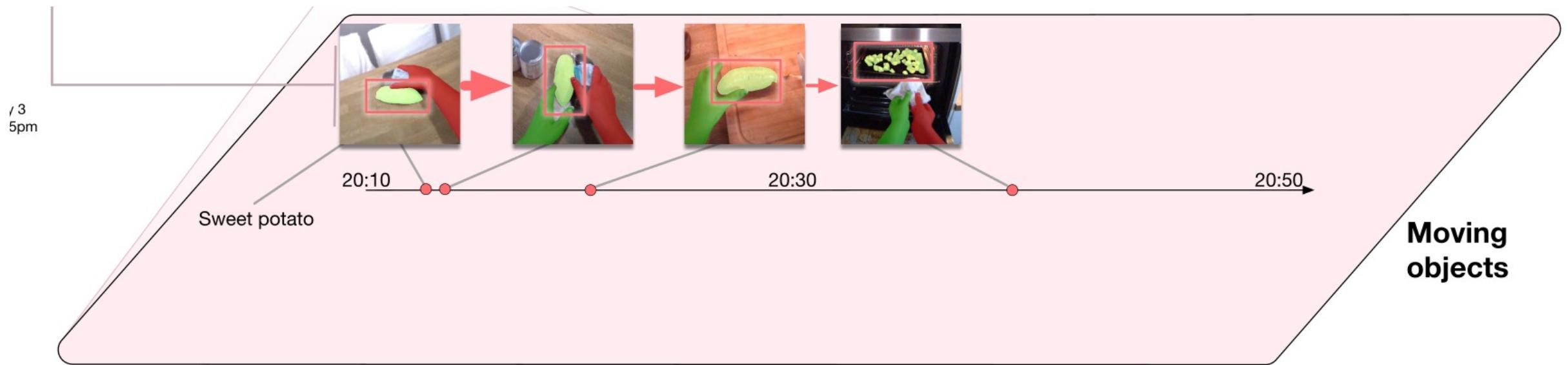
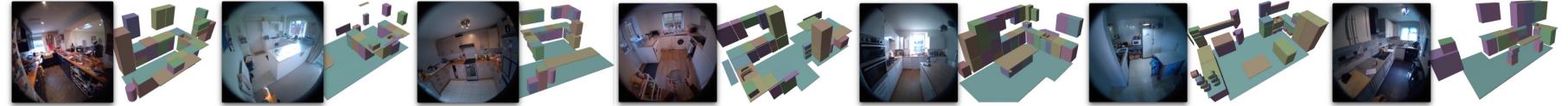


# HD-EPIC



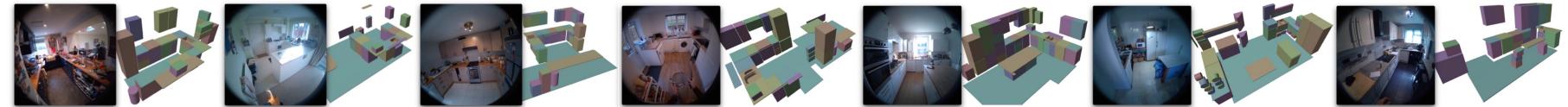


# HD-EPIC

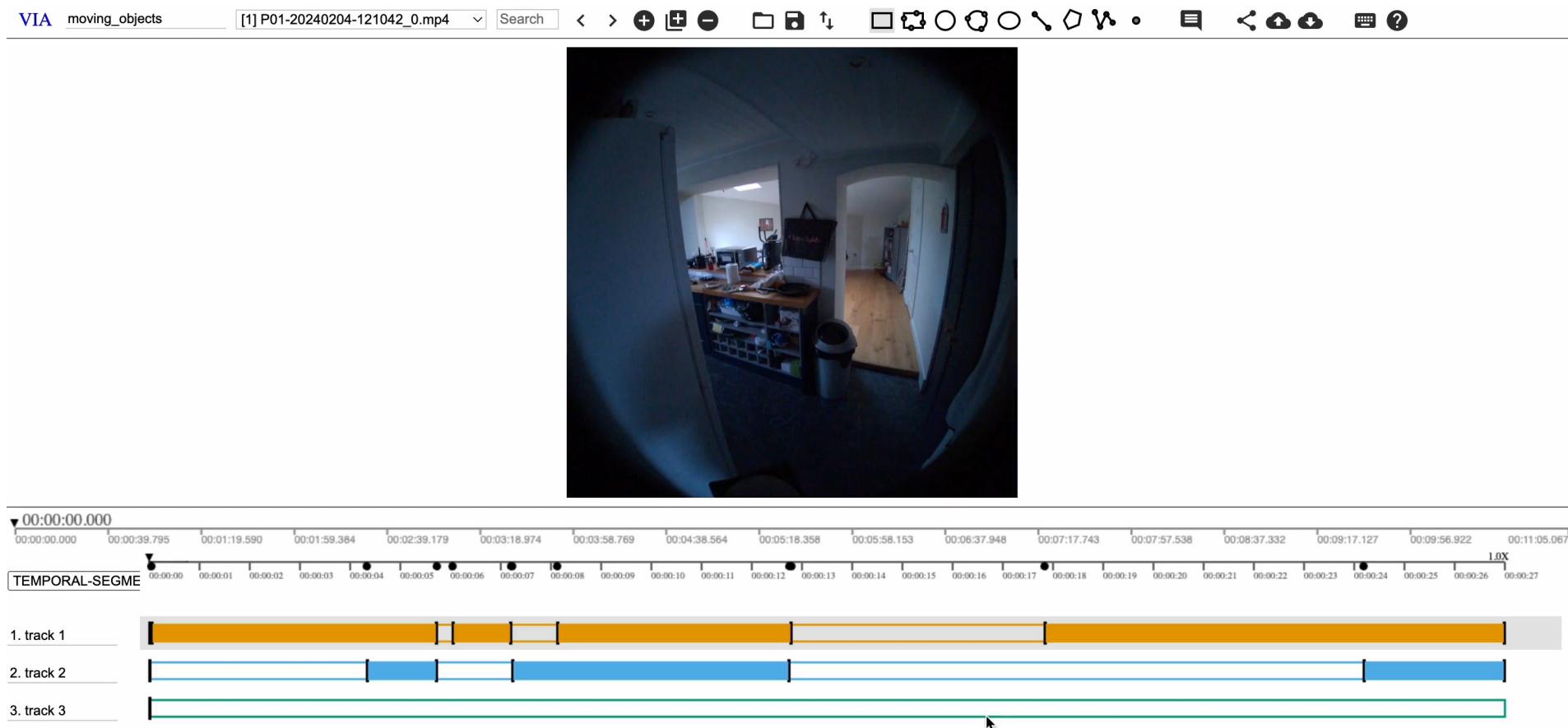




# HD-EPIC

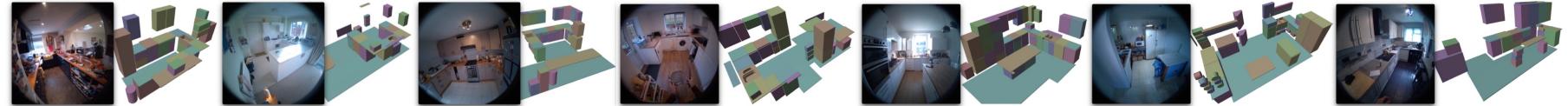


- How to minimize the annotations for tracking objects...

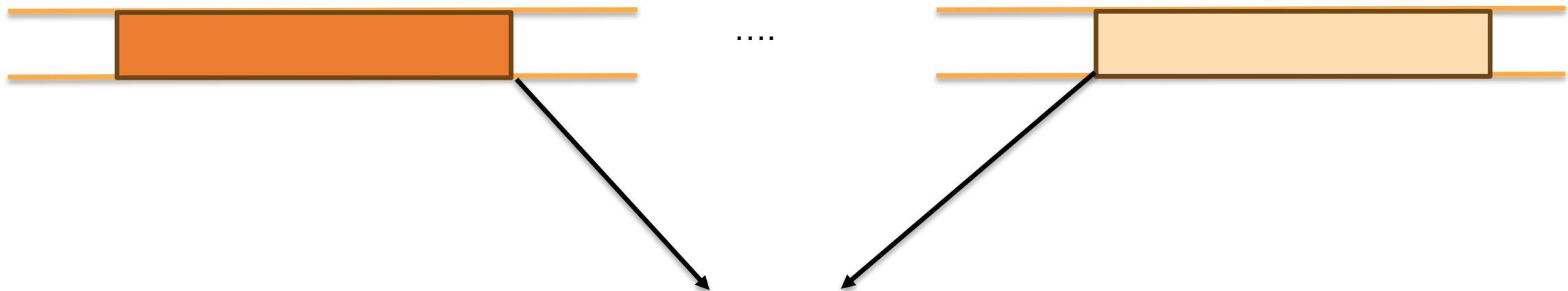




# HD-EPIC

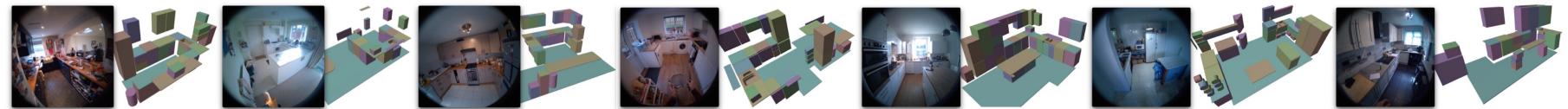


- How to minimize the annotations for tracking objects...





# HD-EPIC



## Current Track

### Image

| Choose Files | 201 files



42 / 199

← Previous

Next →

Undo

▼ [rubbish bin](#) [box of chicken](#) [wooden chopping board](#)

Enter Track Name (optional)

Create New Track

Inconsistent Query

## Previous Tracks

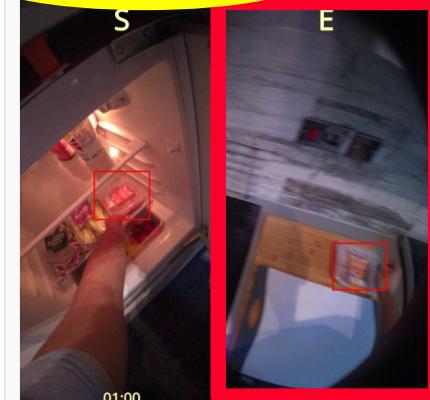
Sort by Distance

Save Tracks

box of chicken (0.0m)

plastic chopping board (0.3m)

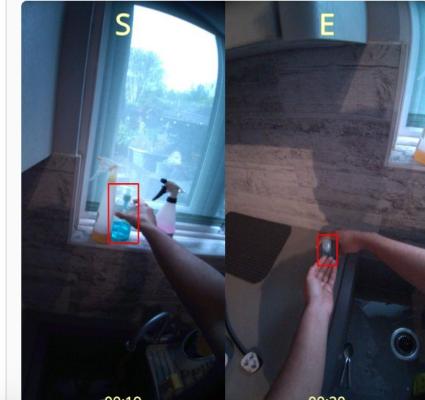
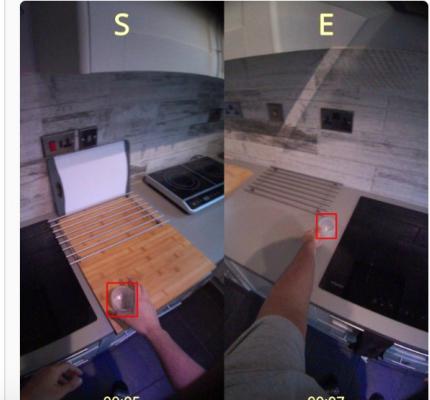
metal cooling rack (0.6m)



plastic measuring cup (1.0m)

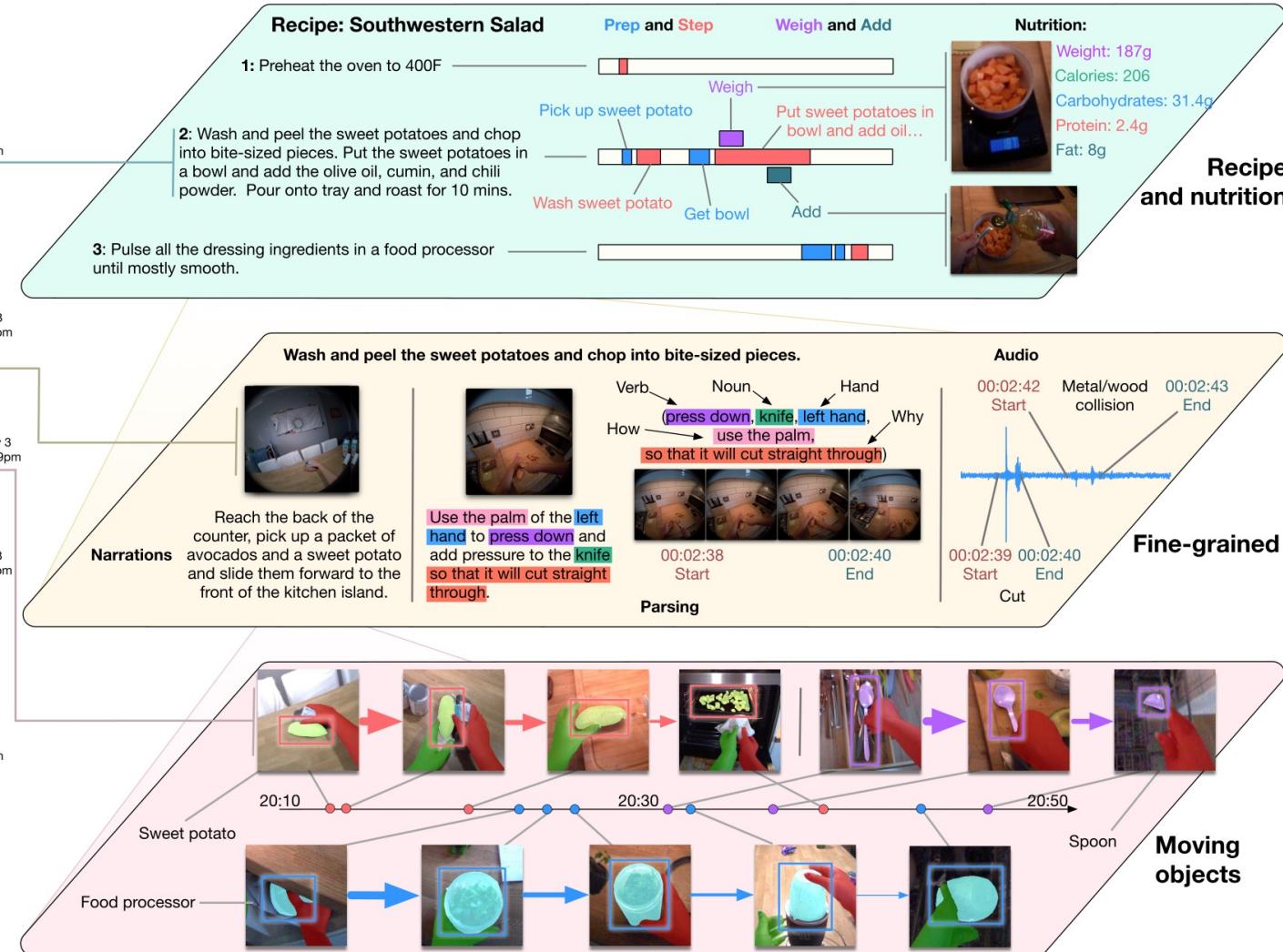
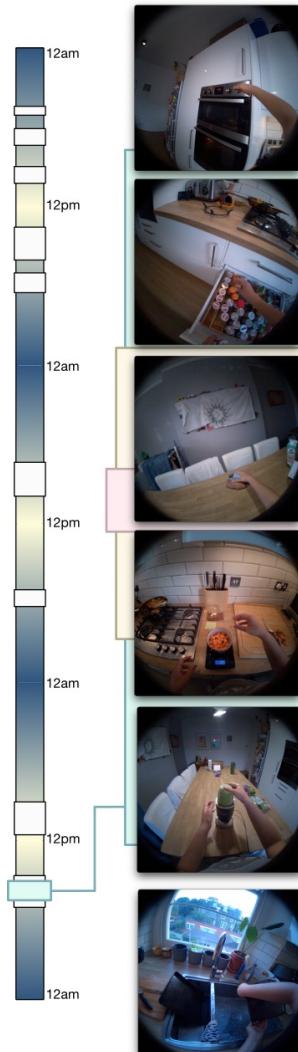
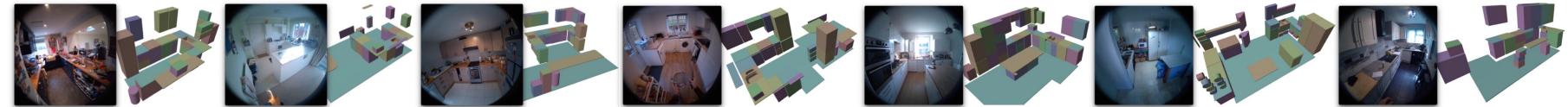
hand washing liquid (1.3m)

kitchen towel (1.5m)



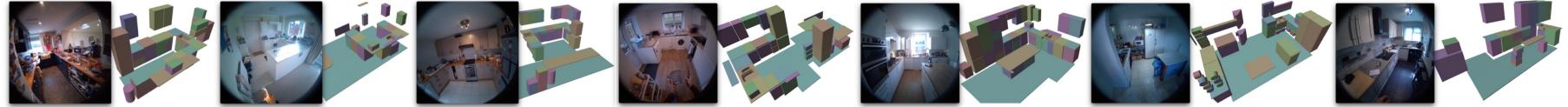


# HD-EPIC

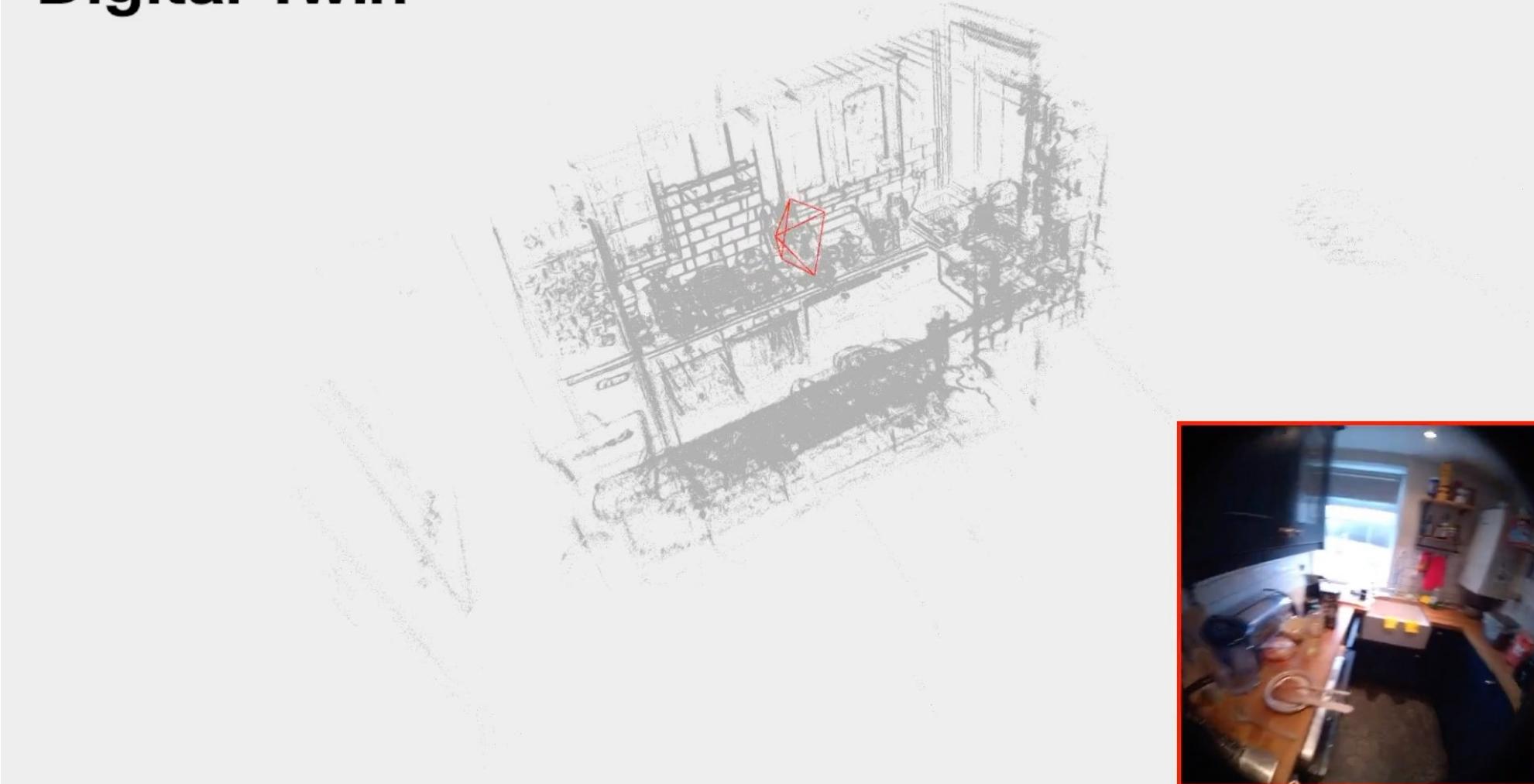




# HD-EPIC

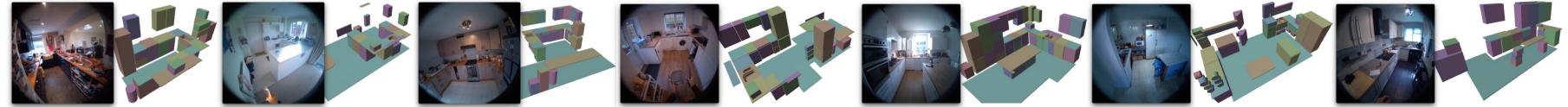


## Digital Twin





# HD-EPIC



fridge.001



counter.002



dishwasher.001



counter.003



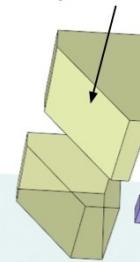
hob.001



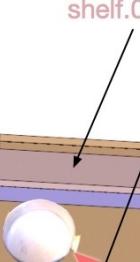
dishwasher.001



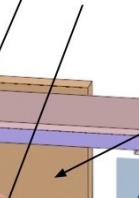
cupboard.004



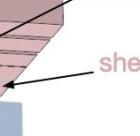
shelf.003



oven.001



hob.001



shelf.001



counter.003



counter.002



drawer.003



drawer.002



hook.001



basket.005



cupboard.001



fridge.001

Pointcloud

Surface mesh

Fixture Annotations



# HD-EPIC

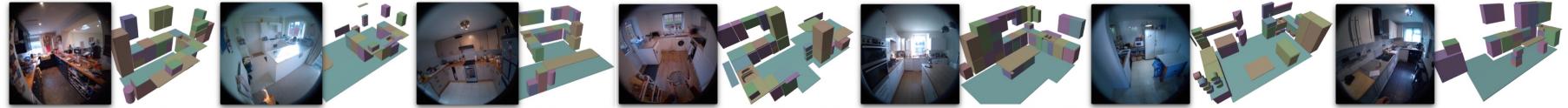
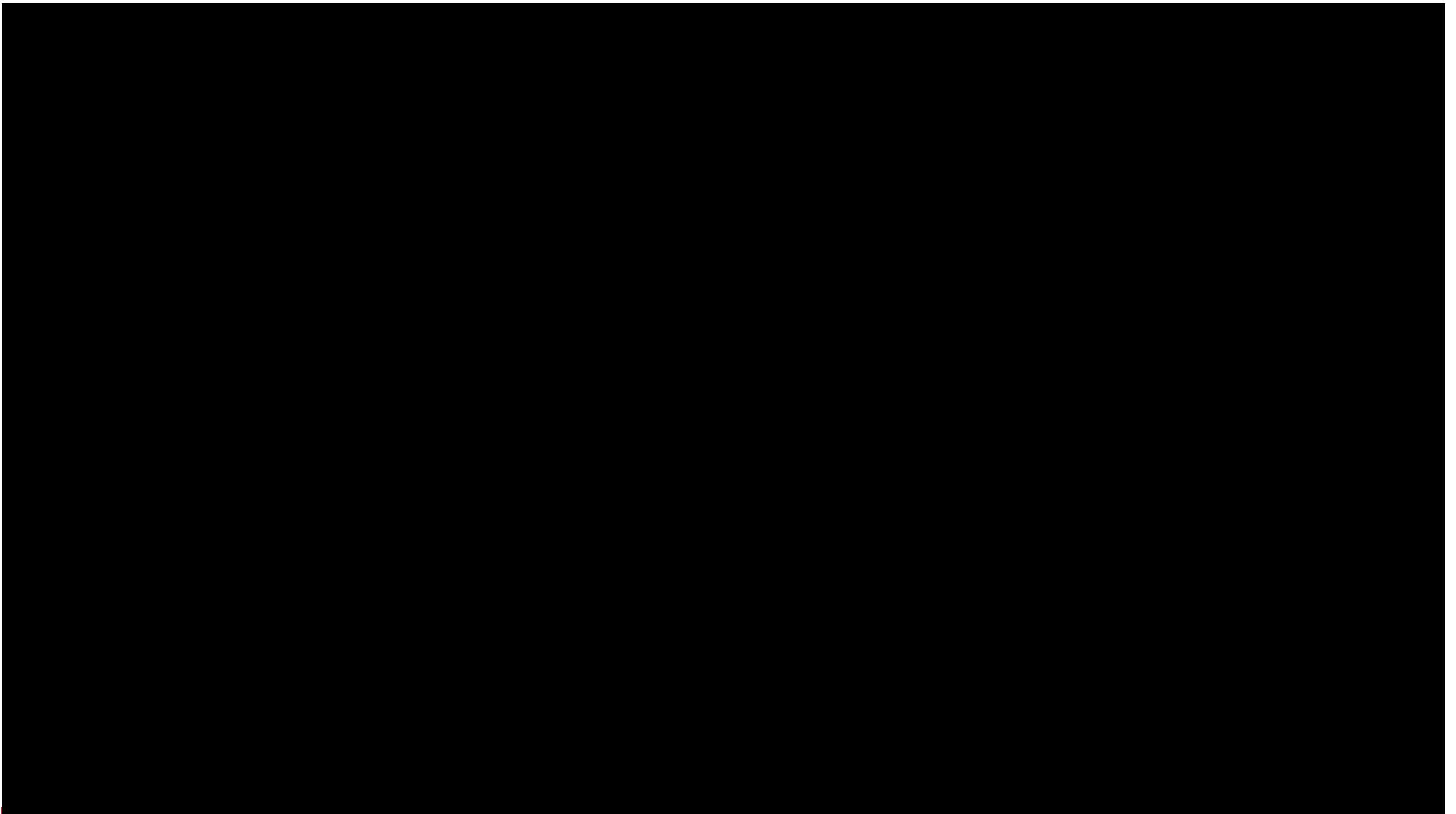
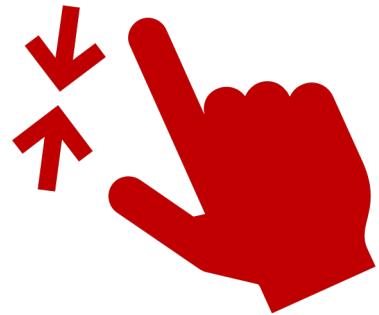
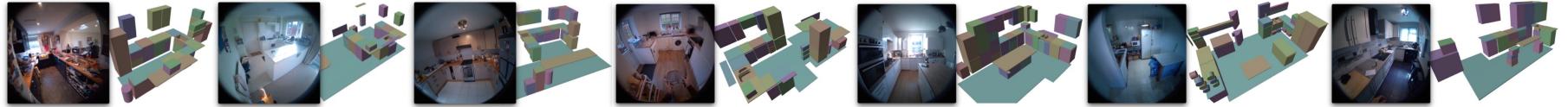


Image Selection

Select Folder: Choose Files 0000000084.jpg

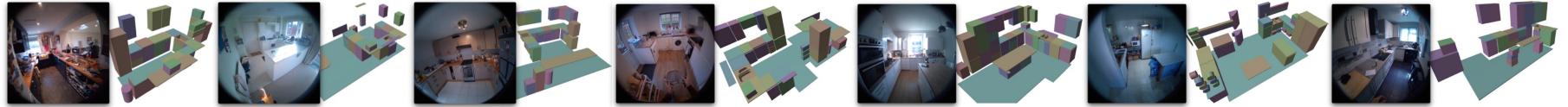




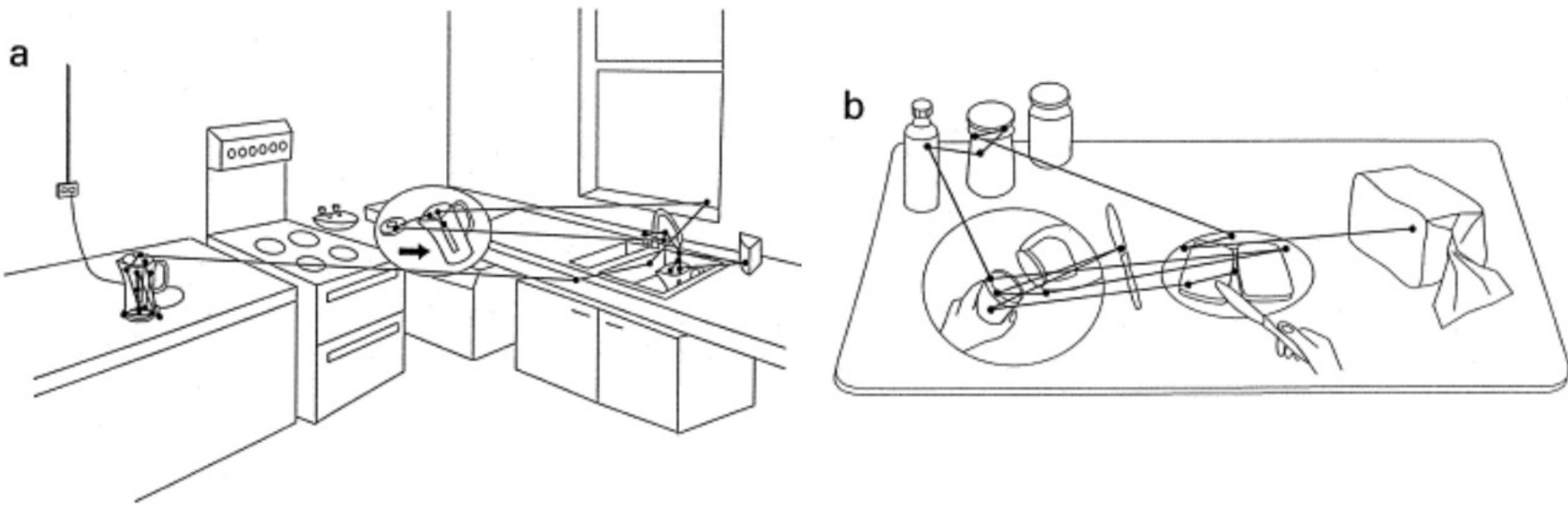
# Video Understanding Out of the Frame

What can we now do with these reconstructions:

- Point Tracking
- Object Tracking
- Gaze Estimation

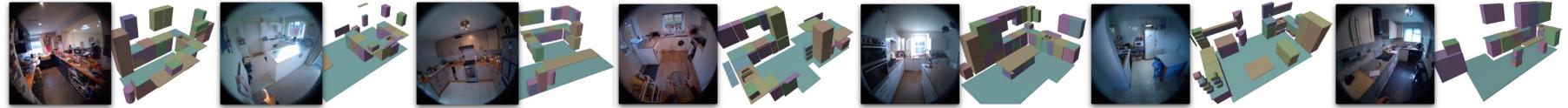


# Gaze and Fixations

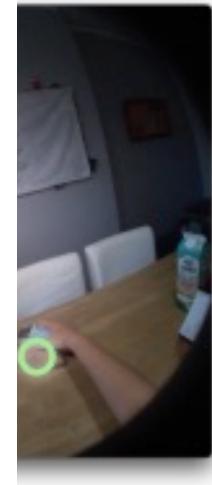
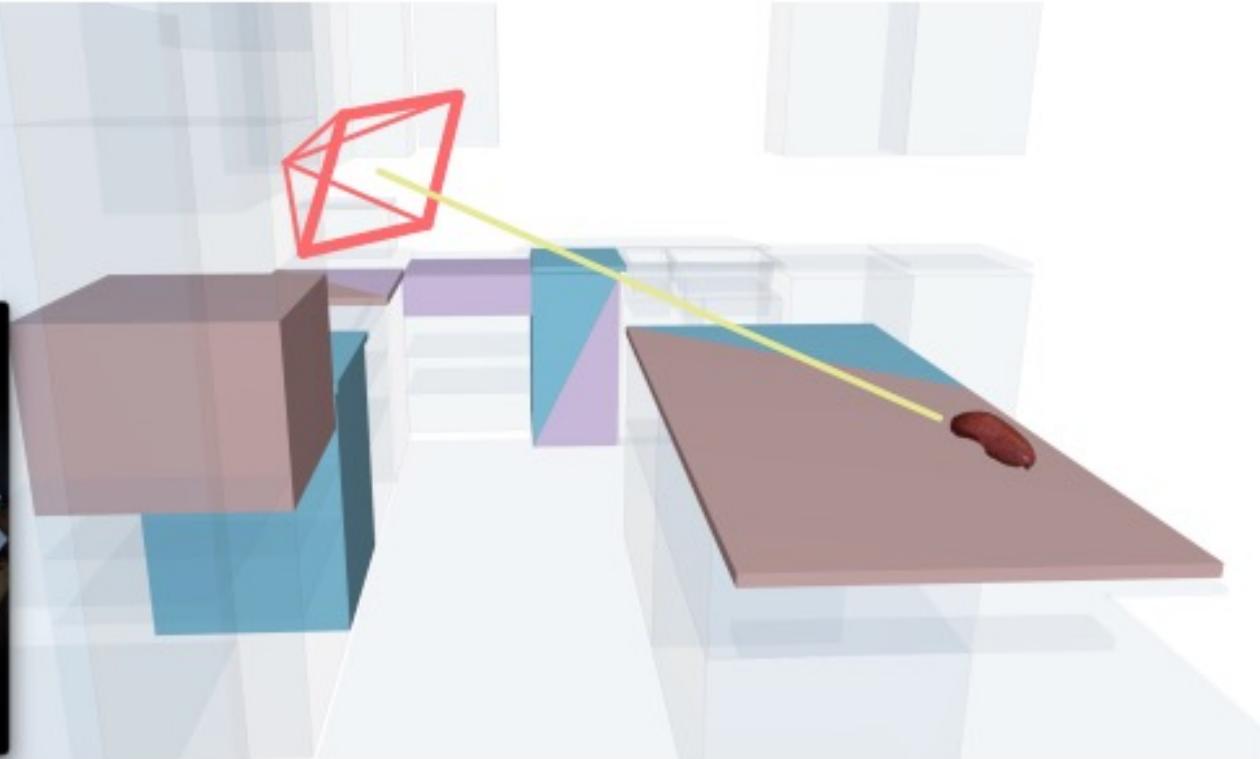
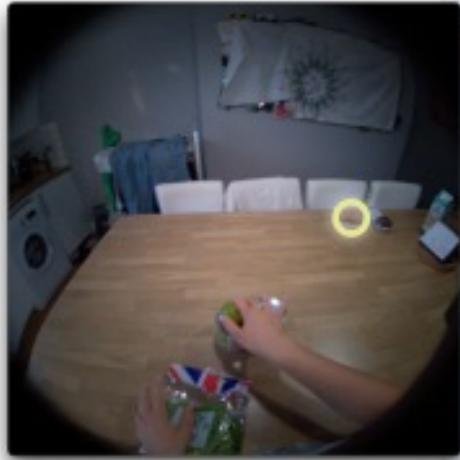




# HD-EPIC

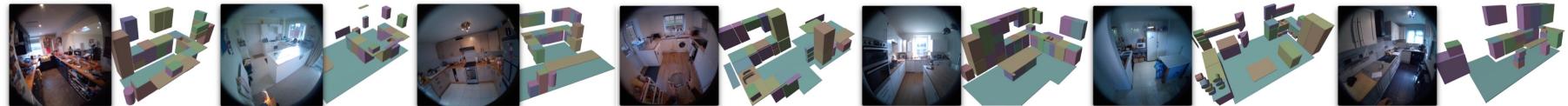


## Gaze priming





# HD-EPIC



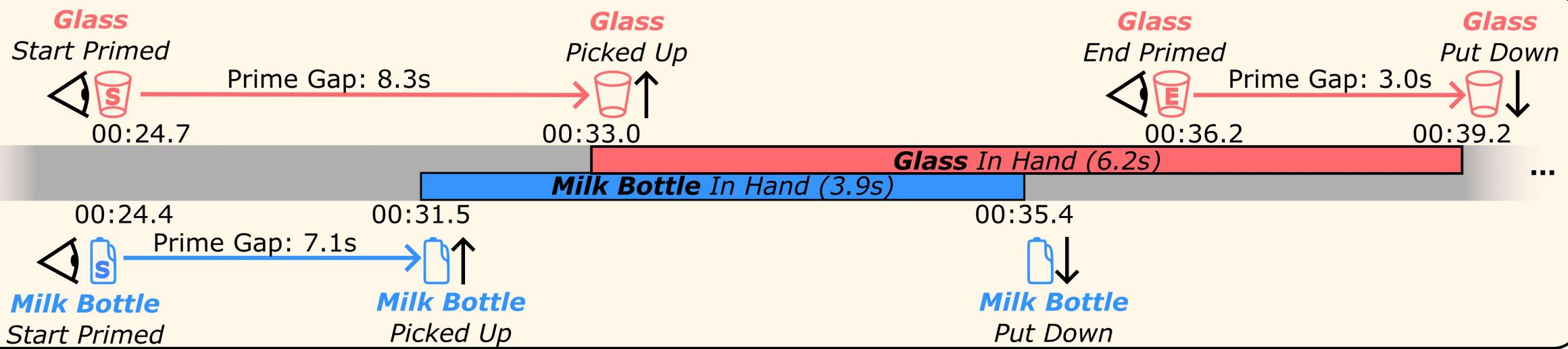
3D Scene



Frames w/ 2D Gaze

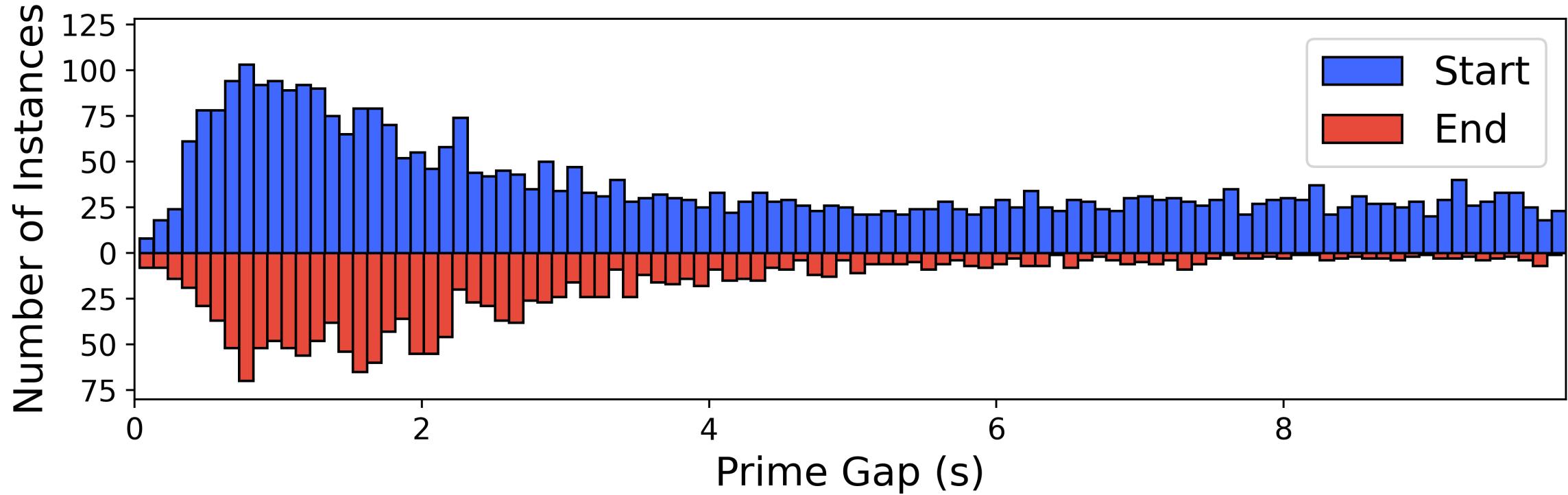
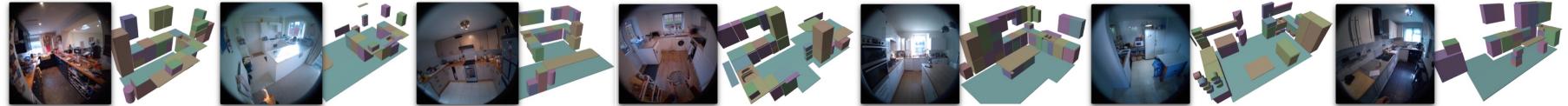


Object Movement





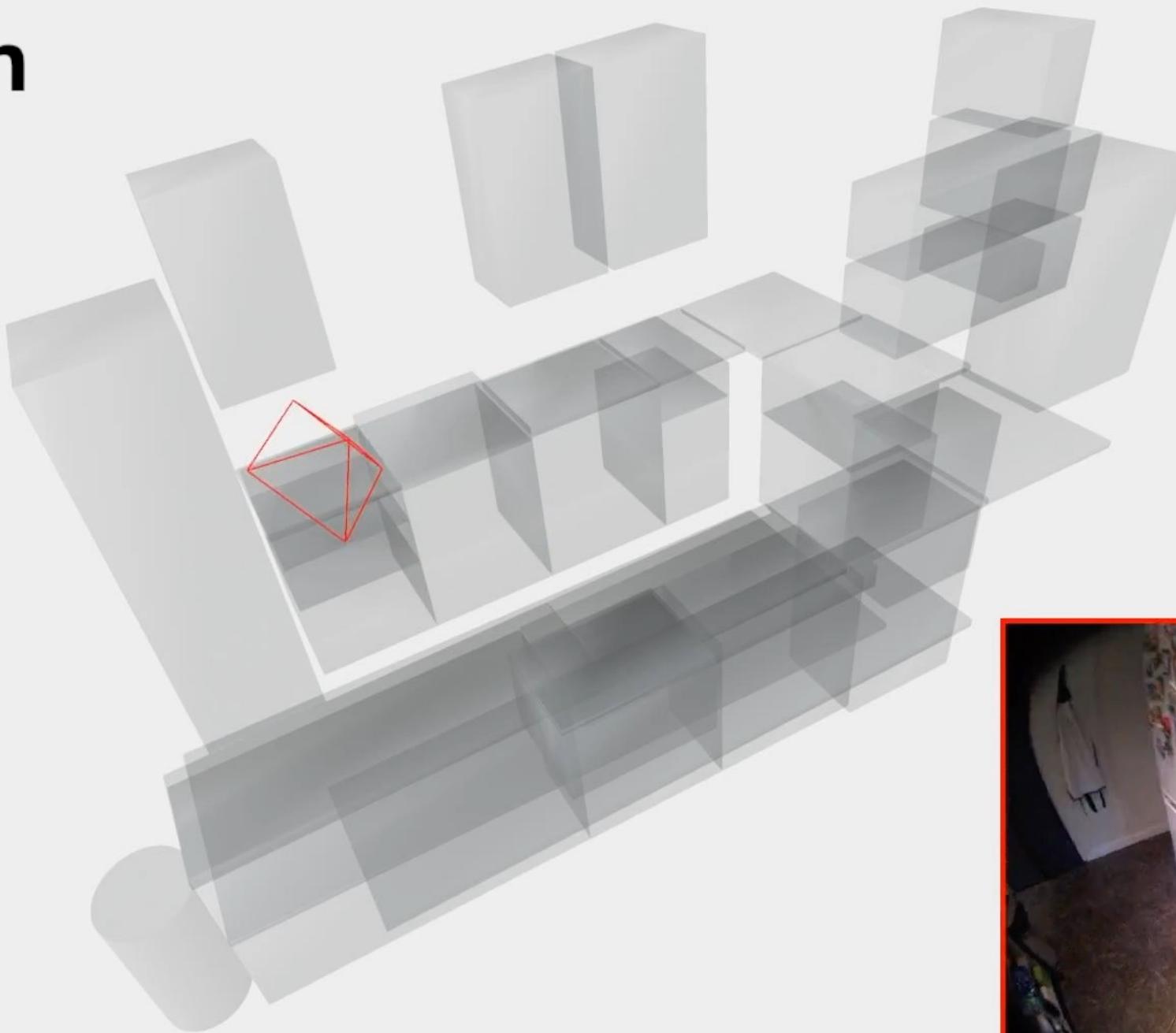
# HD-EPIC



# Digital Twin

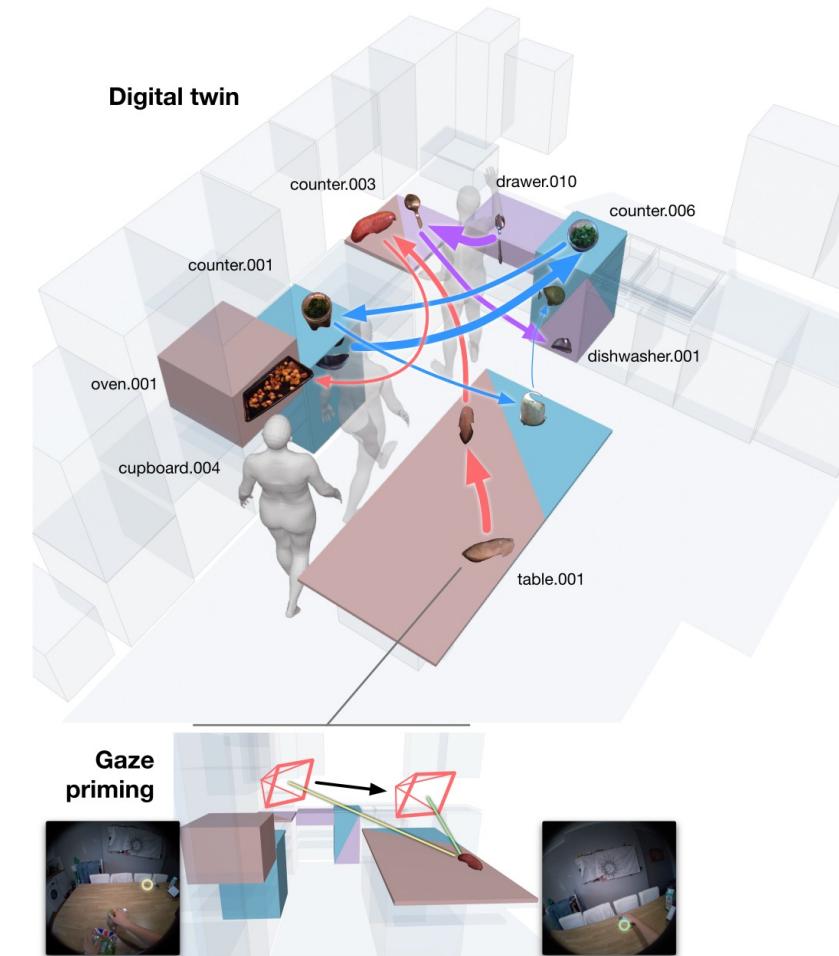
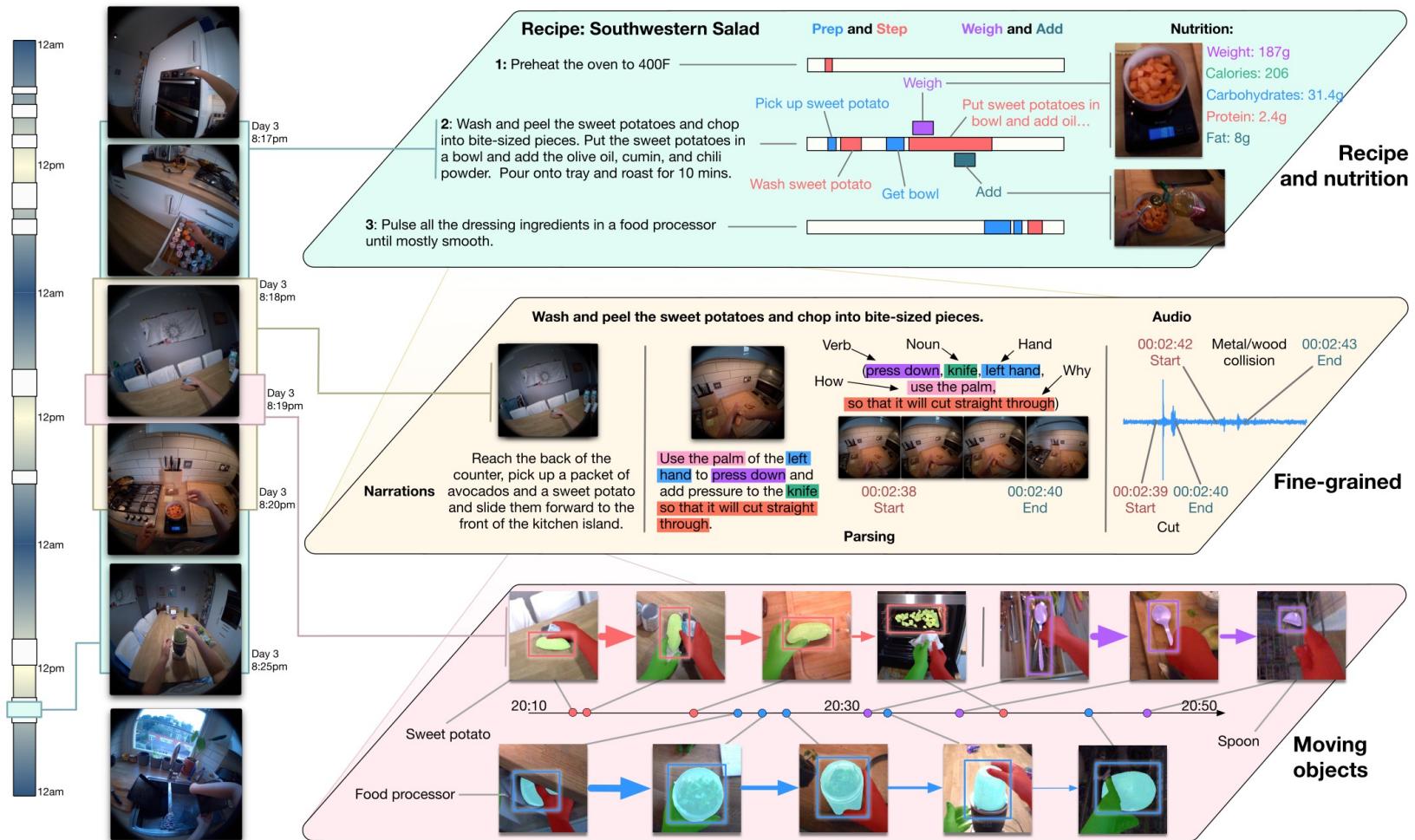
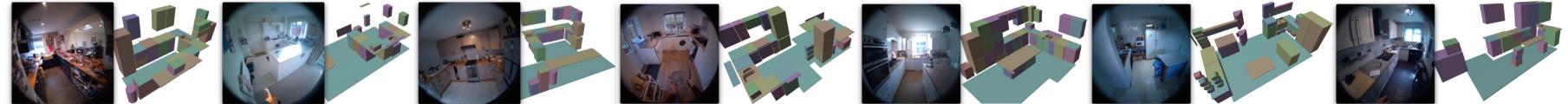
## Fixtures

Open drawer



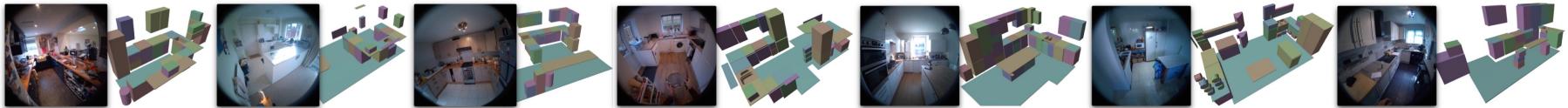


# HD-EPIC





# HD-EPIC

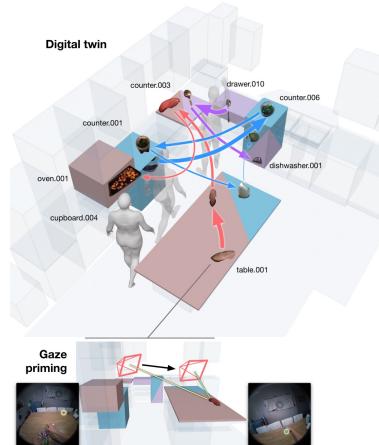
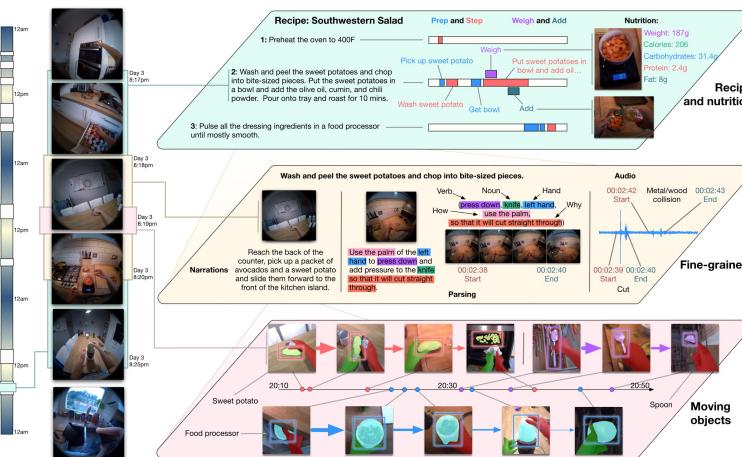
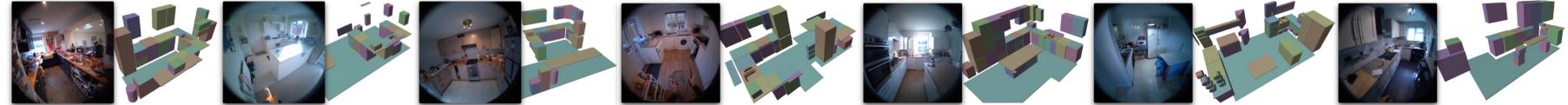


Annotation Type	Total annotations	Annotations/min
Narrations	59,454	24.0
Parsing (Verbs + Nouns + Hands + How + Why)	303,968	122.7
Recipes (Preps + Steps)	4,052	1.6
Sound	50,968	20.6
Action boundaries	59,454	24.0
Object Motion (Pick up + Put down + Fixtures + Bboxes + Masks)	153,480	62.0
Object Itinerary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

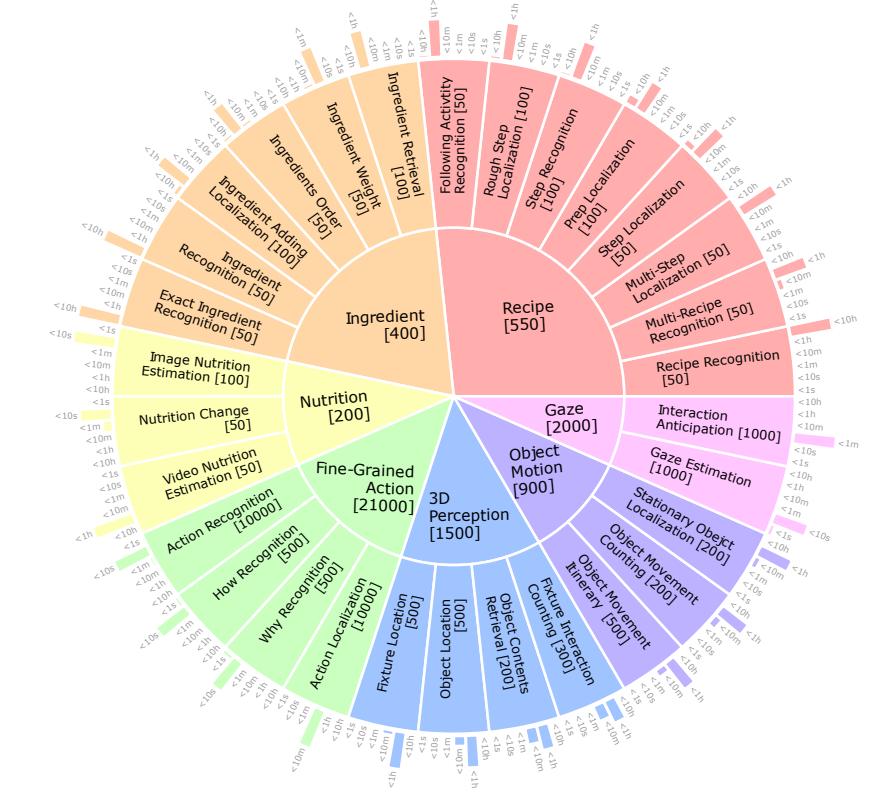
Table A3. HD-EPIC annotations per minute



# HD-EPIC



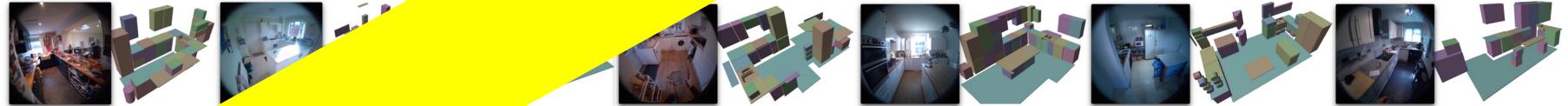
## Sec 1: Highly-Detailed Dataset



## Sec 2: HD-EPIC VQA Benchmark



# HD-EPIC



Try it Yourself

Use Wise to Search  
through HD-EPIC



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



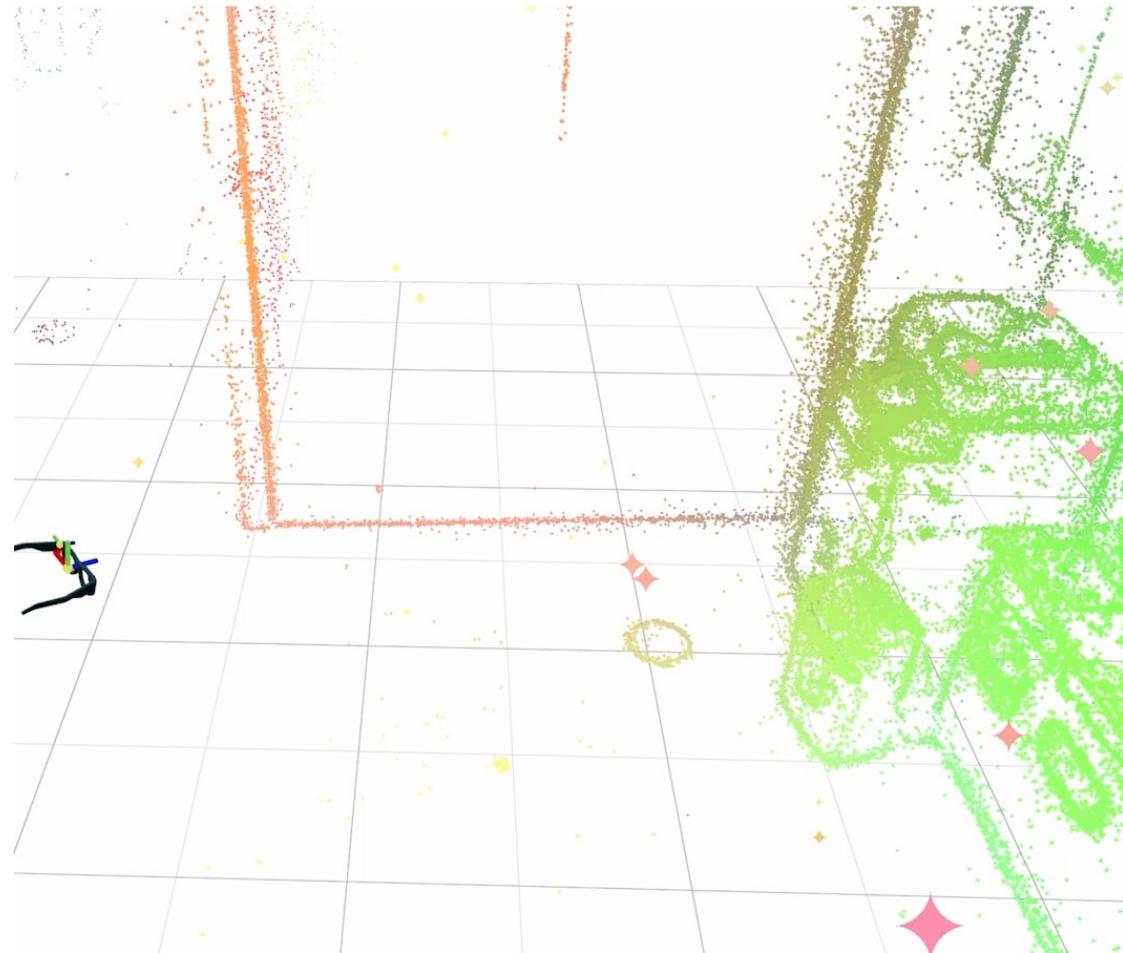
# Video Understanding Out of the Frame

Body and Hand Motion Estimation “out of the frame”



# EgoBody

EgoAllo uses egocentric (6d) SLAM poses and images

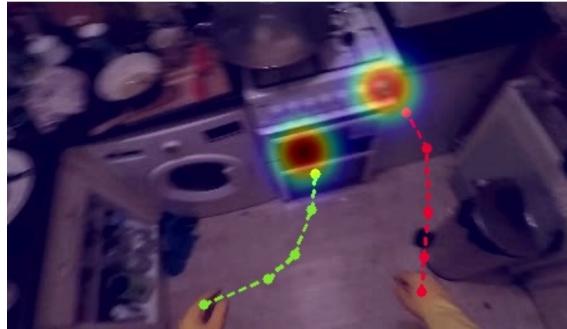


# EgoHand Forecasting – Previous Works

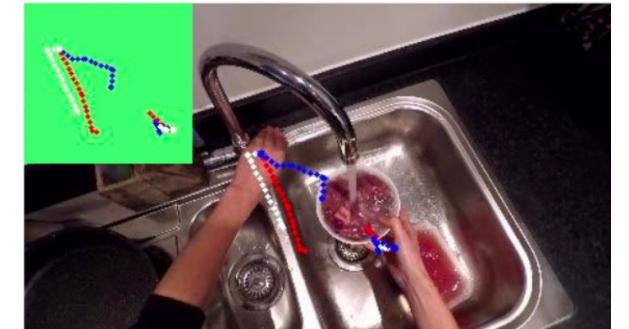
with: Masashi Hatano  
Zhifan Zhu  
Hideo Saito

## 2D Hand Forecasting

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



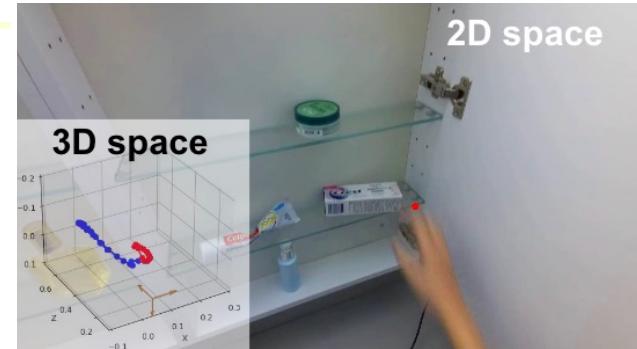
OCT [CVPR'22]



Diff-IP2D [IROS'25]

## 3D Hand Forecasting

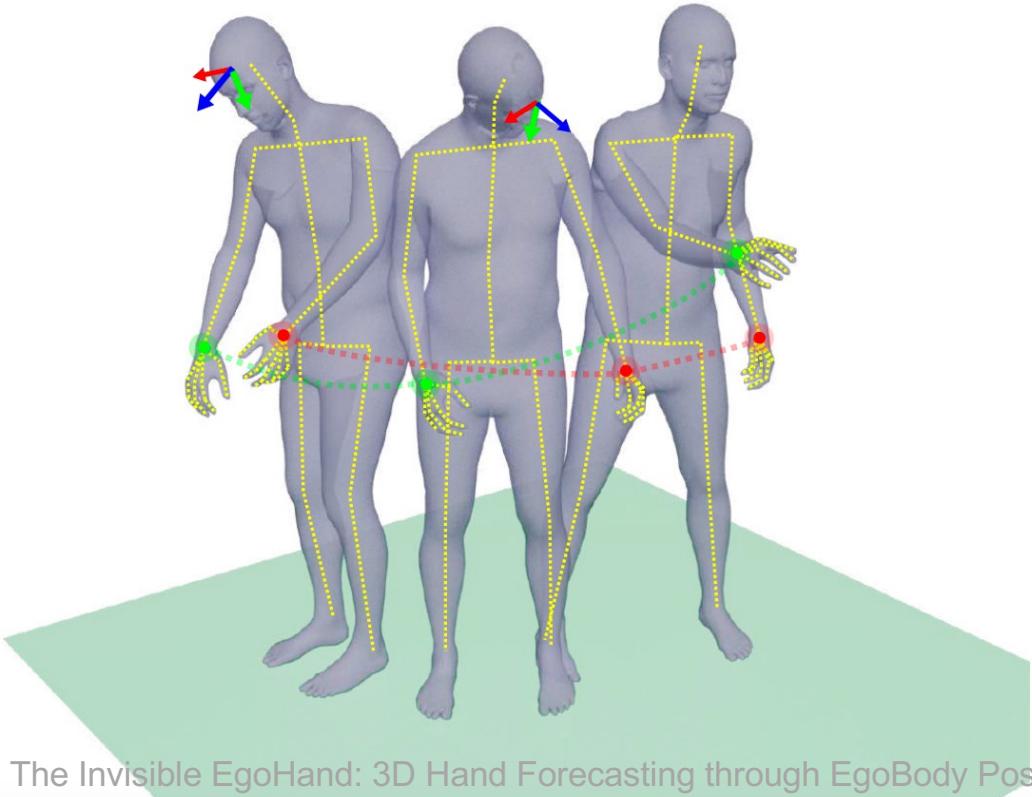
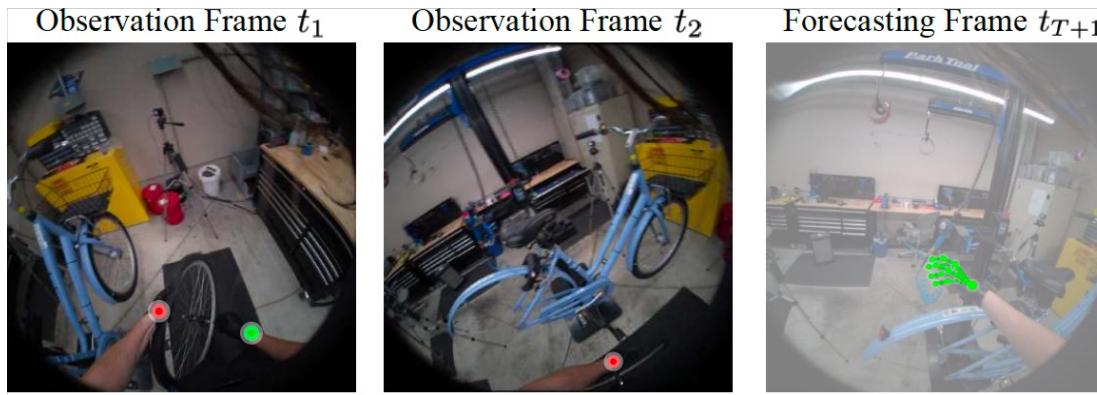
Given an egocentric video & 3D hand trajectory,  
forecast 3D hand positions of one hand



USST [ICCV'23]

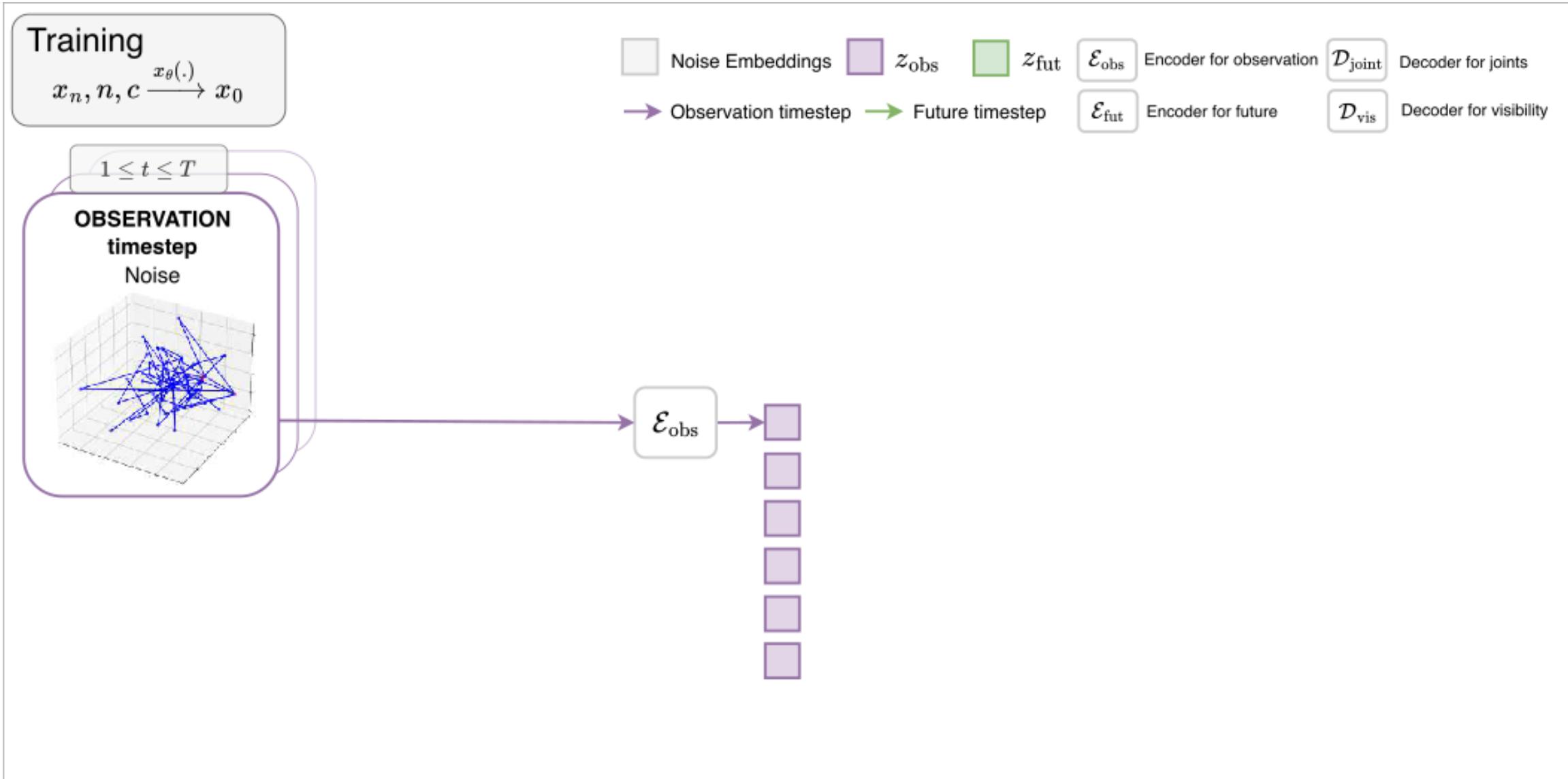
# The Invisible EgoHand

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

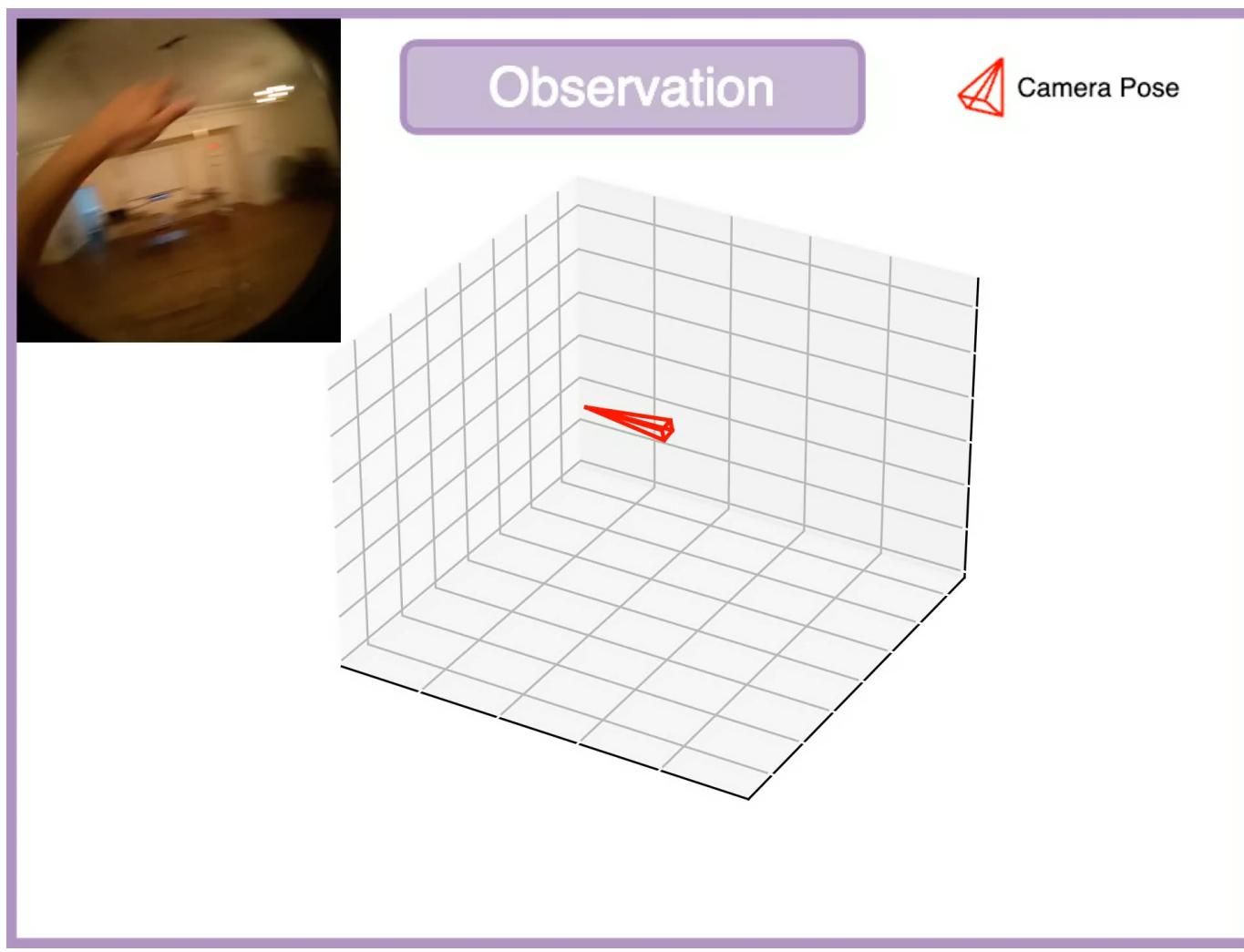
with: Masashi Hatano  
Zhifan Zhu  
Hideo Saito

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
Ours	<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



with: Masashi Hatano  
Zhifan Zhu  
Hideo Saito





# Video Understanding Out of the Frame

From First-Point View to Second- and Third-





# Ego-Exo4D

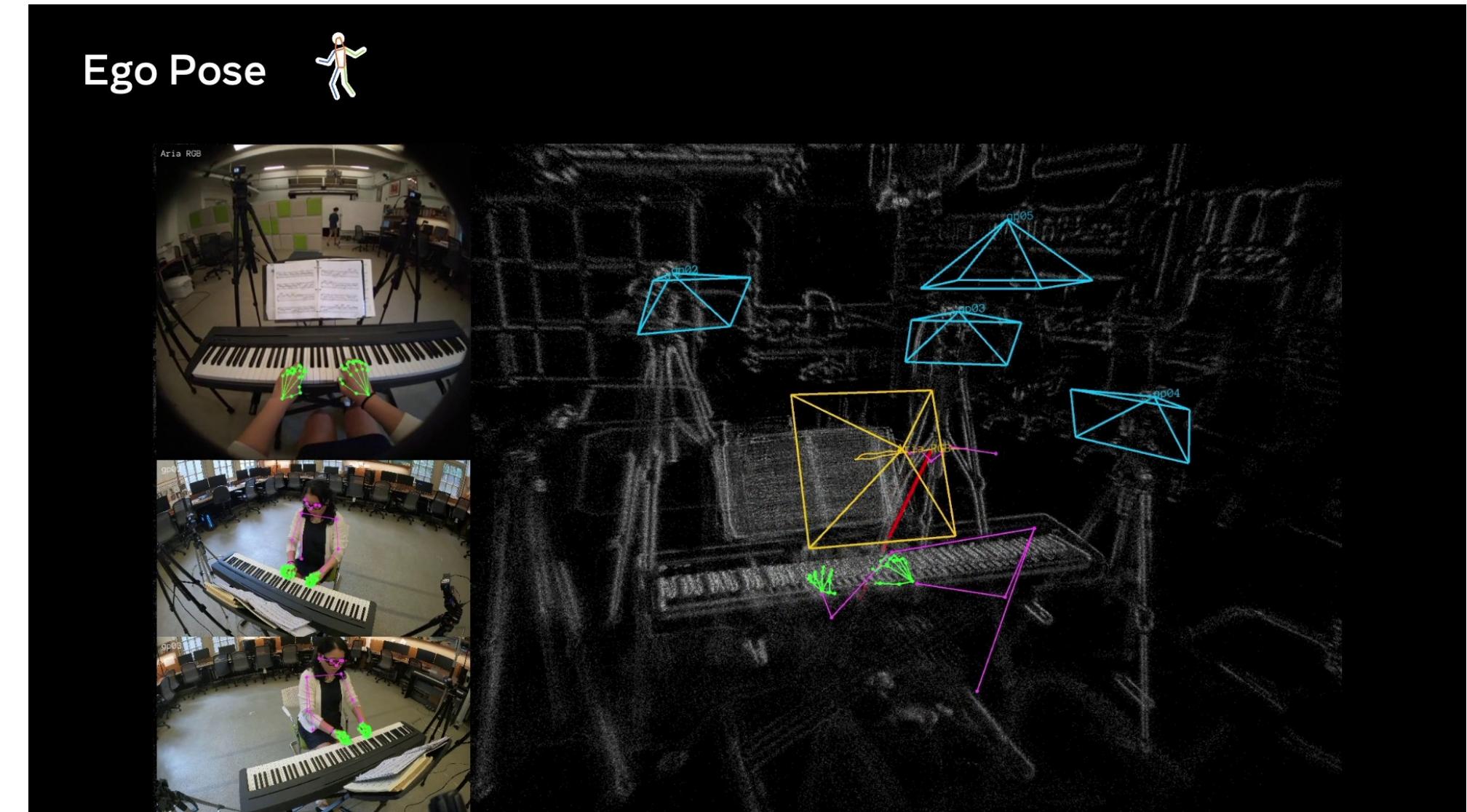
with: Kristen Grauman  
+102 authors

## Ego-Exo Relation

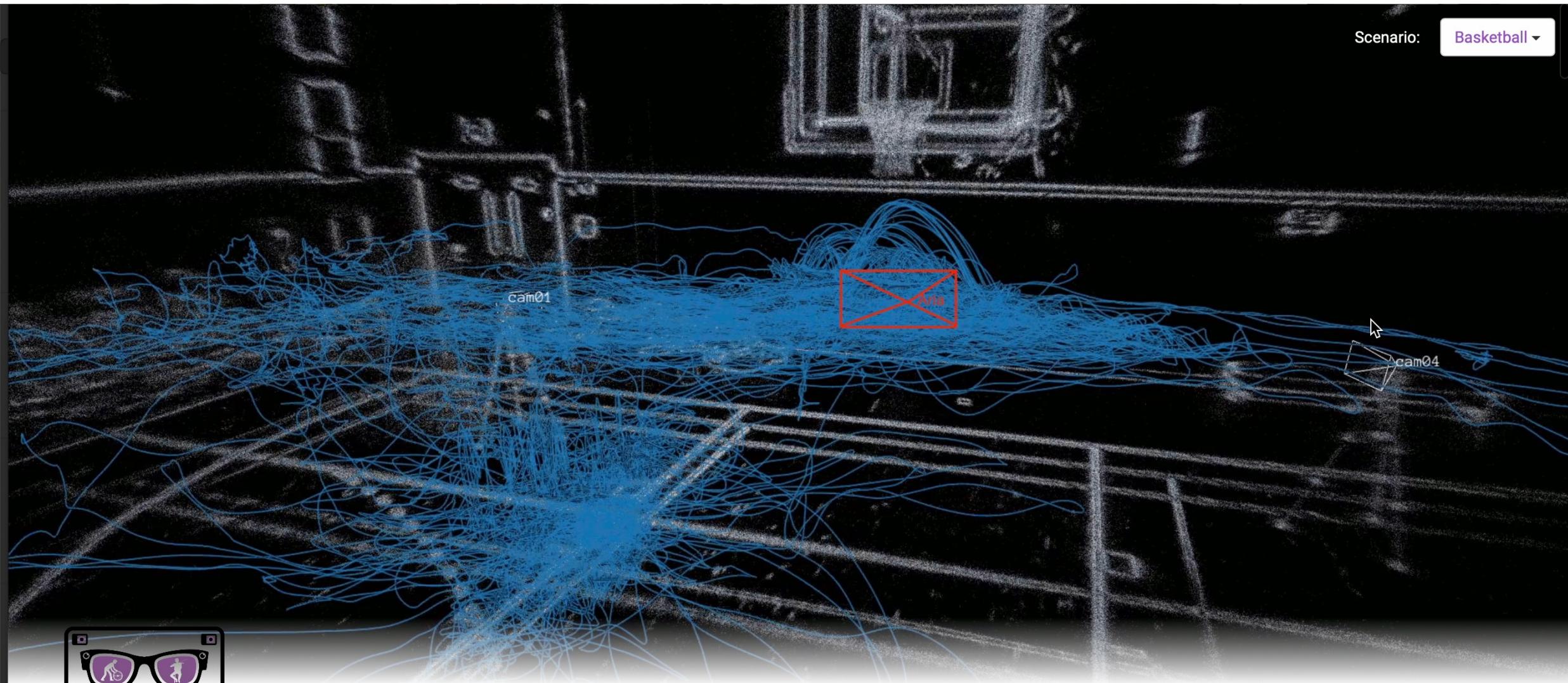




with: Kristen Grauman  
+102 authors



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 ↗](#)

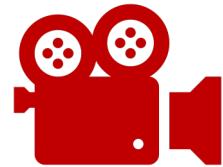
[Start Here ↘](#)



Motivation and Datasets in  
Egocentric Video Understanding



Video Understanding  
Out of the Frame



Video Understanding:  
Data and Tasks



Teaser: The Wizard of Oz  
& Genie 3



Videos are Multimodal



Outlook into the Future of  
Egocentric Vision



Connected Videos of One's Life



Conclusion

# The Wizard of Oz at the Sphere

Sold out tickets – August 2025

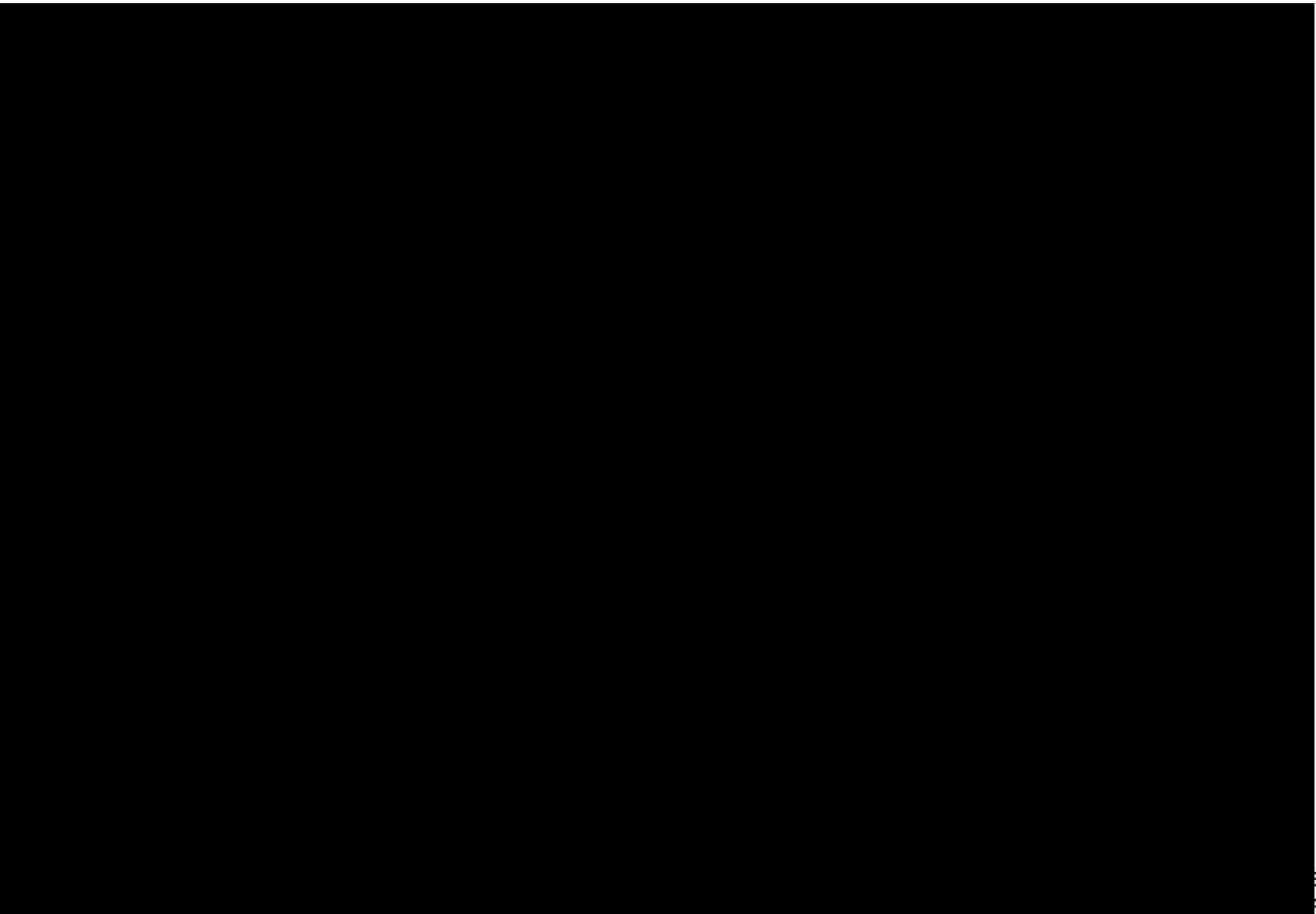


<https://behindthecurtain.withgoogle.com>

Dima Damen  
PAISS 2025

# The Wizard of Oz @ The Sphere

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



# The Wizard of Oz @ The Sphere

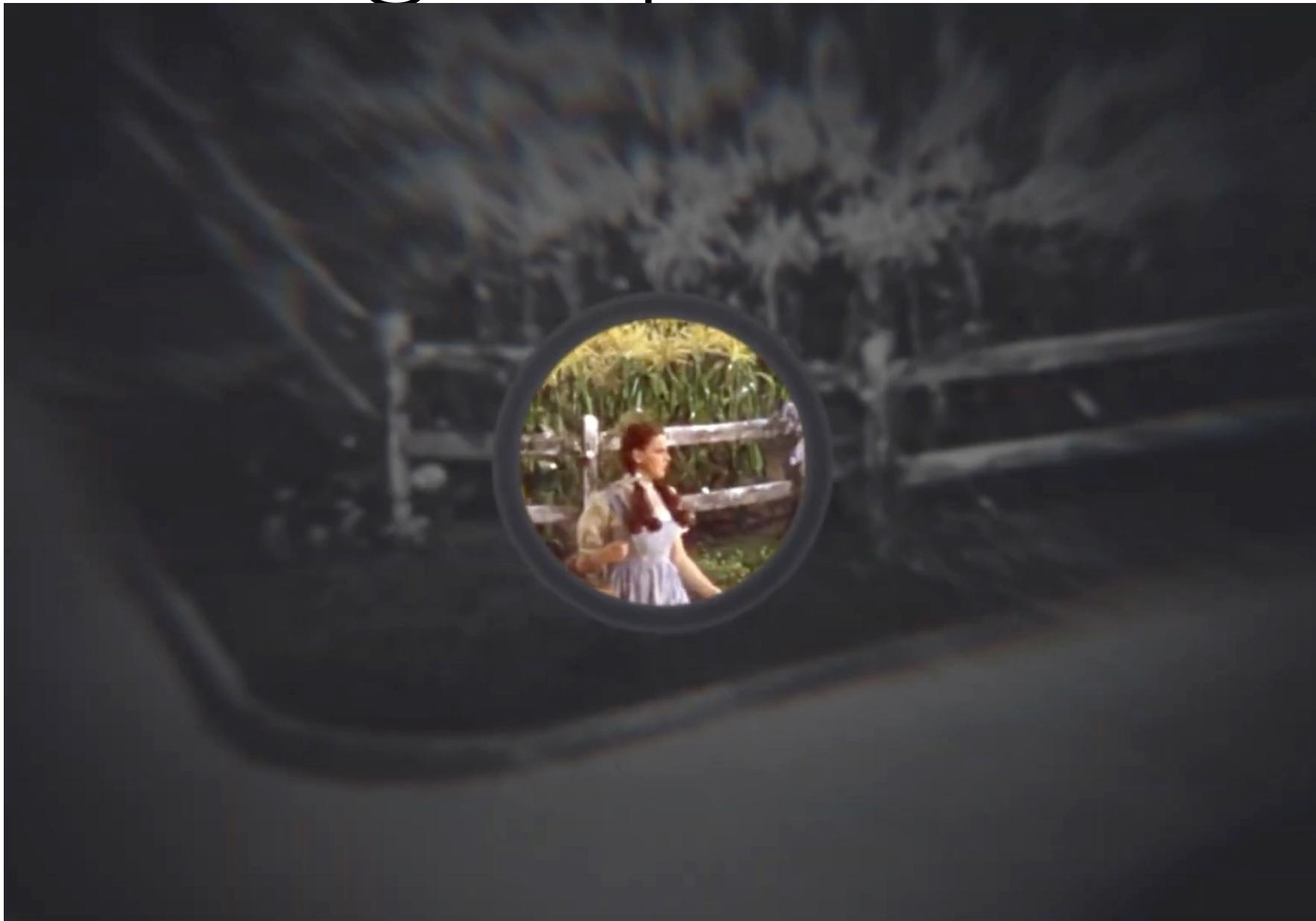


**Ralph Winte**

Head of Physical Production -  sphere

Dima Damen  
PAISS 2025

# The Wizard of Oz @ The Sphere



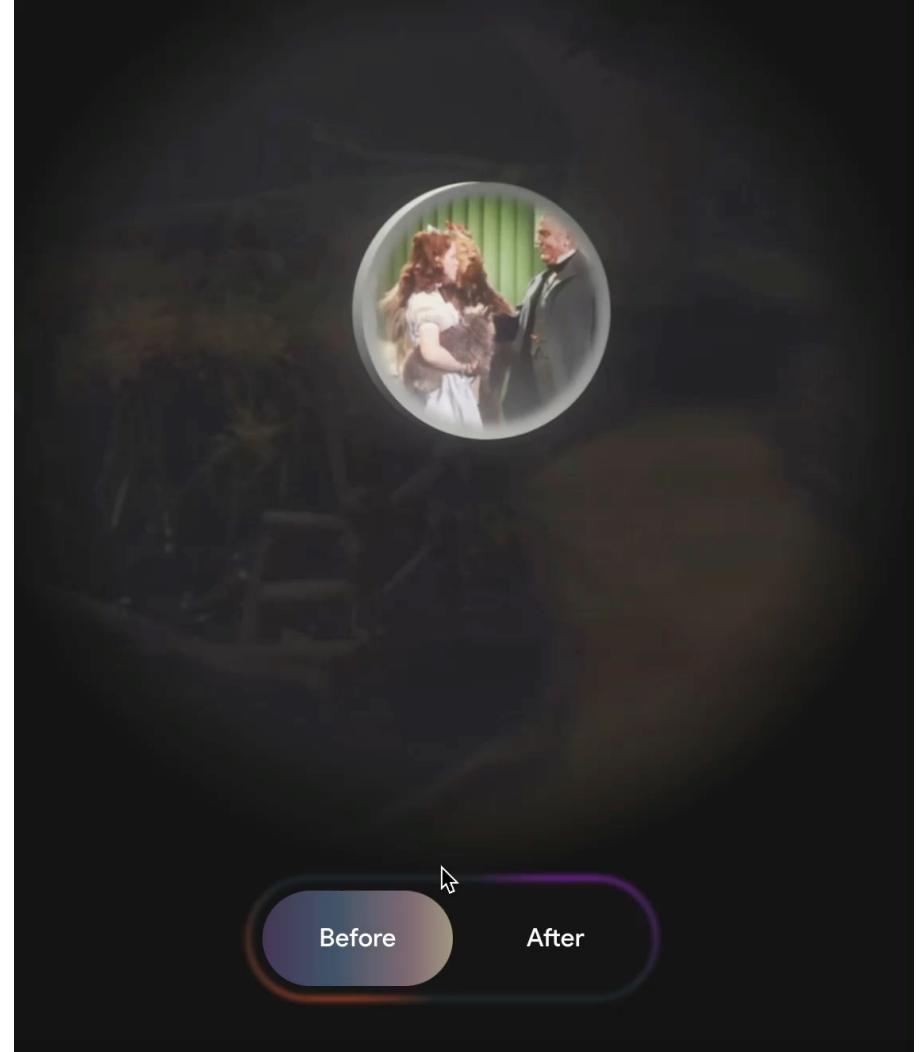
Dima Damen  
PAISS 2025

# The Wizard of Oz @ The Sphere

- Super-resolution,
- Outpainting...



<https://behindthecurtain.withgoogle.com>



Dima Damen  
PAISS 2025

# The Wizard of Oz @ The Sphere

- Performance Interpolation,



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



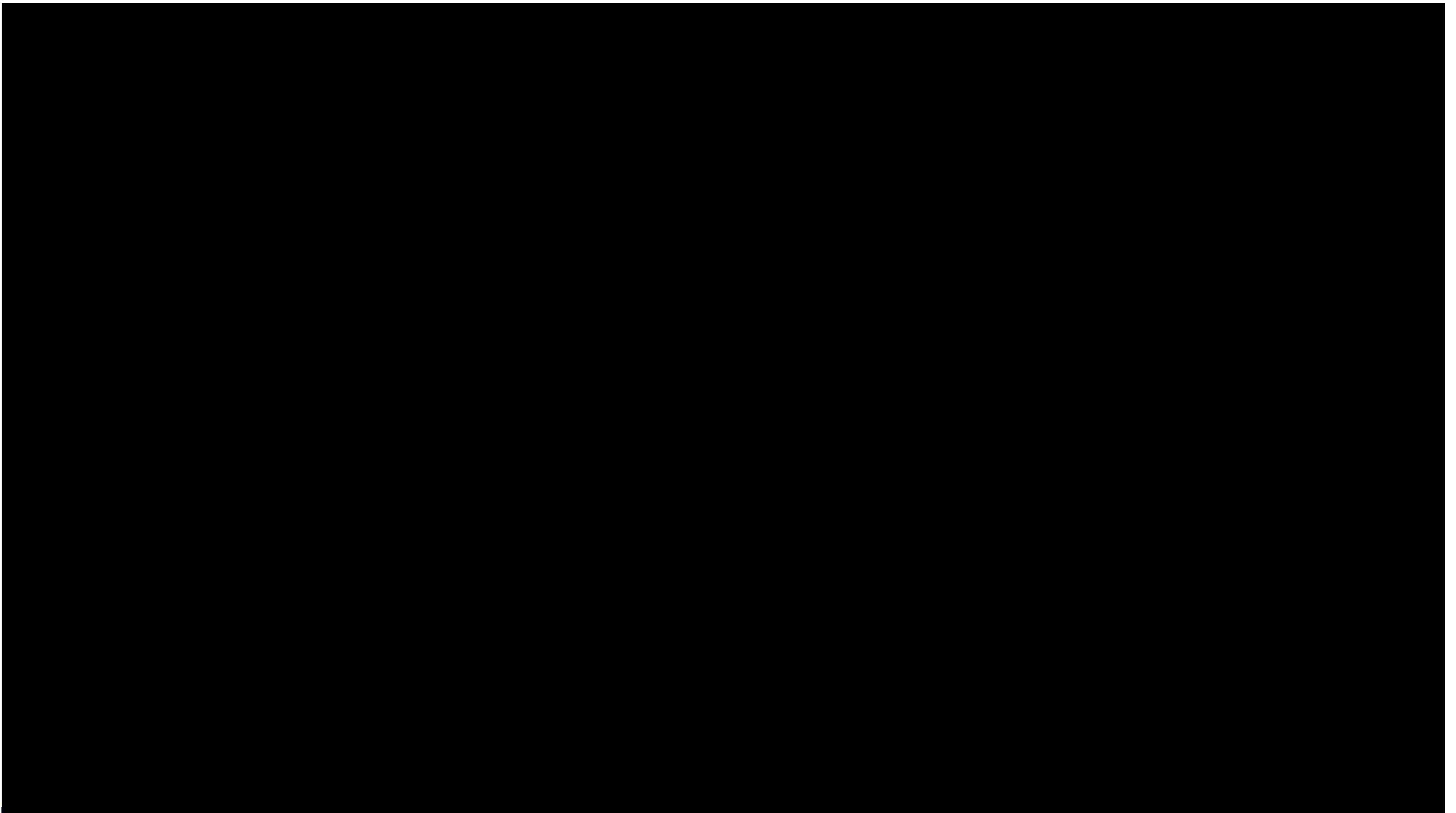
# The Wizard of Oz @ The Sphere

- Auto-Director



<https://behindthecurtain.withgoogle.com>

Dima Damen  
PAISS 2025





# 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  
PAISS 2025

# Genie 3

Released Aug 2025



<https://deepmind.google/discover/blog/genie-3-a-new-frontier-for-world-models/>

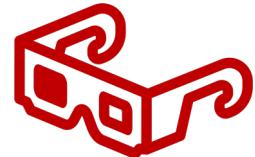
Dima Damen  
PAISS 2025



# Genie 3

- A new frontier of world models
- From a text prompt or a starting image/video
- Generate interactive worlds at 24fps and 720px
- Remains consistent for several minutes
- Genie 3's consistency is an emergent capability





Motivation and Datasets in  
Egocentric Video Understanding



Video Understanding  
Out of the Frame



Video Understanding:  
Data and Tasks



Teaser: The Wizard of Oz  
& Genie 3



Videos are Multimodal



Outlook into the Future of  
Egocentric Vision



Connected Videos of One's Life



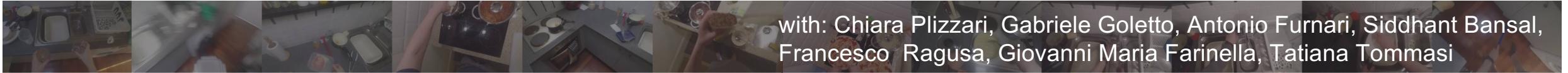
Conclusion



# 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>

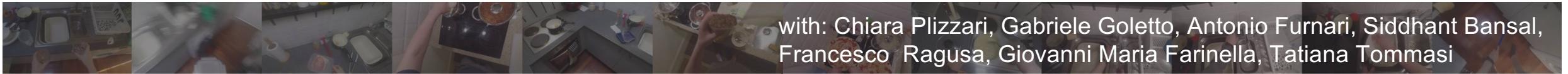




with: Chiara Plizzari, Gabriele Goletto, Antonio Furnari, Siddhant Bansal,  
Francesco Ragusa, Giovanni Maria Farinella, Tatiana Tommasi

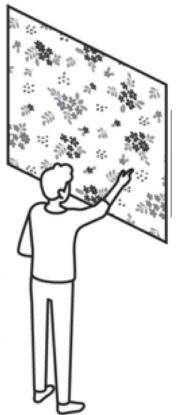
# Envisioning an Ambitious Future and Analysing the Current Status of Egocentric Vision

How did we do this?



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

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



**EGO-Designer**



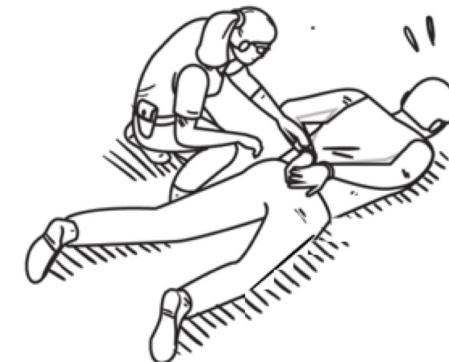
**EGO-Worker**



**EGO-Tourist**



**EGO-Home**

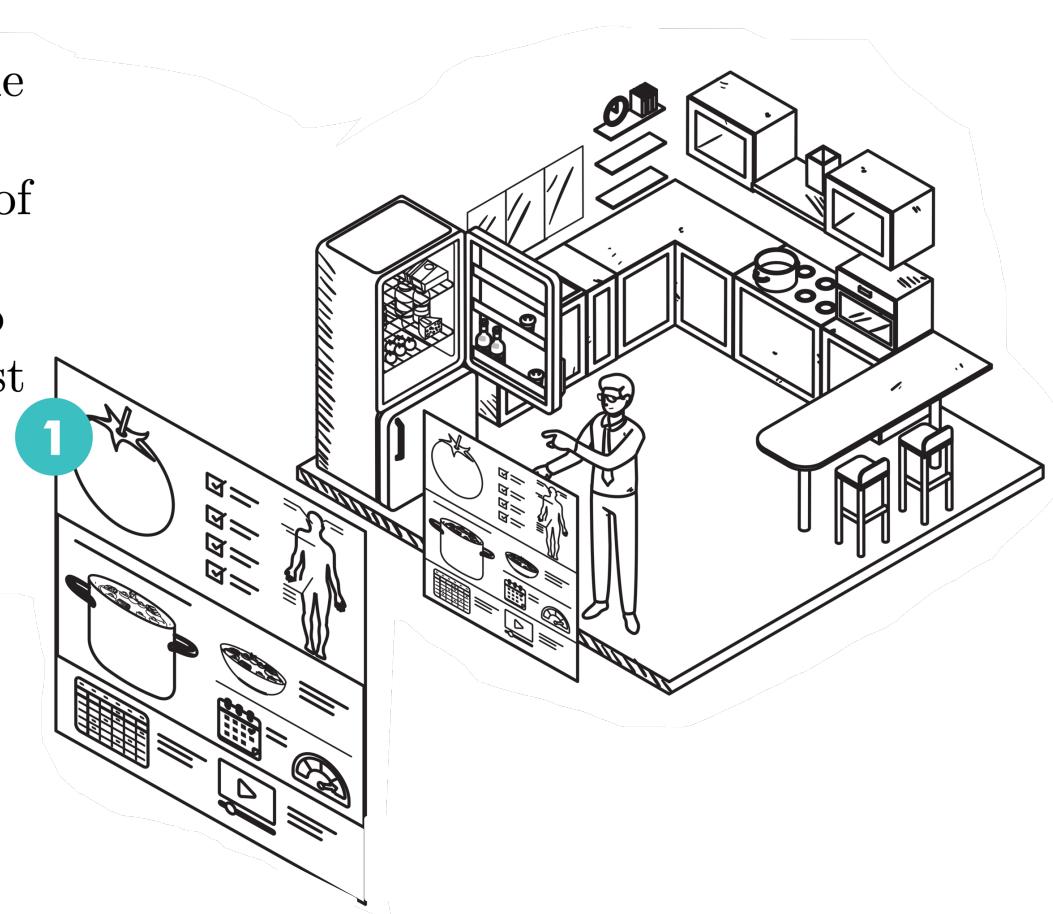


**Ego-Police**

# EGO-Home

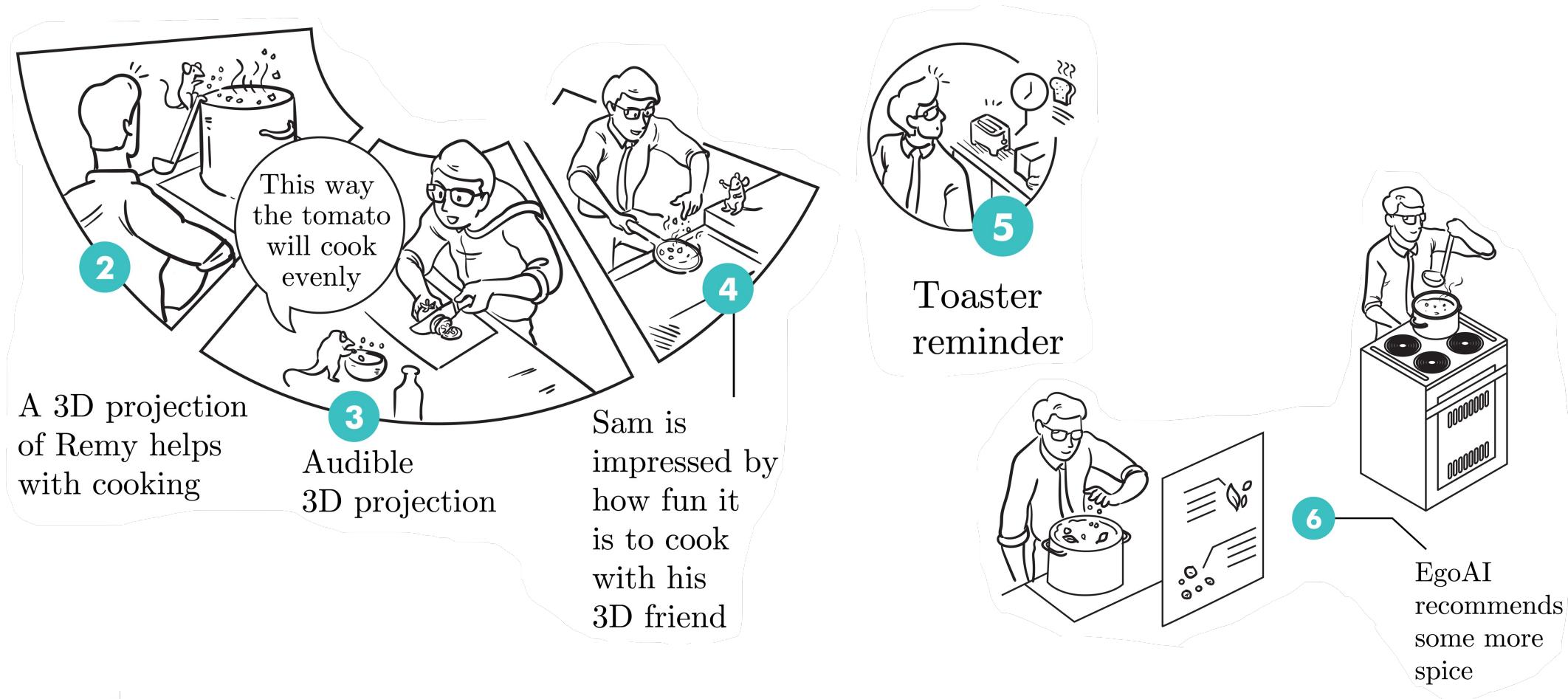
with: Chiara Plizzari, Gabriele Goletto, Antonio Furnari, Siddhant Bansal, Francesco Ragusa, Giovanni Maria Farinella, Tatiana Tommasi

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



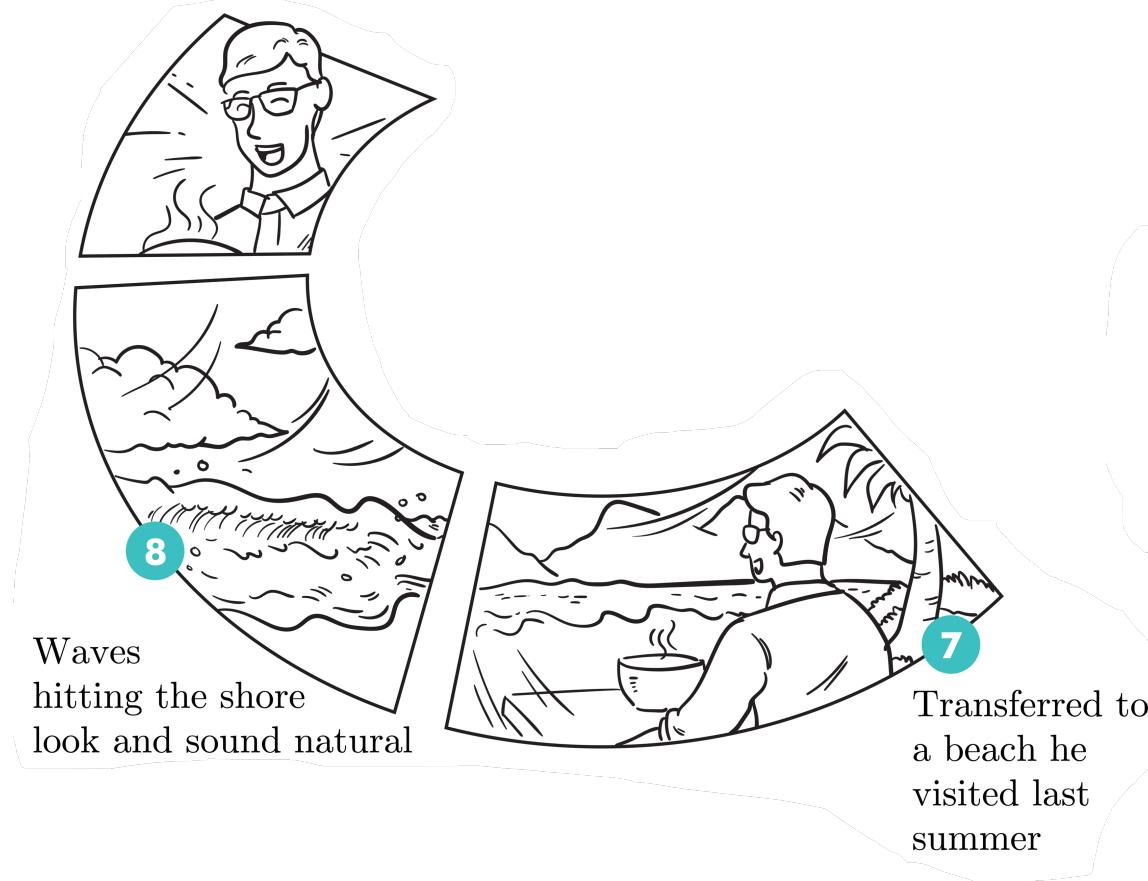
# EGO-Home

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



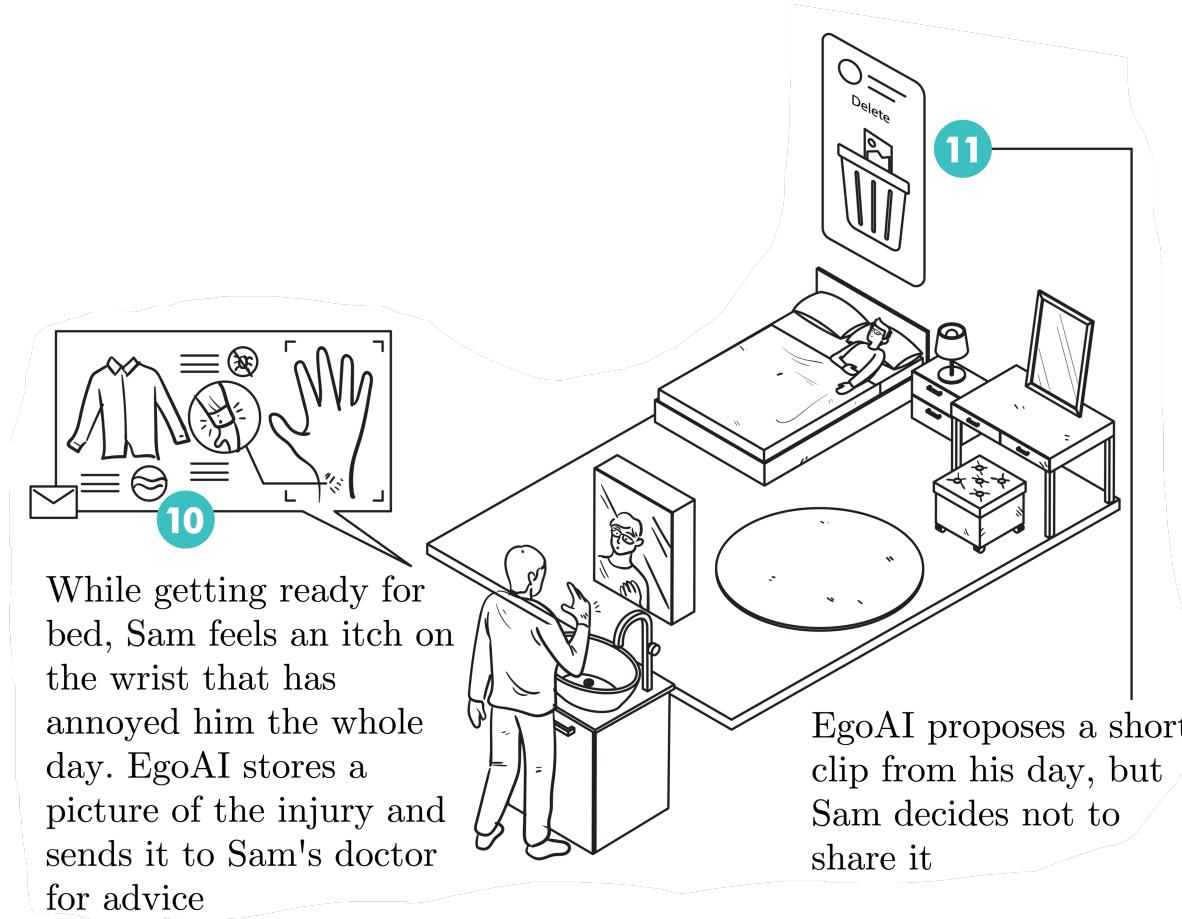
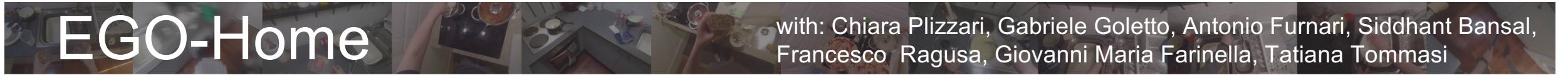
# EGO-Home

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

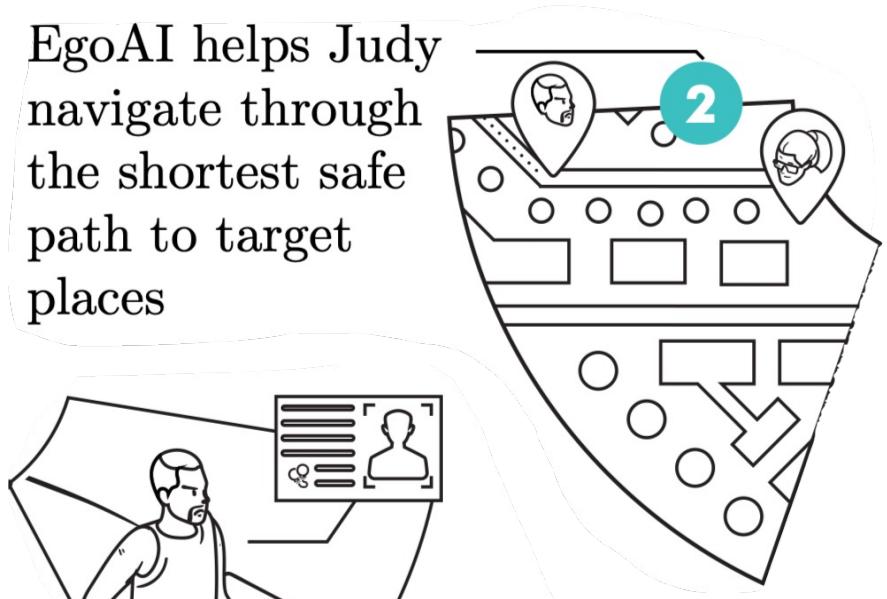




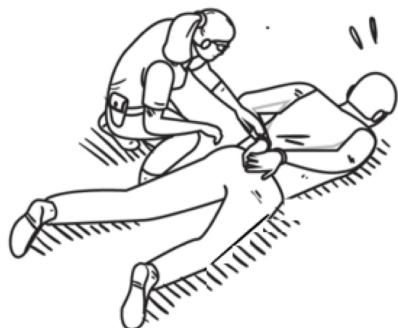
# From Stories to Tasks

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**

**Localisation and Navigation**

1 2

**Messaging**

1 3 11

**Action Recognition**

2 13

**Person Re-ID**

2 4

**Object Detection and Retrieval**

7

**Measuring System**

8 9

**Decision Making**

9

**3D Scene Understanding**

10

**Hand-Object Interaction**

12

**Summarisation**

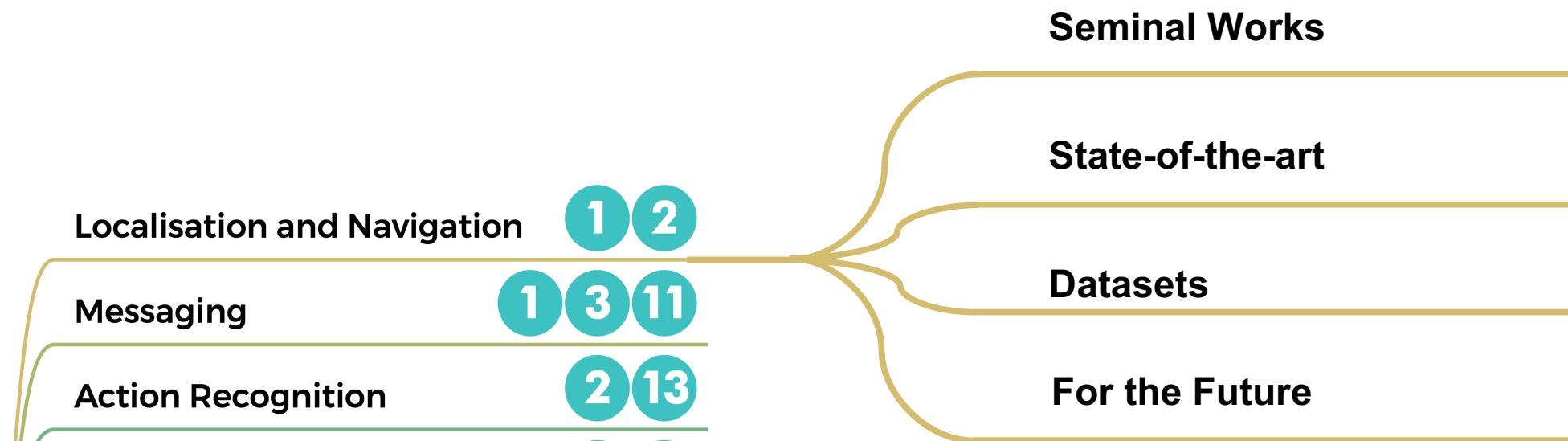
13

**Privacy**

14

# The Survey Part

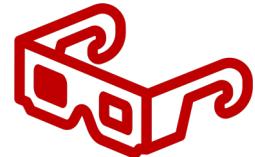
with: Chiara Plizzari, Gabriele Goletto, Antonio Furnari, Siddhant Bansal, Francesco Ragusa, Giovanni Maria Farinella, Tatiana Tommasi



# The Survey Part

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

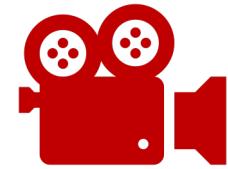
- 12 tasks
- 46 pages (excluding references)
- 462 references



Motivation and Datasets in  
Egocentric Video Understanding



Video Understanding  
Out of the Frame



Video Understanding:  
Data and Tasks



Teaser: The Wizard of Oz  
& Genie 3



Videos are Multimodal



Outlook into the Future of  
Egocentric Vision



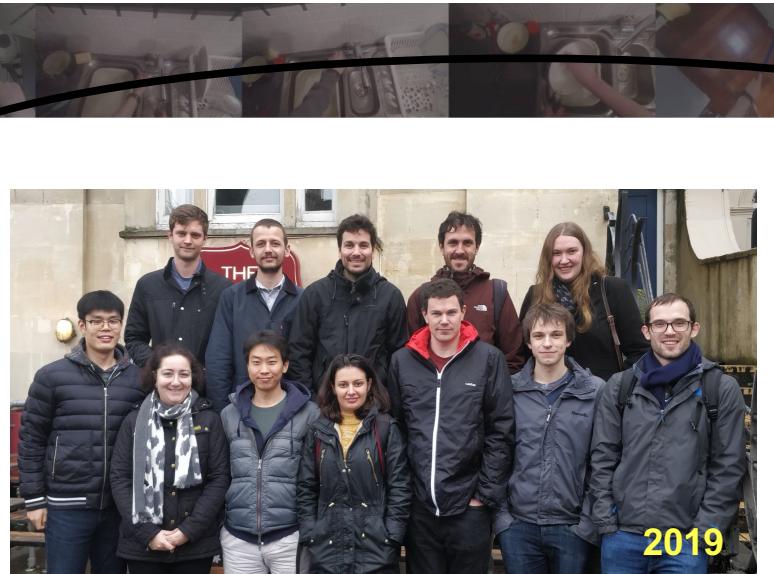
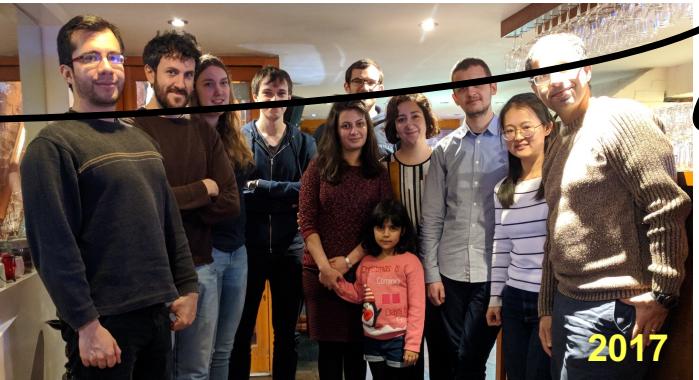
Connected Videos of One's Life



Conclusion

# The Team

grateful

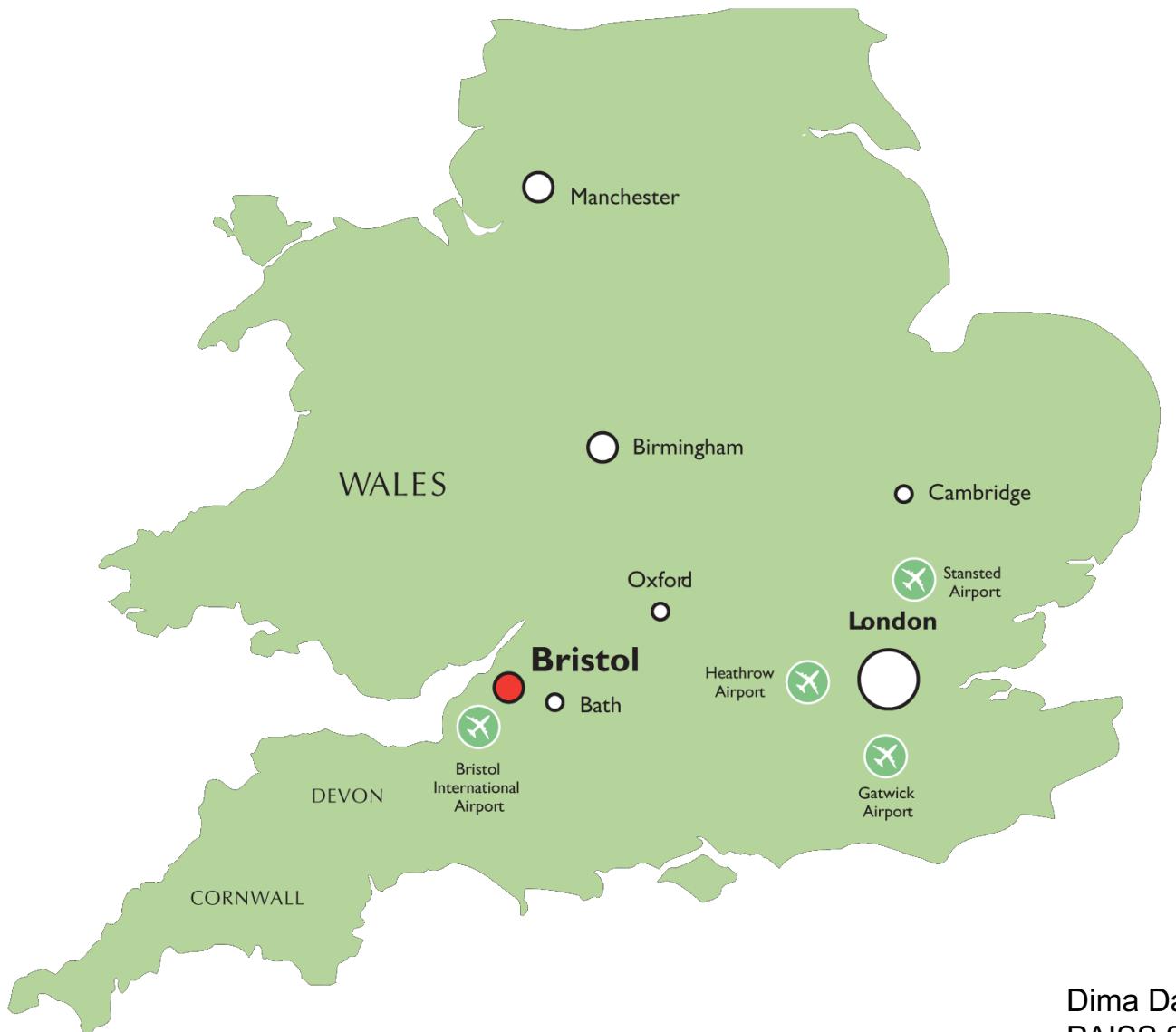




My research team...

*grateful*







# Thank you

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

<http://dimadamen.github.io>



@dimadamen



@dimadamen.bsky.social



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

## Q&A