



# 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



# 360 vs Egocentric...



Ego

360



# 360 vs Egocentric...



# In today's talk...



Motivation and Datasets in  
Egocentric Video Understanding



Teaser: The Wizard of Oz  
& Genie 3



Video Understanding  
Out of the Frame



Outlook into the Future of  
Egocentric Vision



Point Tracking



Object Tracking



Gaze Priming



Hand Tracking



Conclusion

In today's talk...



Motivation and Datasets in  
Egocentric Video Understanding



Teaser: The Wizard of Oz  
& Genie 3



Video Understanding  
Out of the Frame



Outlook into the Future of  
Egocentric Vision



Point Tracking



Object Tracking



Gaze Priming



Hand Tracking

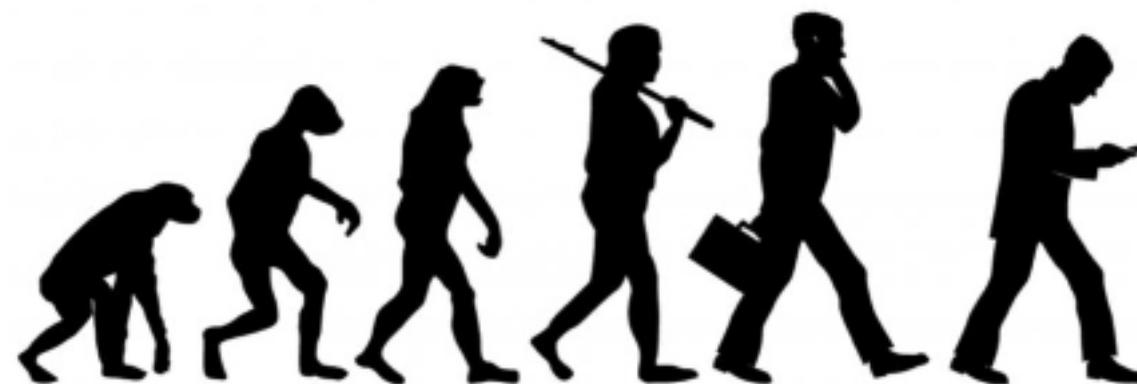


Conclusion

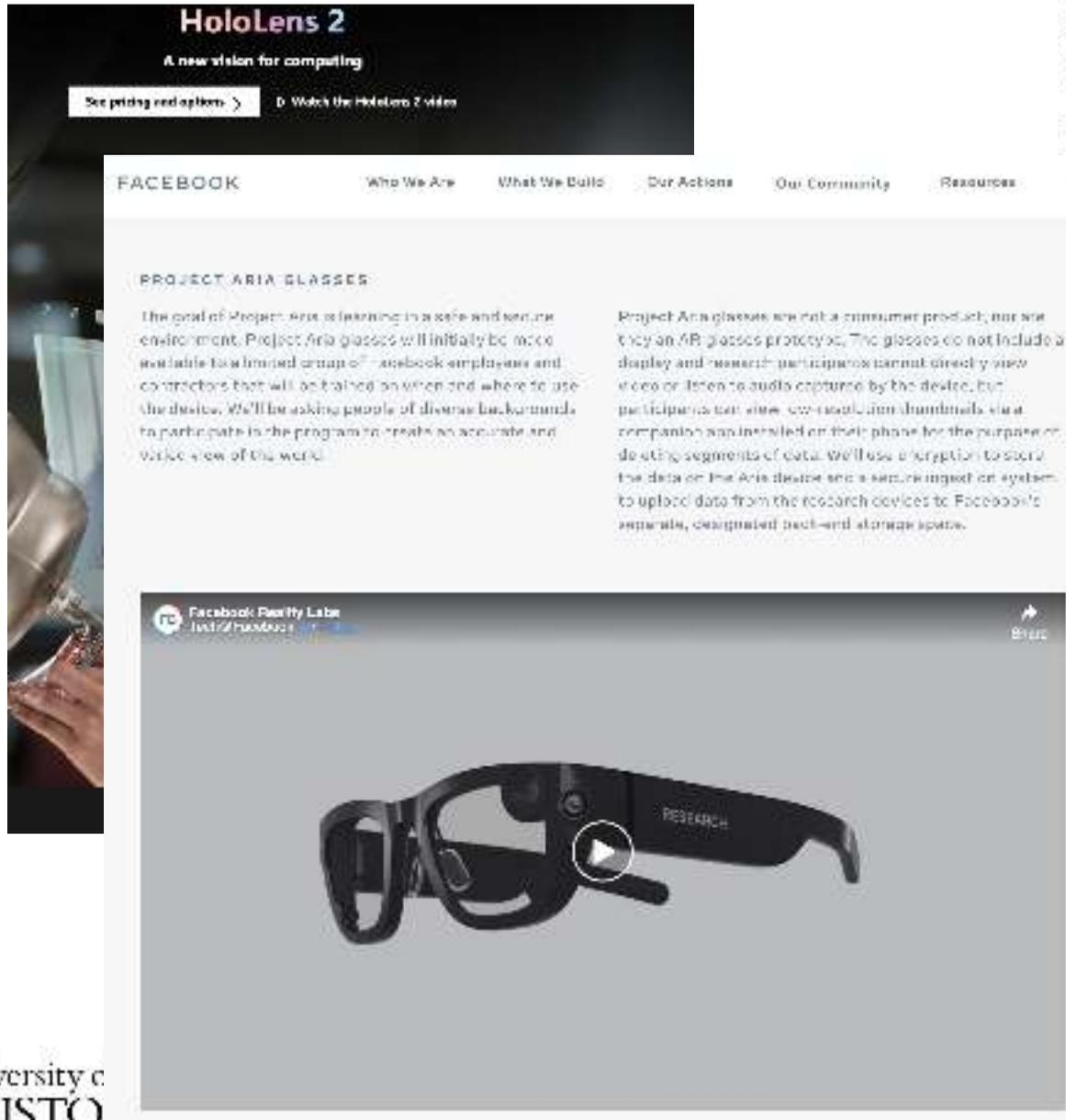
# The present...



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 diverse 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 cloud system to upload data from the research devices to Facebook's separate, designated back-end storage space.

Facebook Faculty Lab  
1/17/2018



## Samsung patent application reveals augmented reality headset design

It comes as the Gear VR slowly fades away

By Jan Porter  
C



# The future is here...



The future can be imagined...



# Egocentric Videos?



In today's talk...



Motivation and Datasets in  
Egocentric Video Understanding



Teaser: The Wizard of Oz  
& Genie 3



Video Understanding  
Out of the Frame



Outlook into the Future of  
Egocentric Vision



Point Tracking



Object Tracking



Gaze Priming



Hand Tracking



Conclusion



**EPIC-KITCHENS**

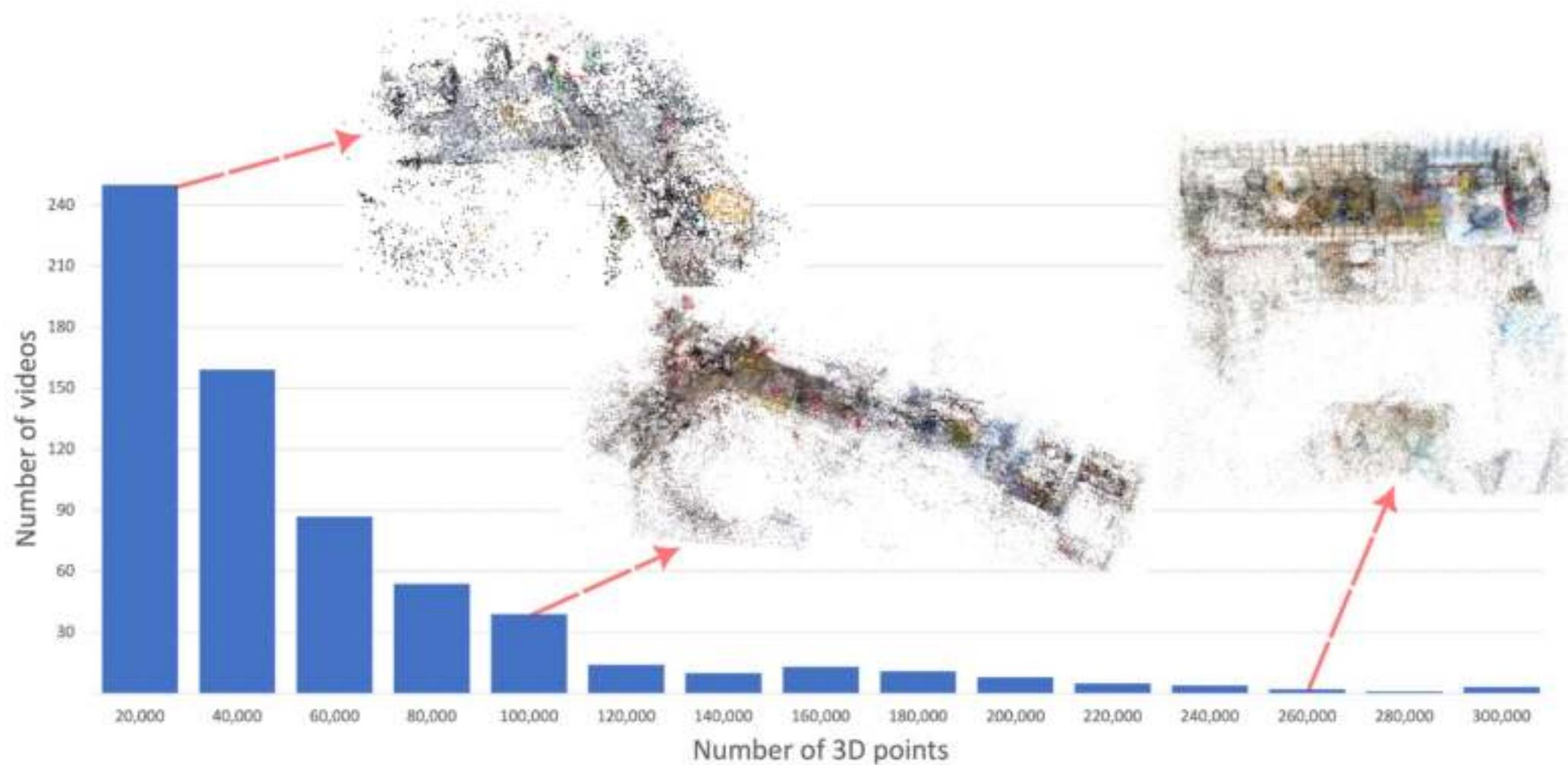


Figure 4: **Number of 3D points histogram.** The majority of our reconstructions generate less than 40,000 points that are enough to represent the kitchen. However, some reconstructions have more than 100,000, we include the point clouds for each points range showing the fine details covered by having 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>	<b>50</b>	<b>6–37 min (Avg 22)</b>	<b>50 / 50</b>	<b>✓</b>

# In today's talk...



Motivation and Datasets in  
Egocentric Video Understanding



Teaser: The Wizard of Oz  
& Genie 3



Video Understanding  
Out of the Frame



Outlook into the Future of  
Egocentric Vision



Point Tracking



Object Tracking



Gaze Priming



Hand Tracking



Conclusion



# 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

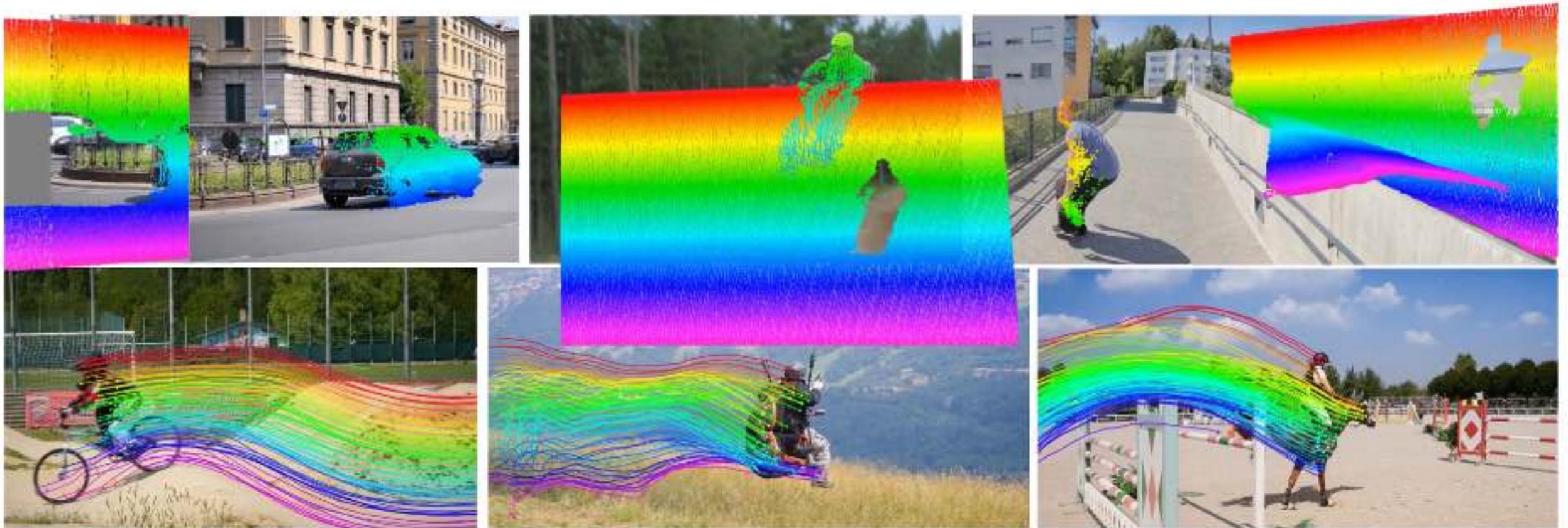


Dima Damen  
BinEgo-360 Workshop @ICCV2025

# What is point tracking?

with: Ahmad Darkhalil  
Rhodri Guerrier  
Adam W Harley

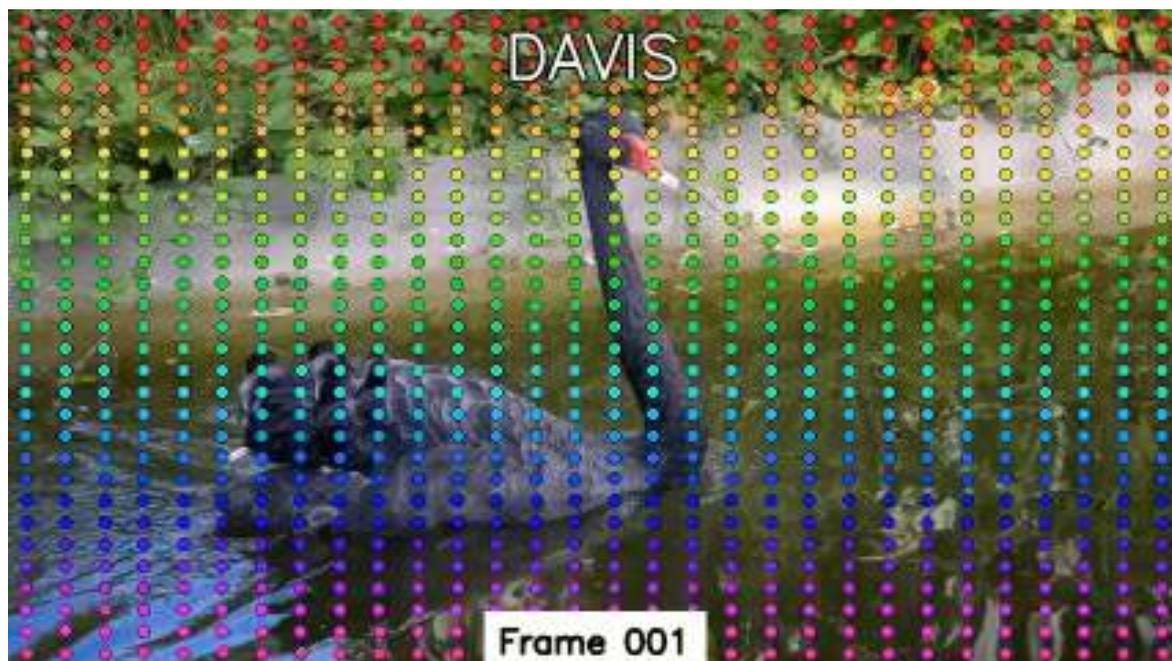
- Given: Query points in one frame
- Track these points throughout the video



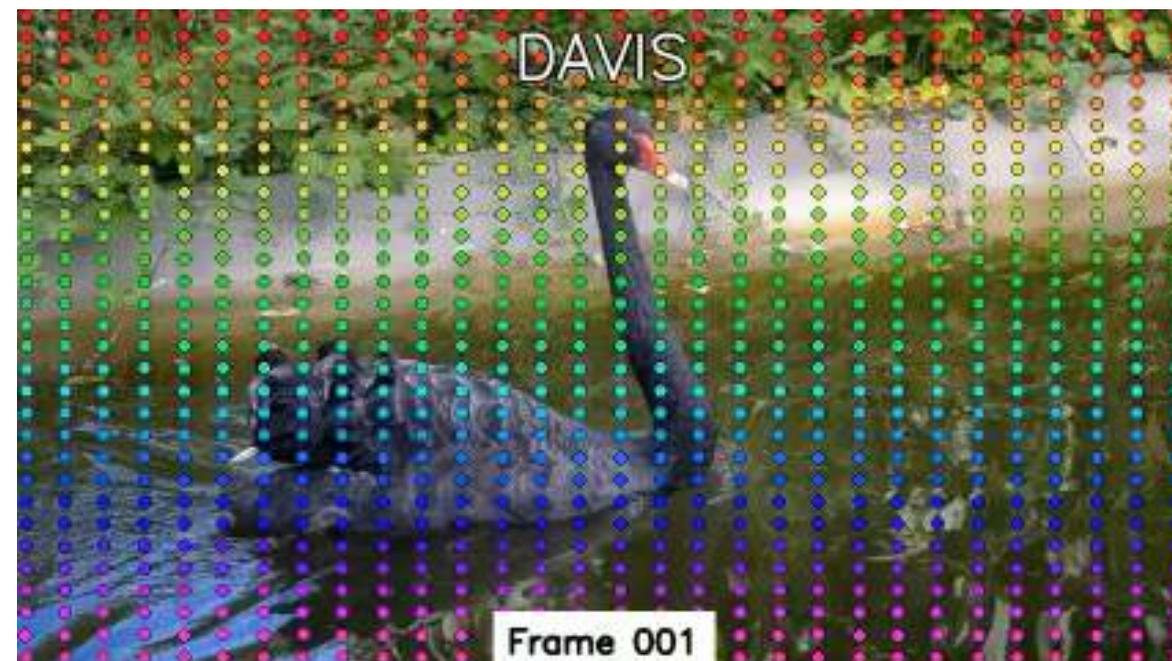
# SOTA on current benchmarks

with: Ahmad Darkhalil  
Rhodri Guerrier  
Adam W Harley

## LocoTrack



## CoTracker3

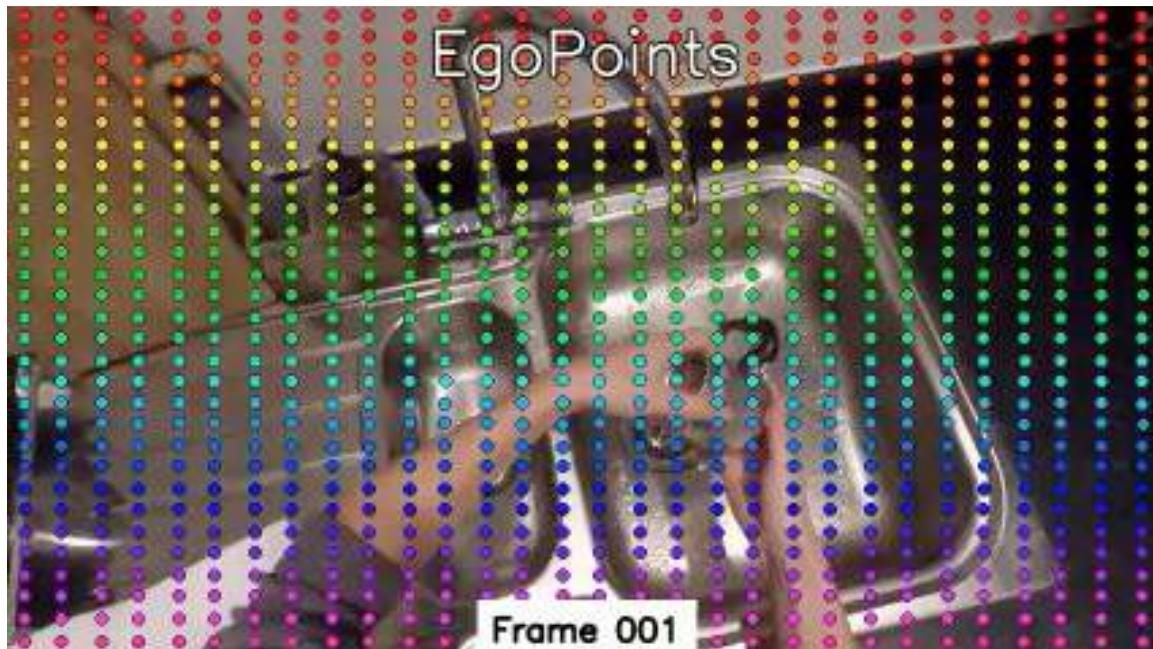


# Current Models Struggle with Egocentric Videos

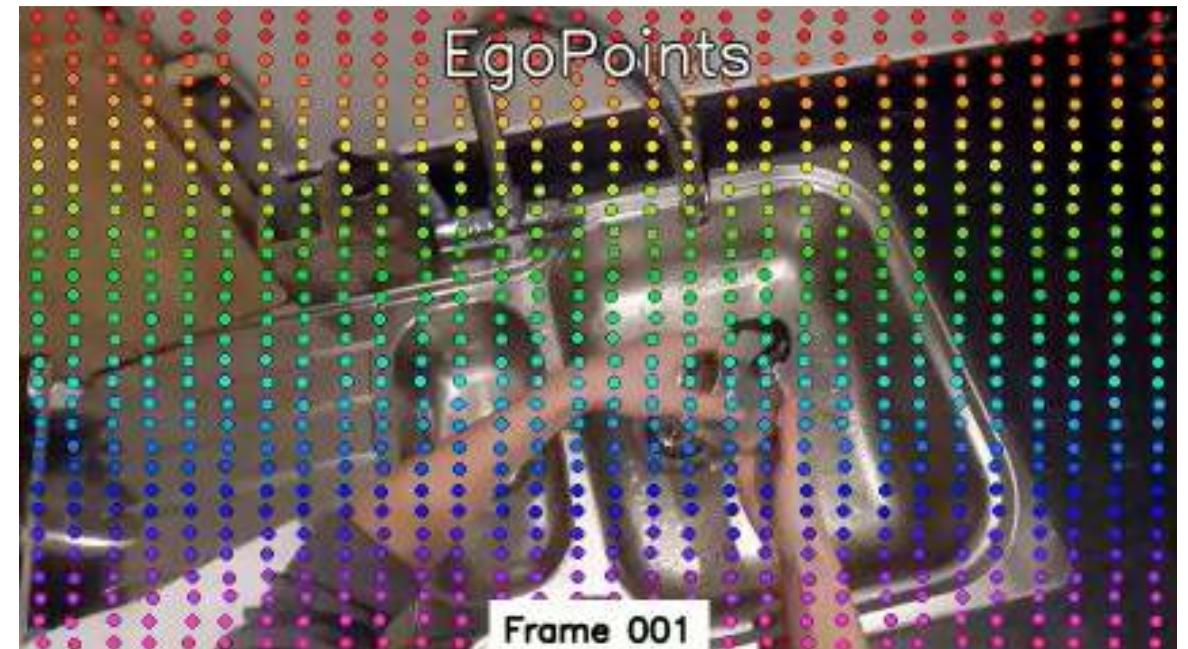
with: Ahmad Darkhalil  
Rhodri Guerrier  
Adam W Harley

- Head motion and motion blur
- Frequent re-identification

LocoTrack



CoTracker3



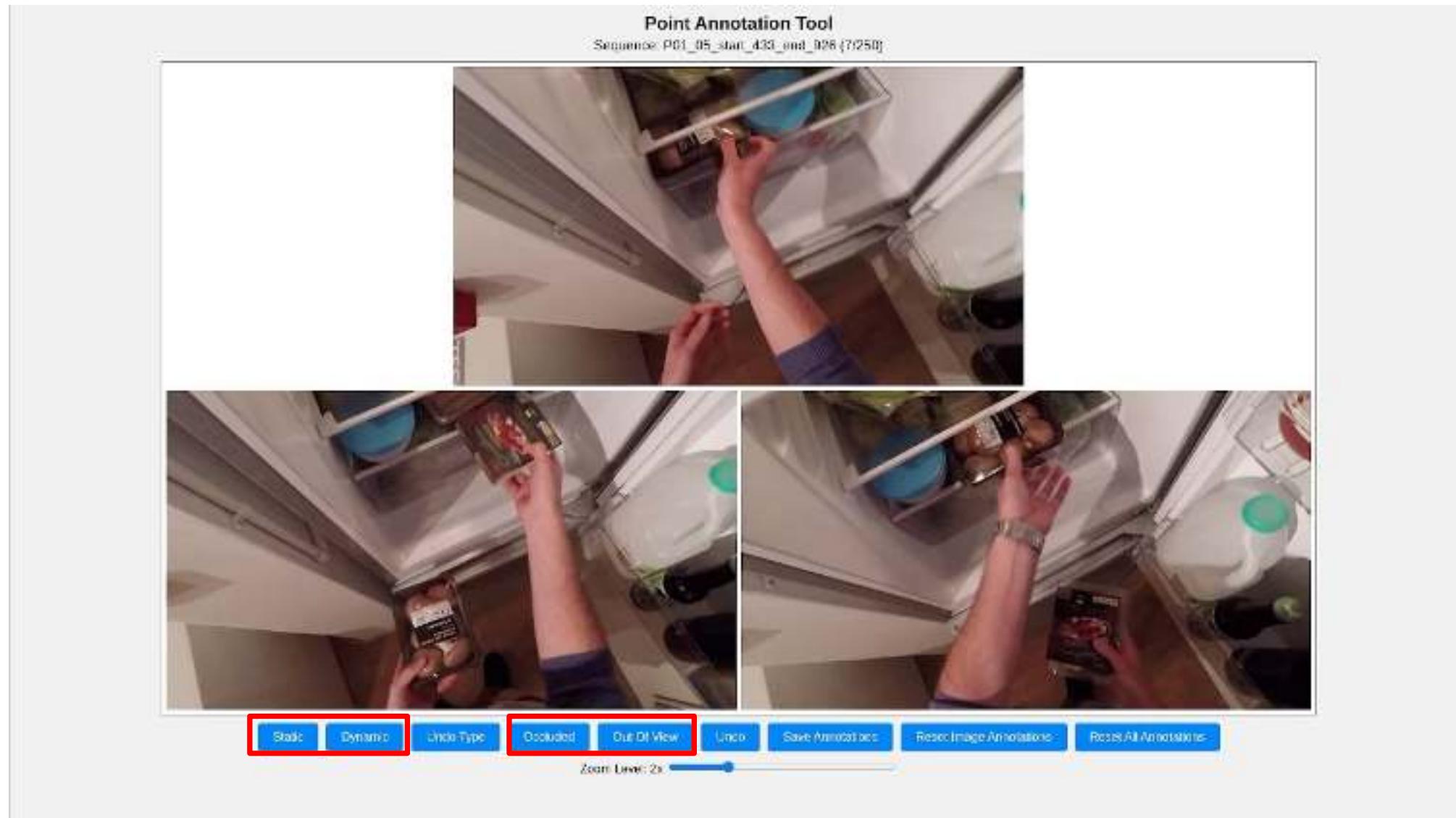
# Main Contributions

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



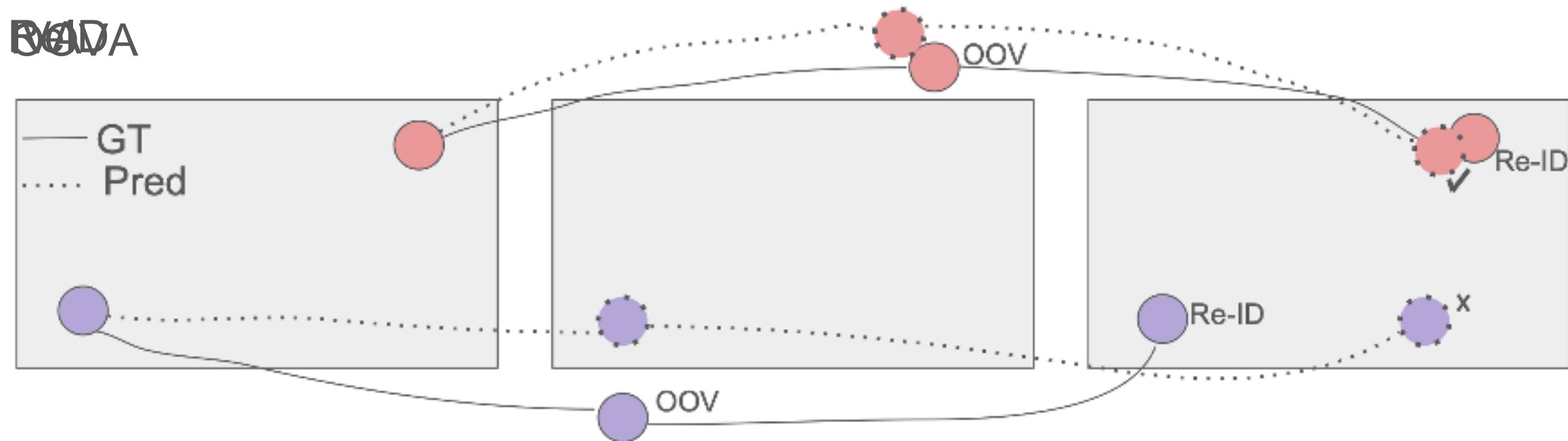
# EgoPoints Benchmark

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

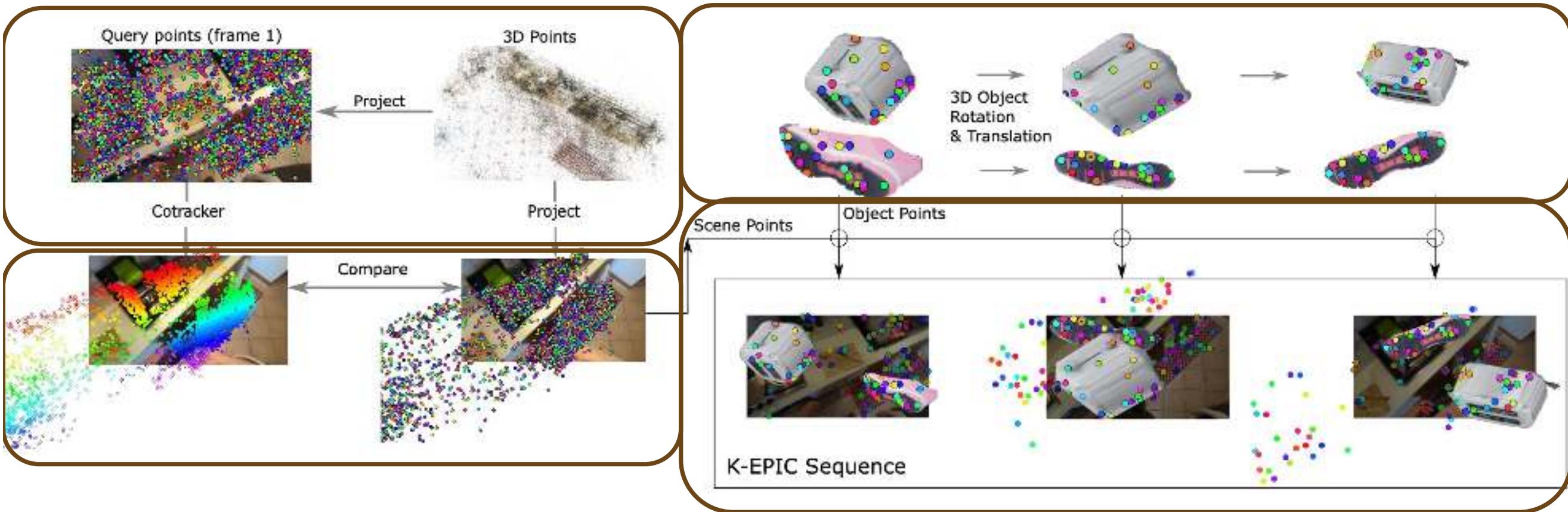
# SOTA Models Struggle on EgoPoints

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 $\uparrow$	IVA $\uparrow$
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

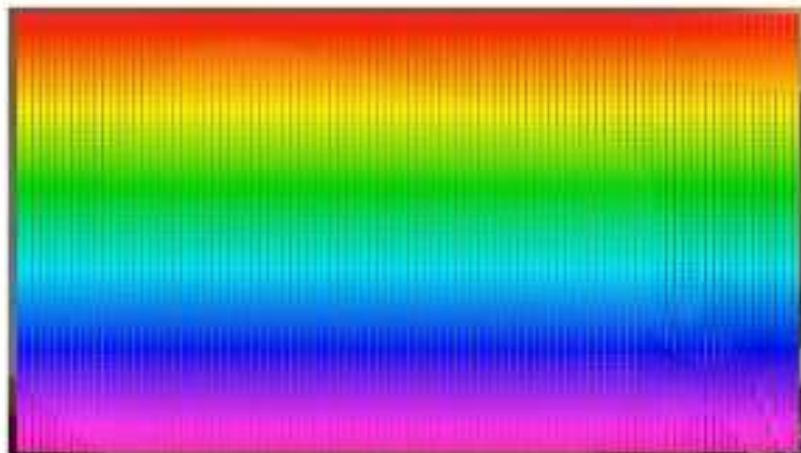
with: Ahmad Darkhalil  
Rhodri Guerrier  
Adam W Harley

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
<u>PIPs++ w. K-EPIC FT (scene points only)</u>	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
<u>CoTracker w. K-EPIC FT (scene points only)</u>	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>

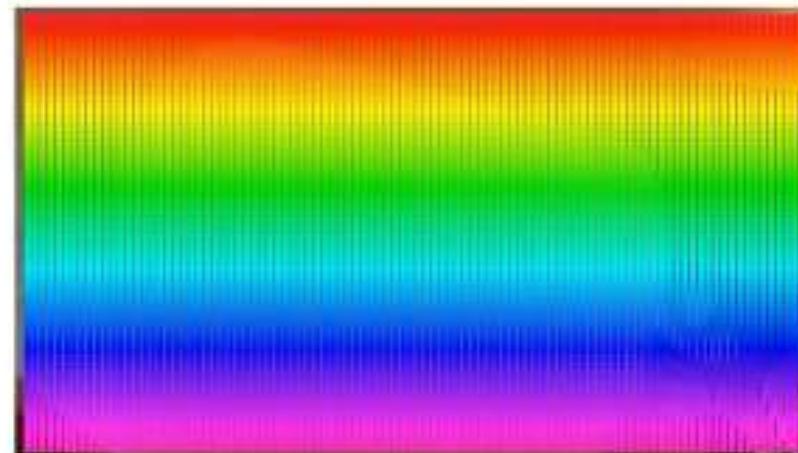
# Qualitative Examples of CoTracker

with: Ahmad Darkhalil  
Rhodri Guerrier  
Adam W Harley

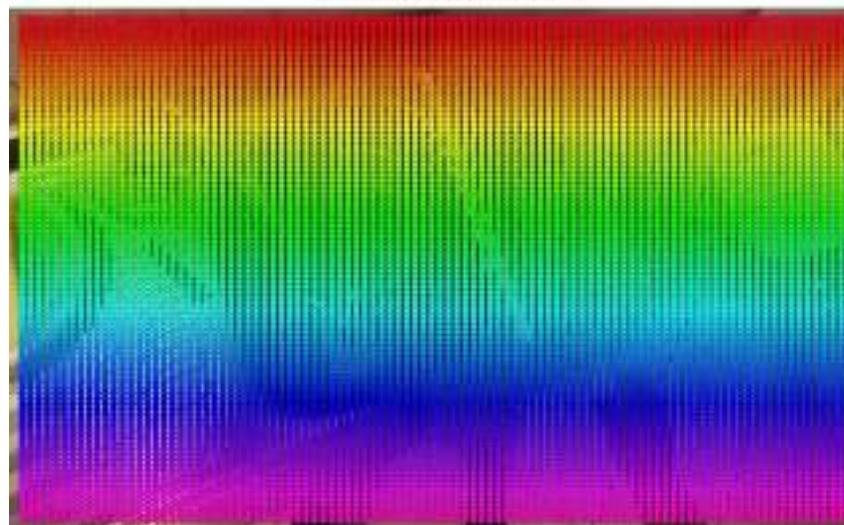
Cotracker Baseline



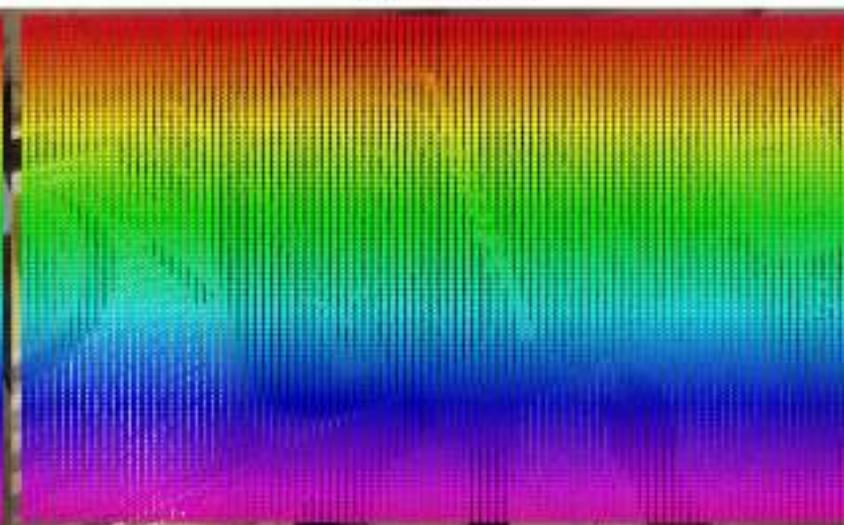
Cotracker FT



Cotracker Baseline



Cotracker FT



# 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	47.5	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>

# In today's talk...



Motivation and Datasets in  
Egocentric Video Understanding



Teaser: The Wizard of Oz  
& Genie 3



Video Understanding  
Out of the Frame



Outlook into the Future of  
Egocentric Vision



Point Tracking



Object Tracking



Gaze Priming



Hand Tracking



Conclusion

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

Chiara Plizzari

Shubham Goel

Toby Perrett

Jacob Chalk

Angjoo Kanazawa

Dima Damen

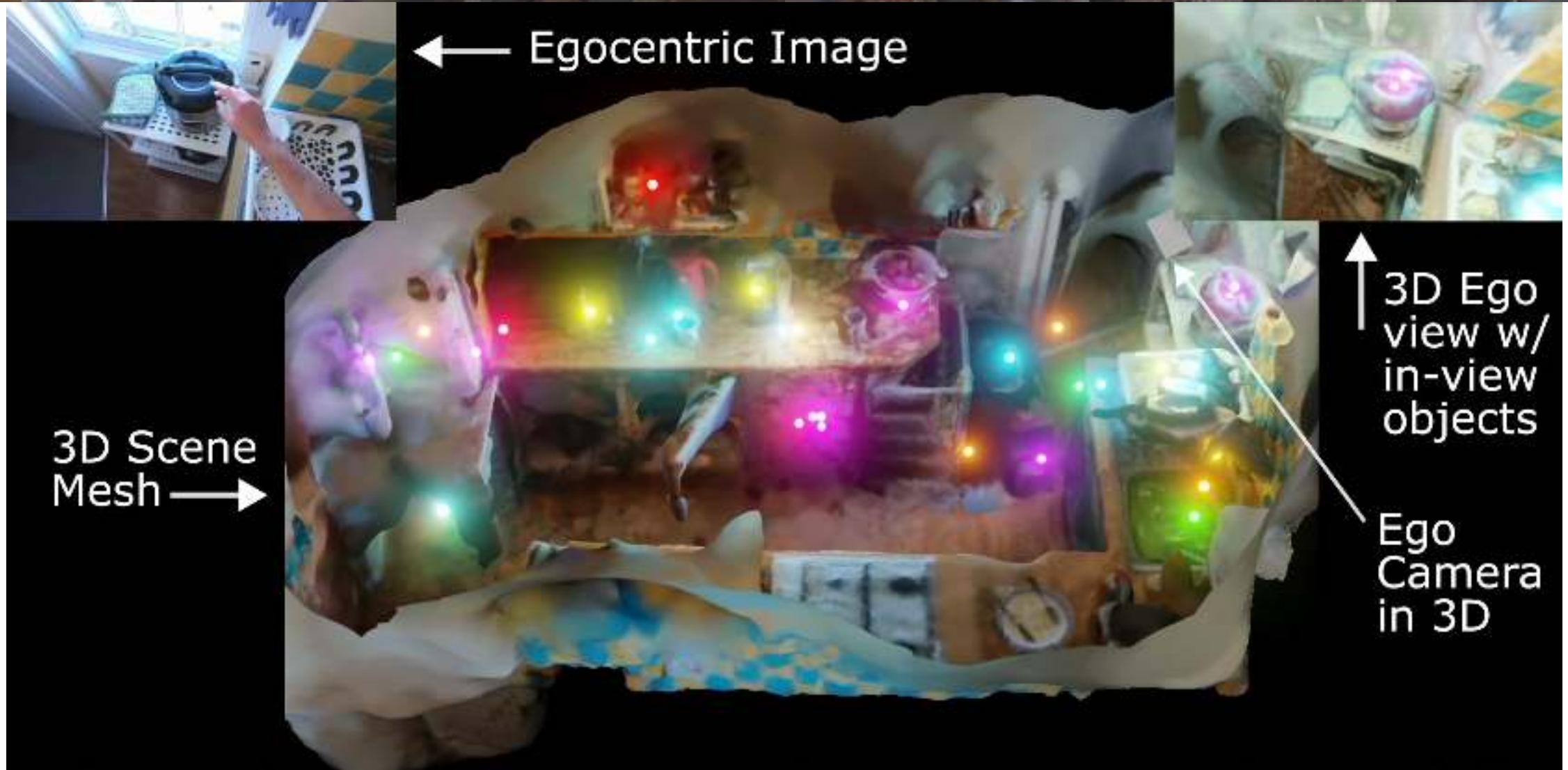
<http://dimadamen.github.io/OSNOM>



# 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



← Egocentric Image



3D Scene Mesh →

↑ 3D Ego view w/  
in-view  
objects

Ego  
Camera  
in 3D

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

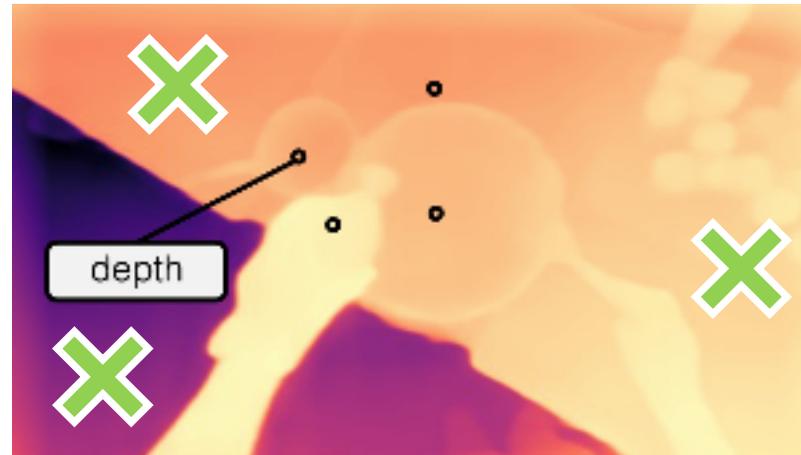
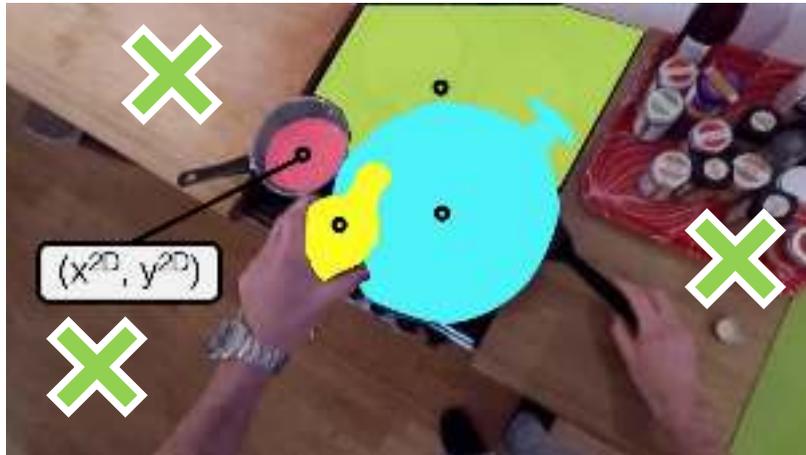
with: Chiara Plizzari  
Toby Perrett

Shubham Goel  
Angjoo Kanazawa

Lift

Match

Keep



0.0 ... 1.0

0.3m ... 1.8m

# Out of Sight, not Out of Mind

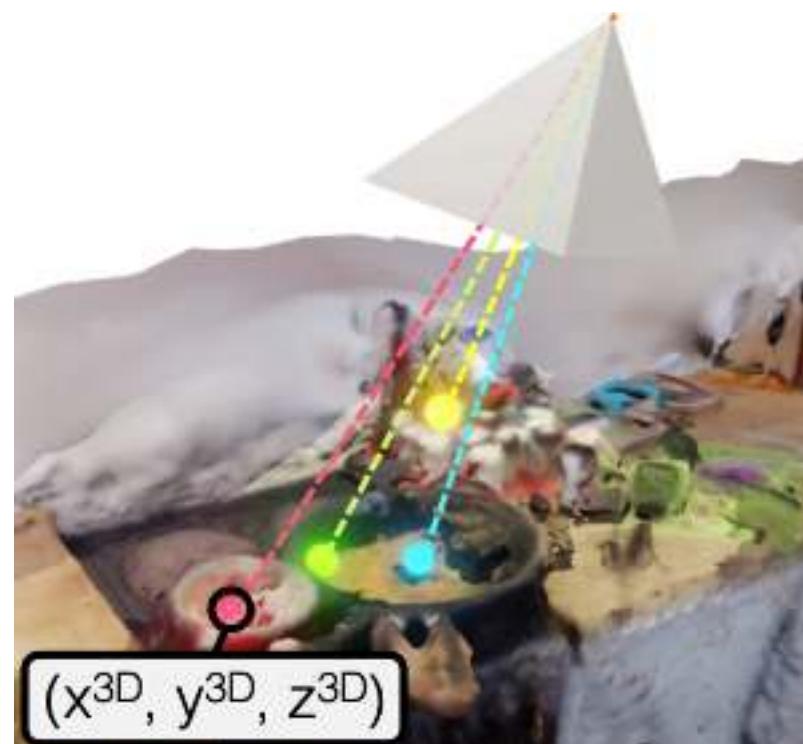
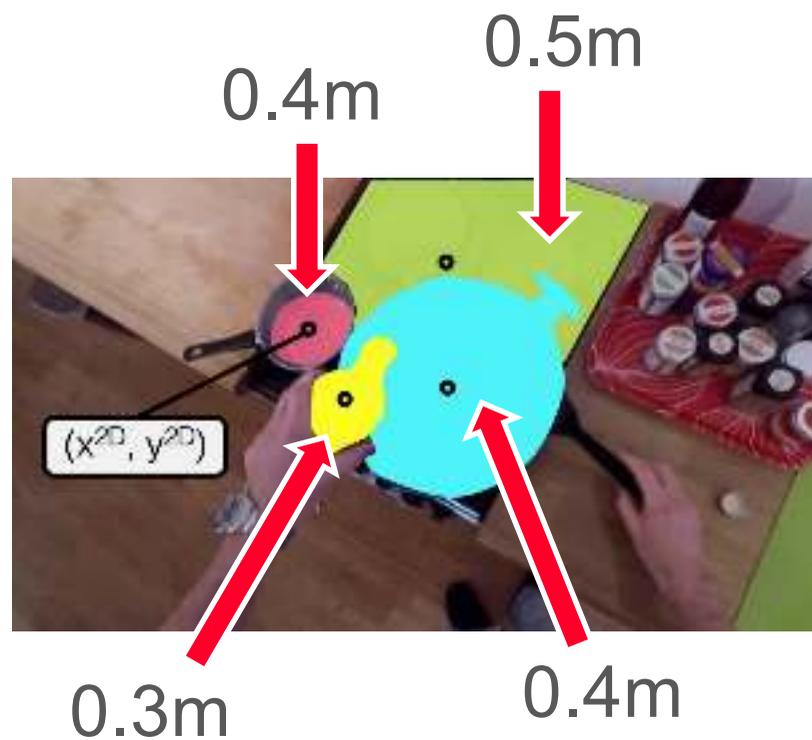
with: Chiara Plizzari  
Toby Perrett

Shubham Goel  
Angjoo Kanazawa

Lift

Match

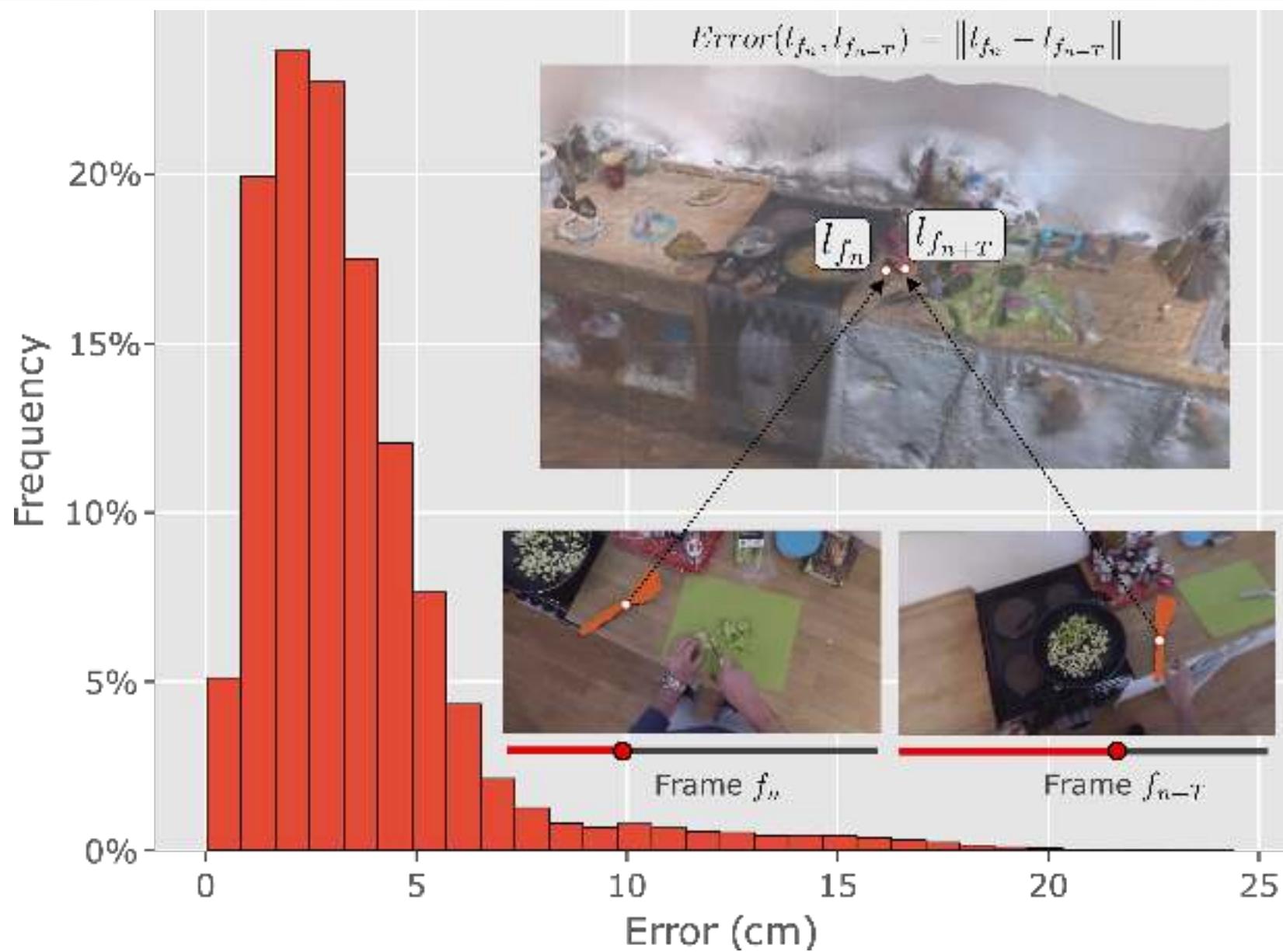
Keep



# 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



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)



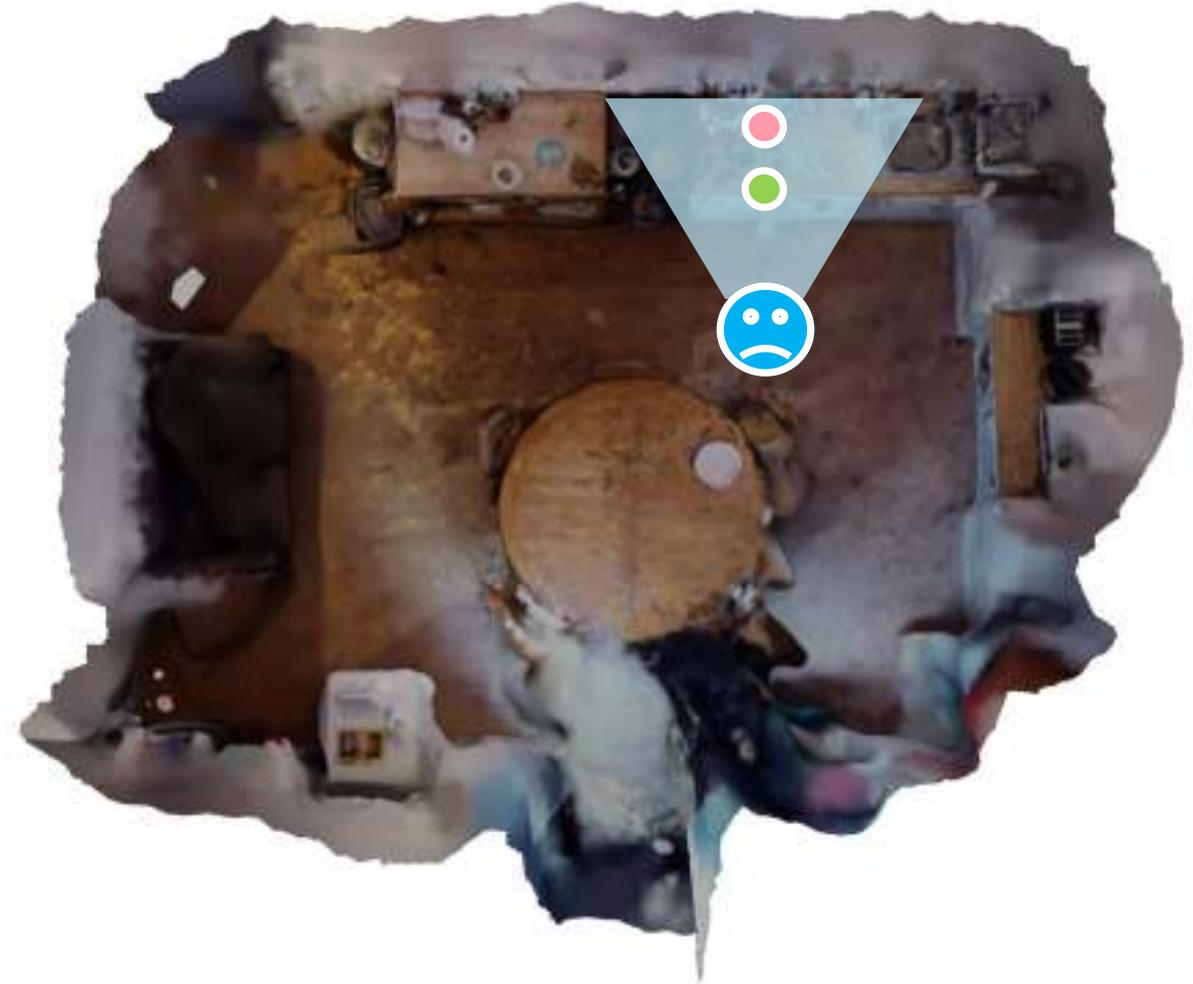
# Out of Sight, not Out of Mind

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)



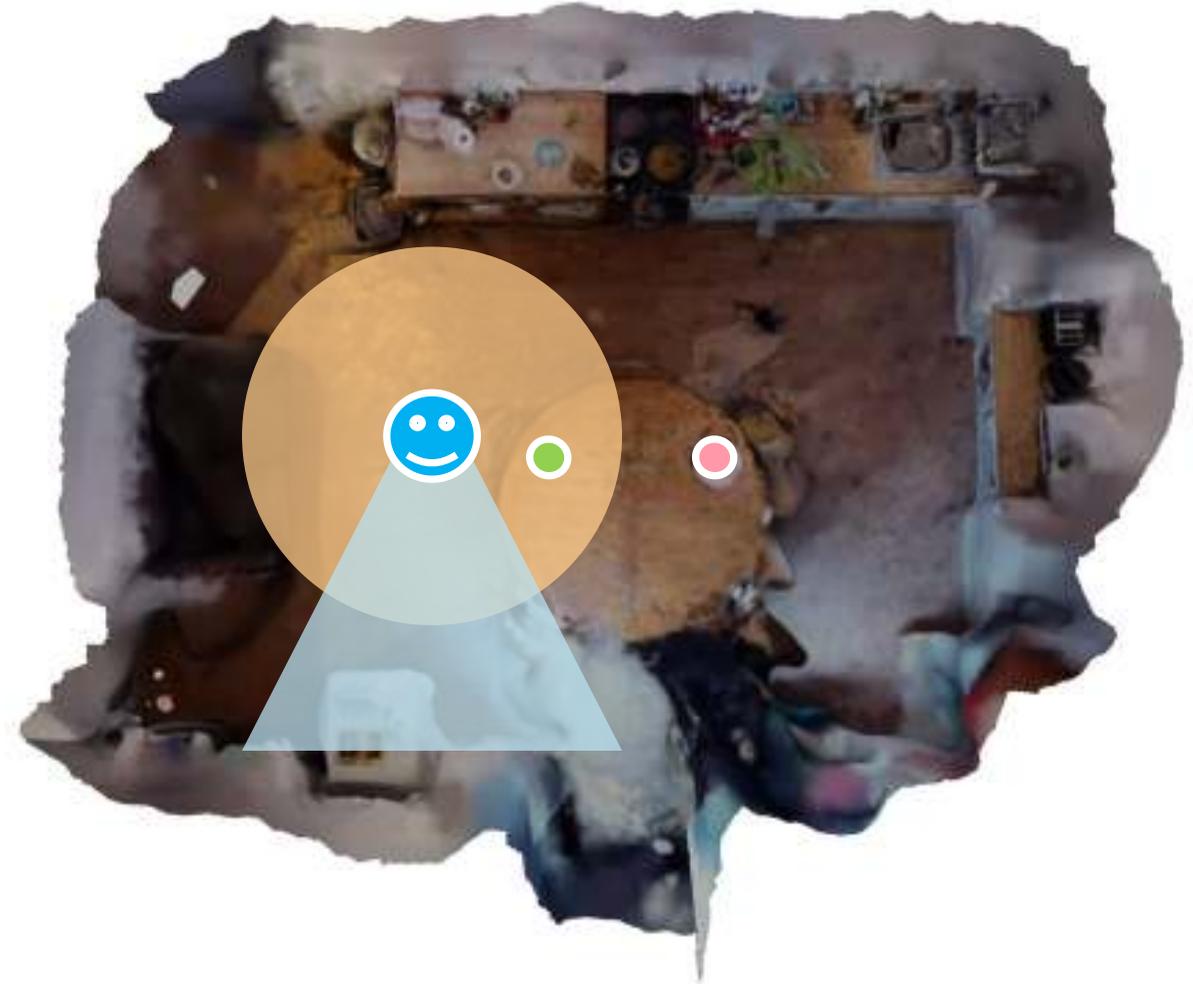
# Out of Sight, not Out of Mind

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

Roby Perrett

Jacob Chalk

Angi

Dima Damen

**Ground-Truth??**

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



Politecnico  
di Torino

Berkeley  
UNIVERSITY OF CALIFORNIA



University of  
BRISTOL



# 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

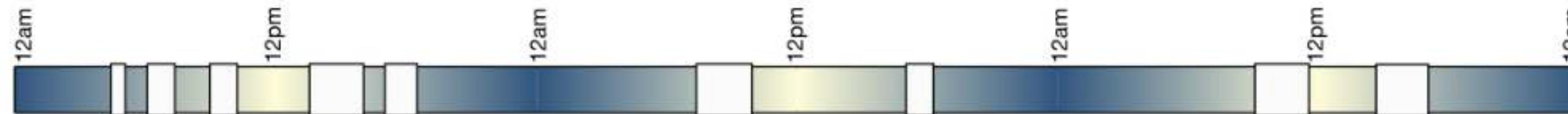


# HD-EPIC





# HD-EPIC





# HD-EPIC



**Recorded over 3 days**



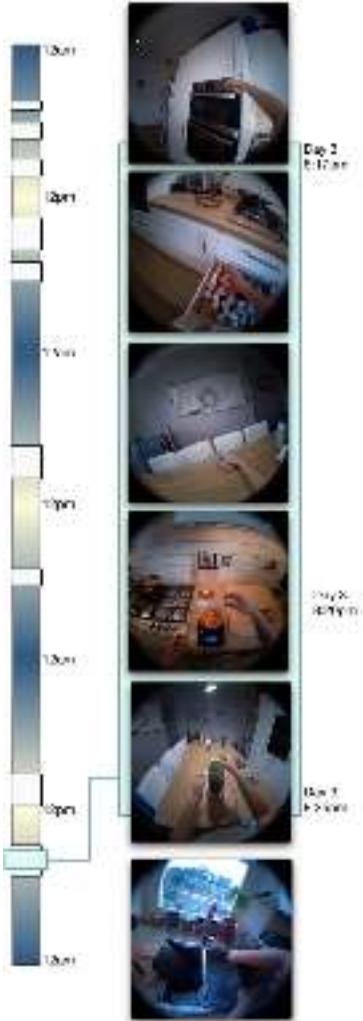


# HD-EPIC





# HD-EPIC





## Recipe: Southwestern Salad

**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**

Day 3  
1:17pm



## Cacio e Pepe (modified)

Ingredients:

200 g

400g of pasta of your choice

(we recommend bucatini)

2 tablespoon of black peppercorn

30 g

200g of freshly grated pecorino cheese

+25g of slightly salted butter

penne

parmigiano



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. <sup>on very low heat</sup>

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





- 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



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

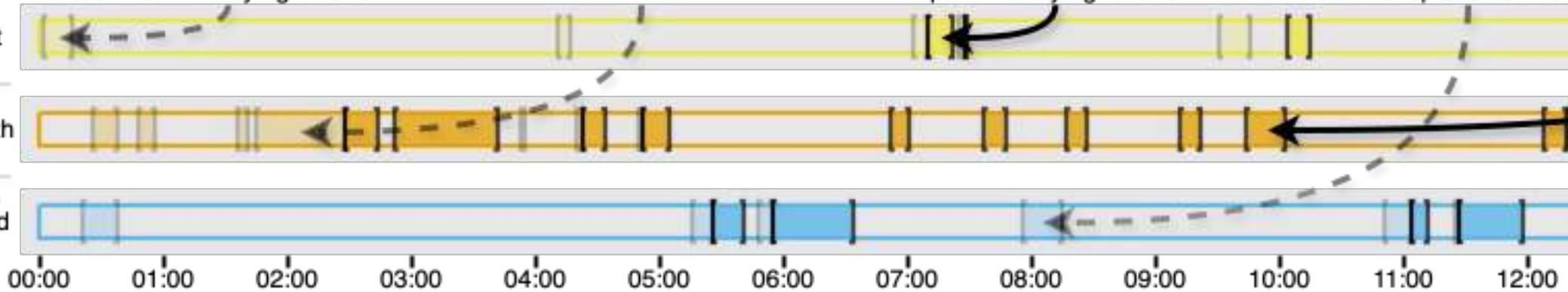


pick up block of cheese the from the top shelf of the

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

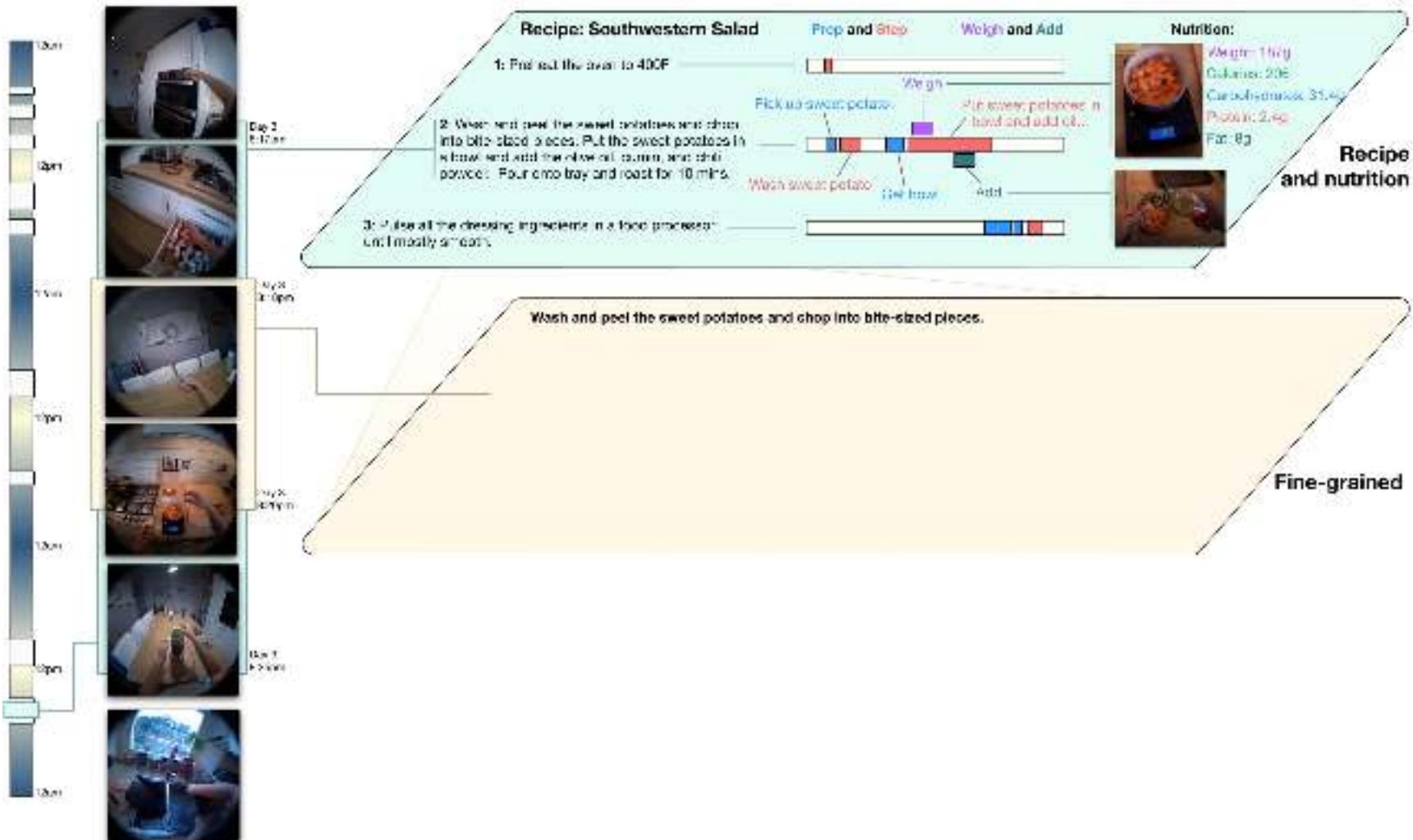
Slice the bacon and place in a non-stick frying pan on a medium heat with half a tablespoon of olive oil and ...

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



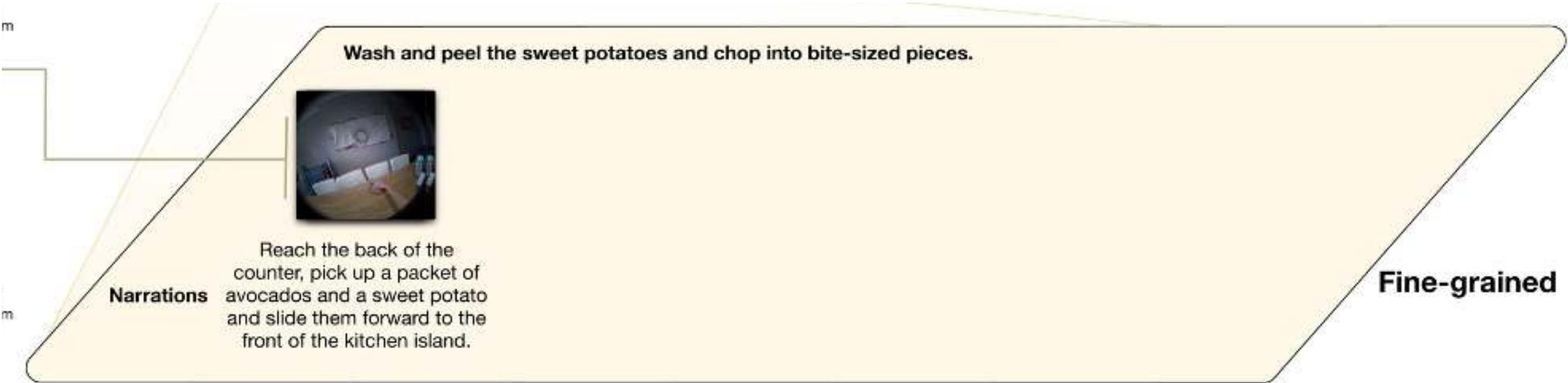


# HD-EPIC





# HD-EPIC



# Highly-Detailed Narrations





# HD-EPIC

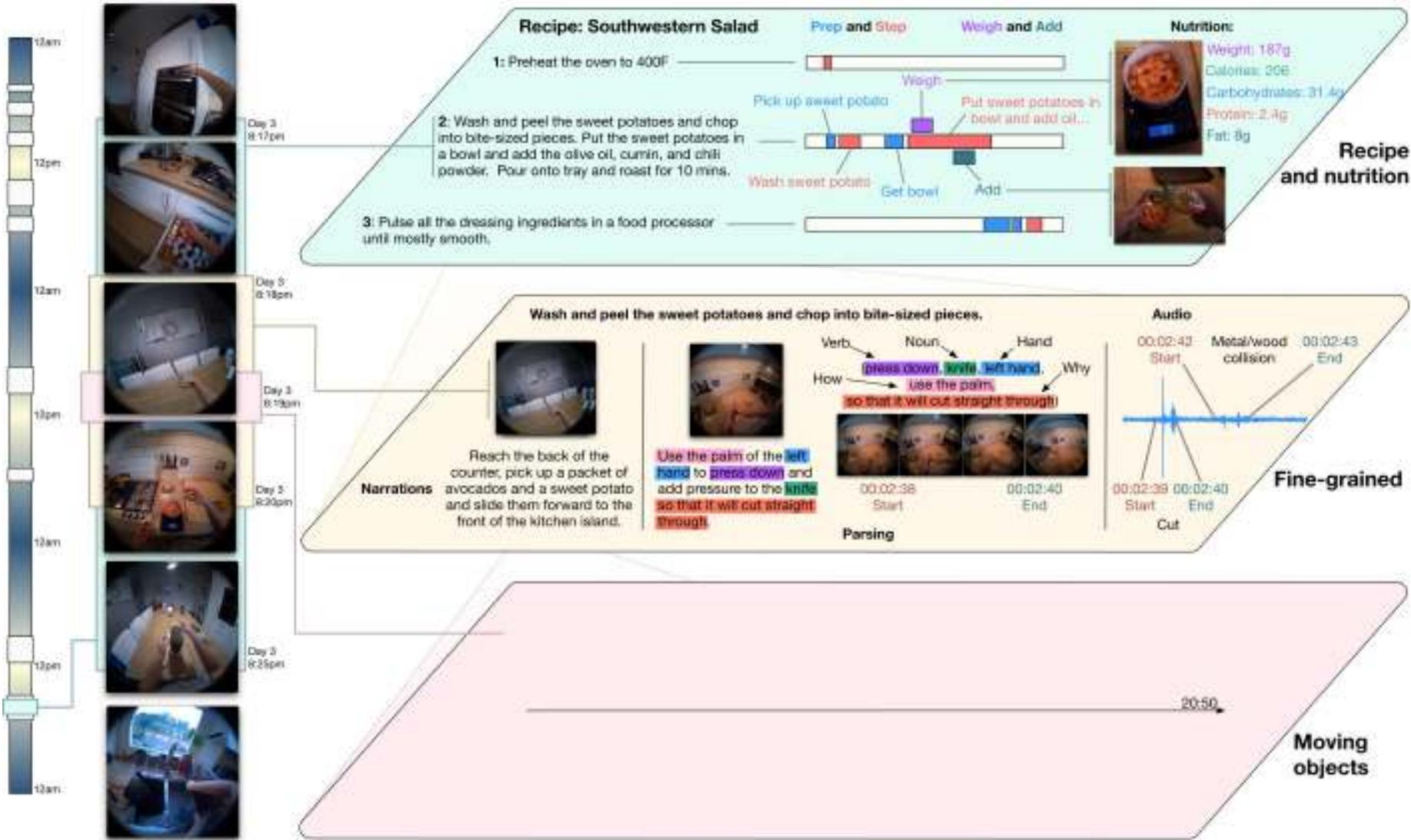


- 59,454 fine-grained actions, with a mean duration of 2.0s ( $\pm 3.4$ s).



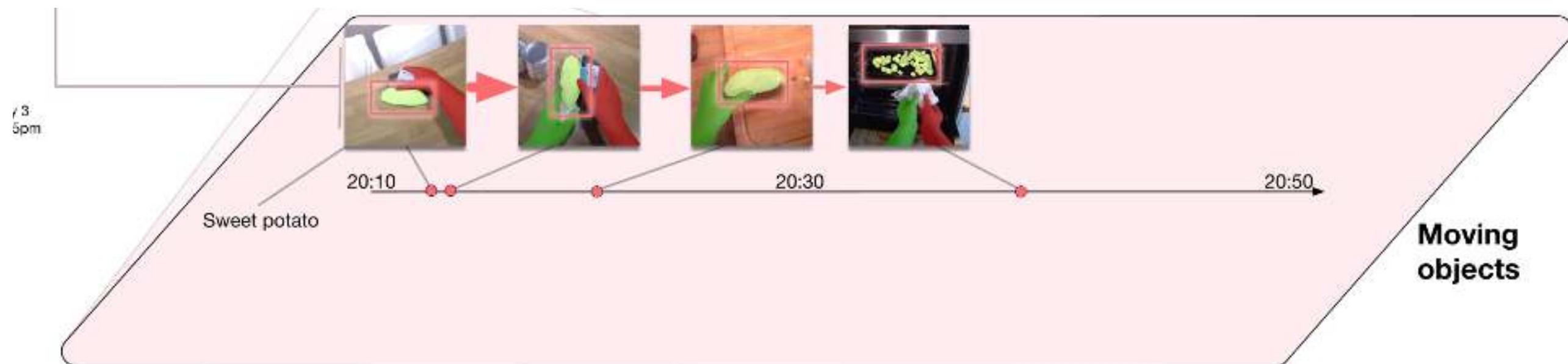


# HD-EPIC





# HD-EPIC

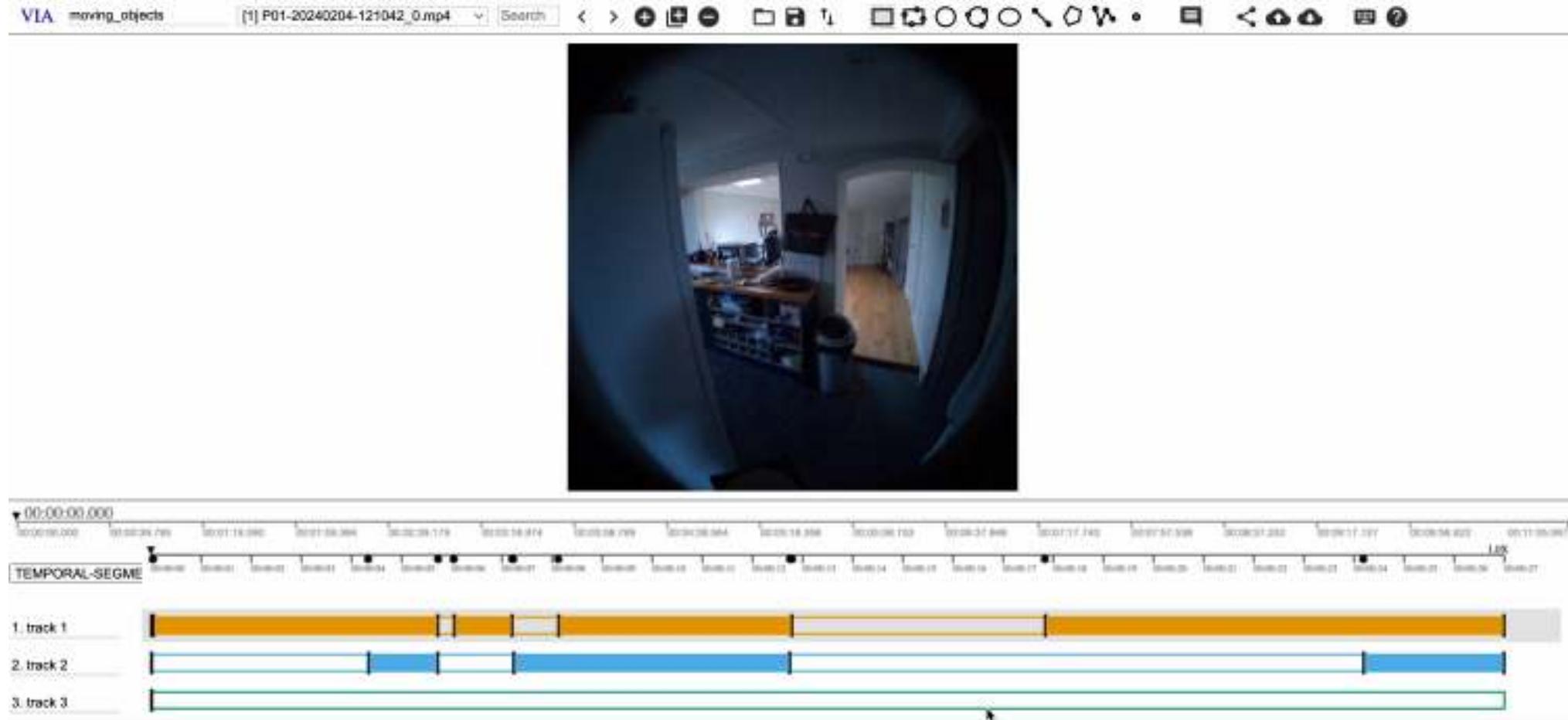




# HD-EPIC



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

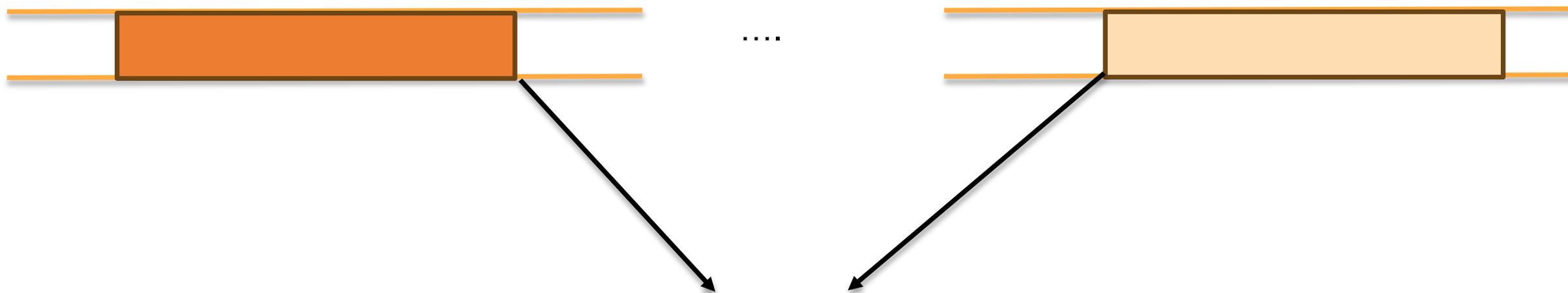




# HD-EPIC



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



Using appearance & 3D location information to match  
Manual confirmation in cases of confusion...





## Current Track Image

Choose Files | 201 files

04:22 04:51

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) Add

plastic chopping board (0.3m) Add

metal cooling rack (0.6m) Add

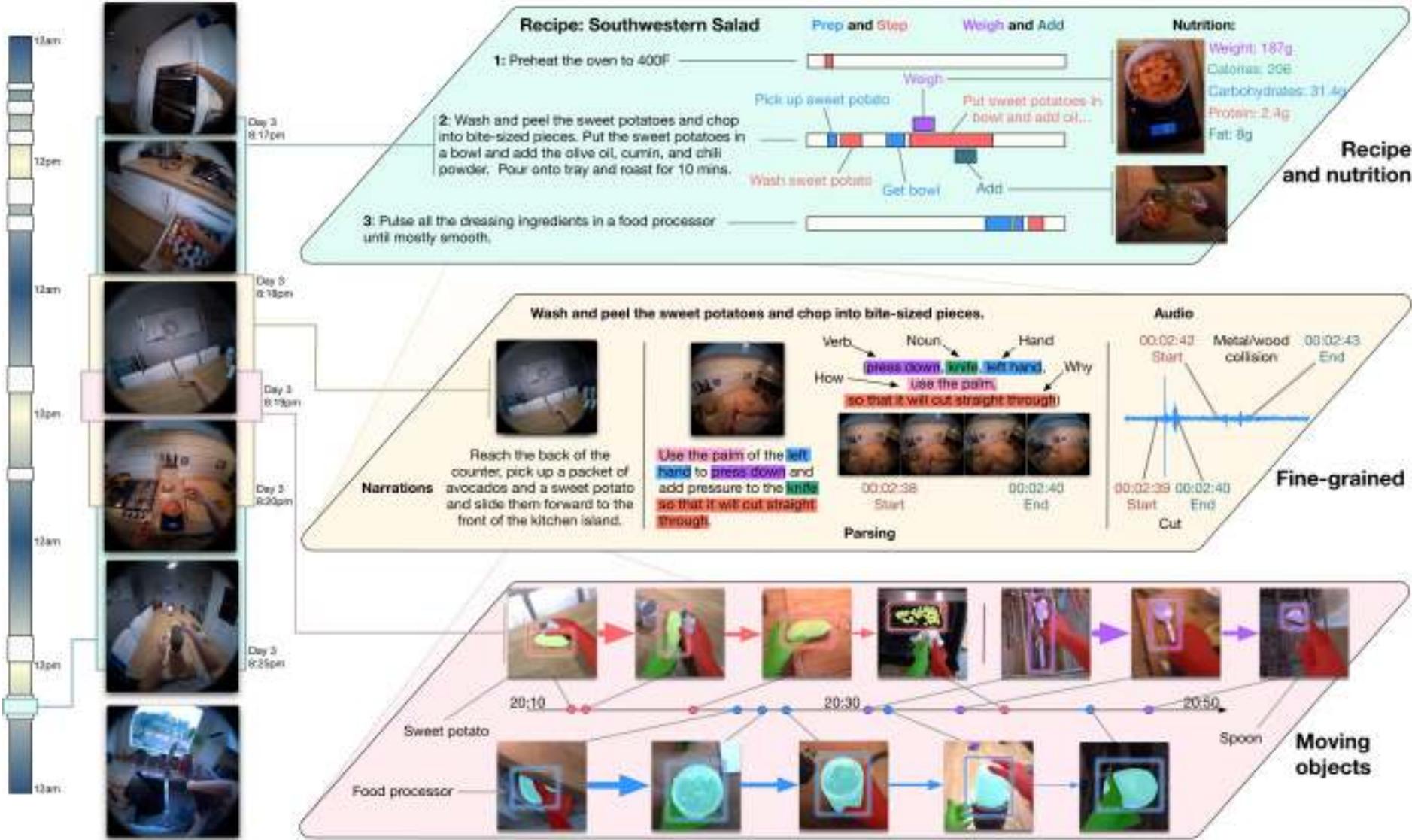
plastic measuring cup (1.0m) Add

hand washing liquid (1.3m) Add

kitchen towel (1.5m) Add



# HD-EPIC

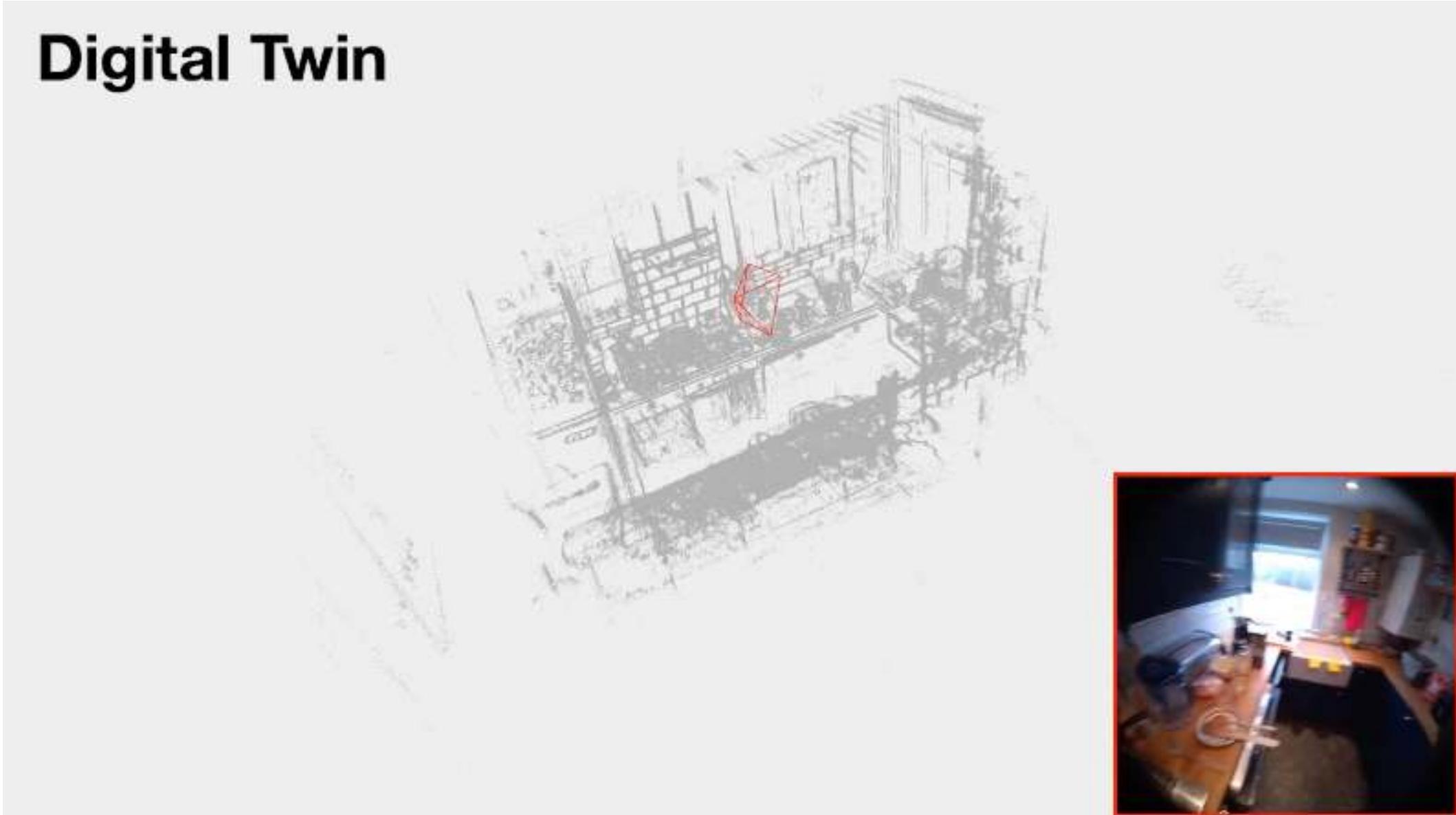




# HD-EPIC

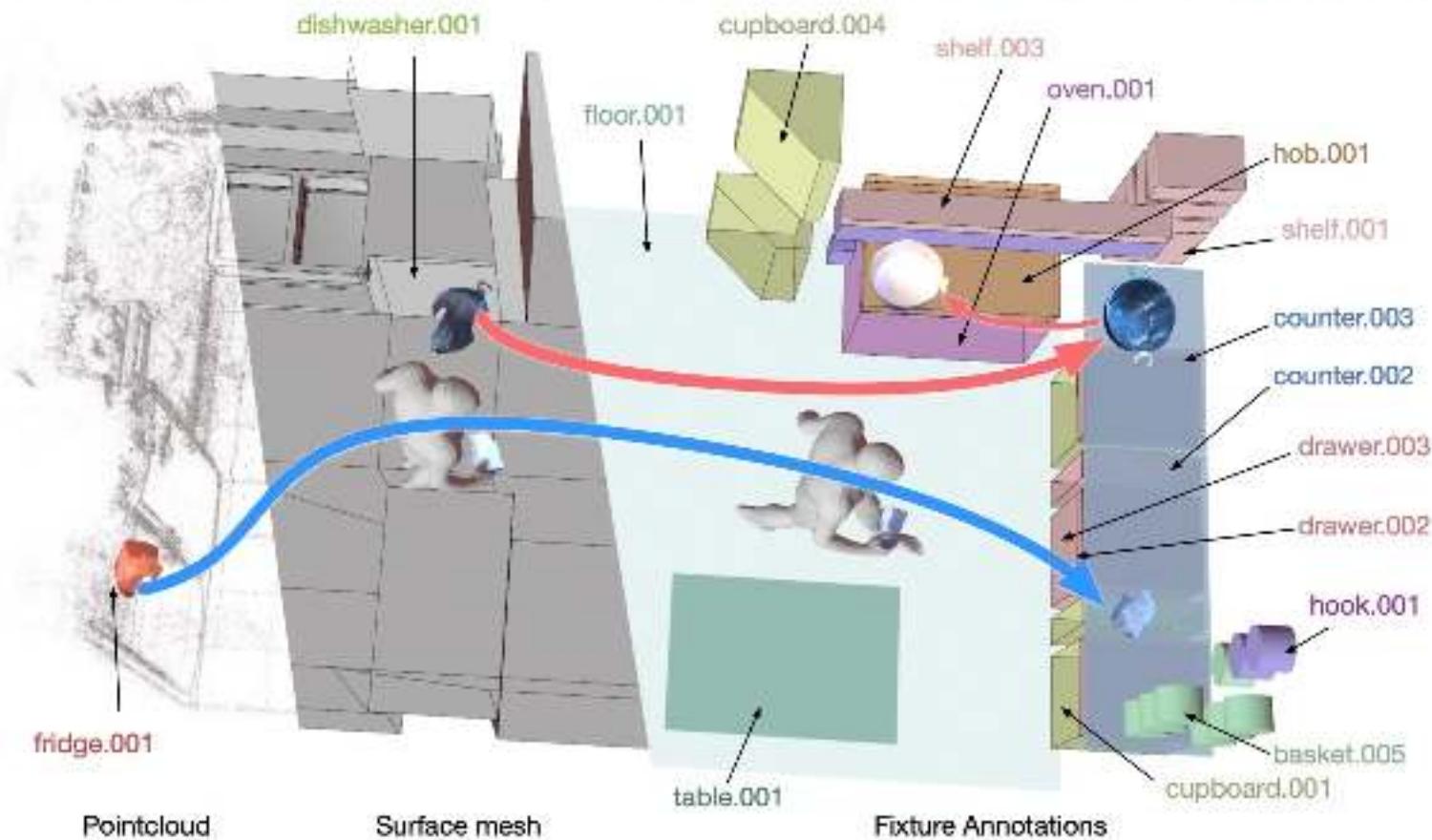


## Digital Twin



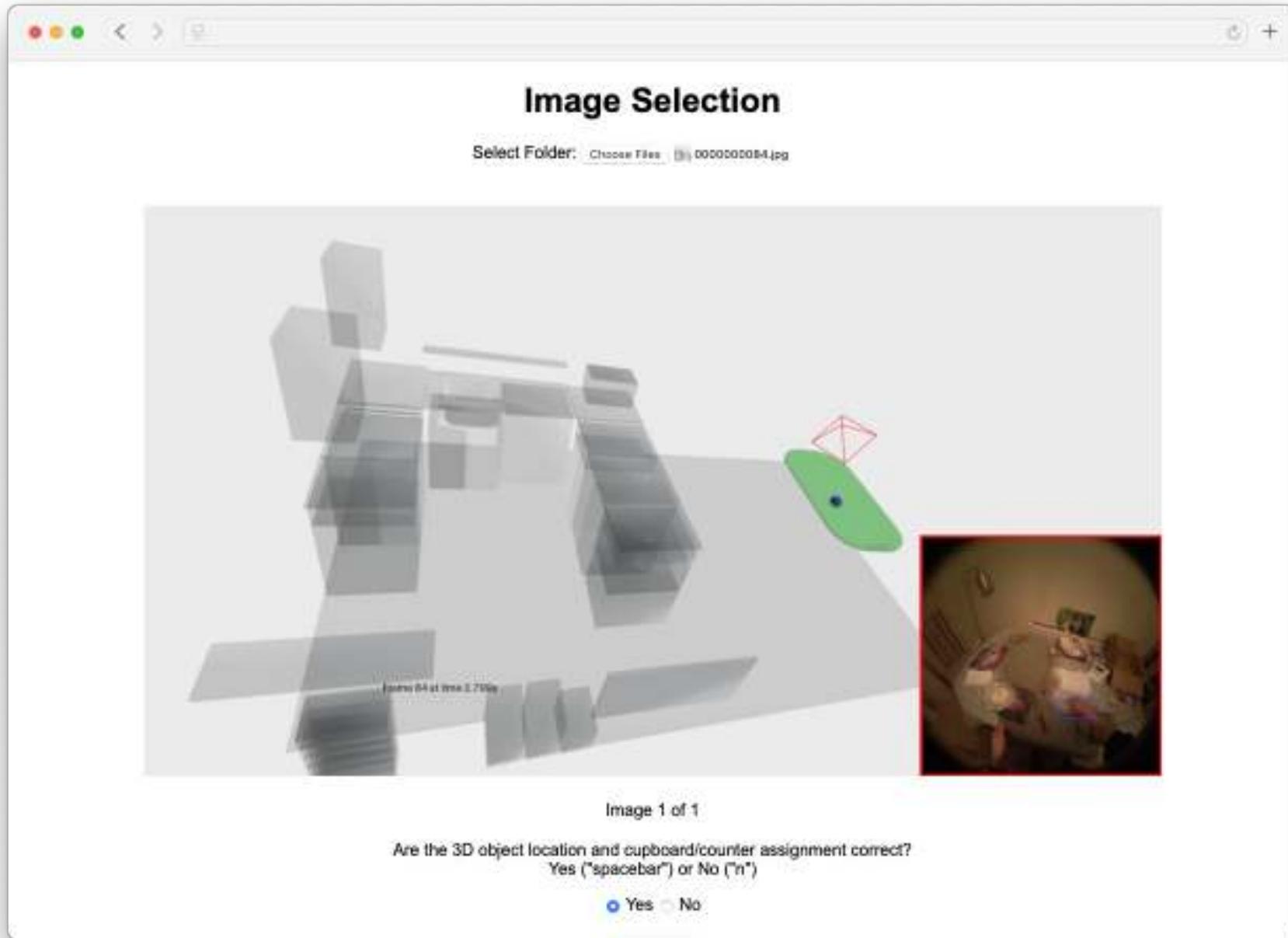


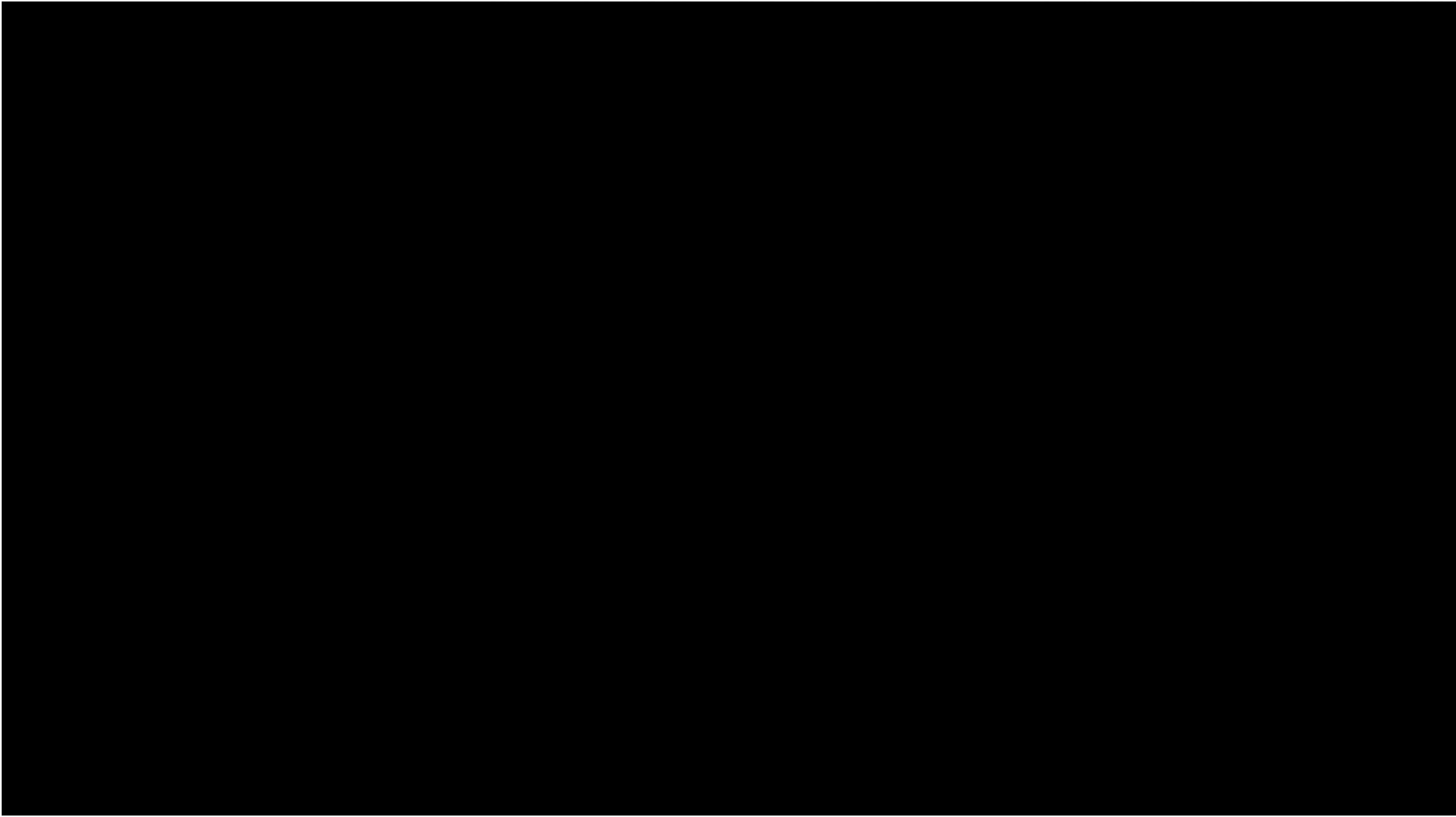
# HD-EPIC





# HD-EPIC





# In today's talk...



Motivation and Datasets in  
Egocentric Video Understanding



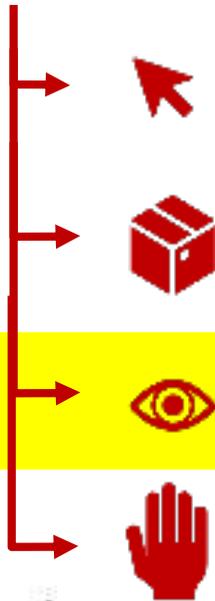
Teaser: The Wizard of Oz  
& Genie 3



Video Understanding  
Out of the Frame



Outlook into the Future of  
Egocentric Vision



Point Tracking



Object Tracking



Gaze Priming



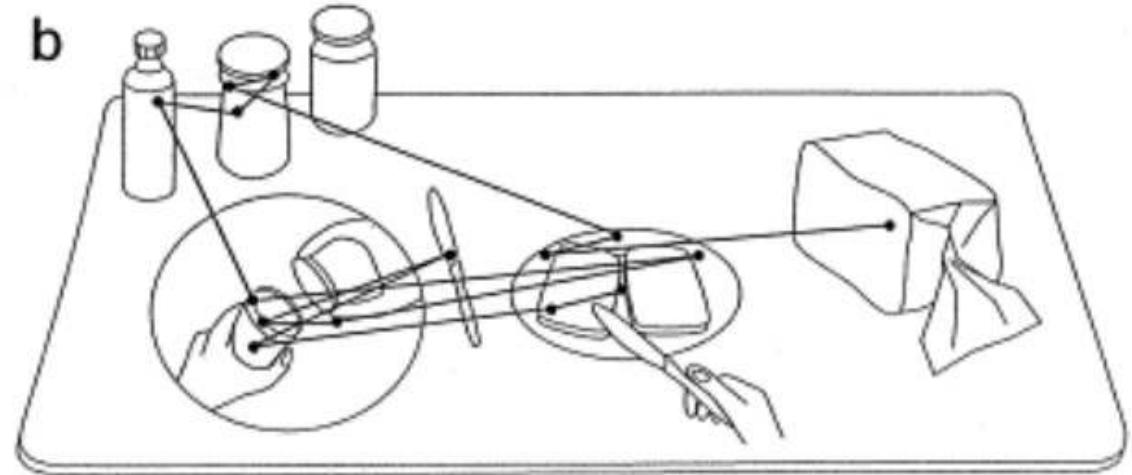
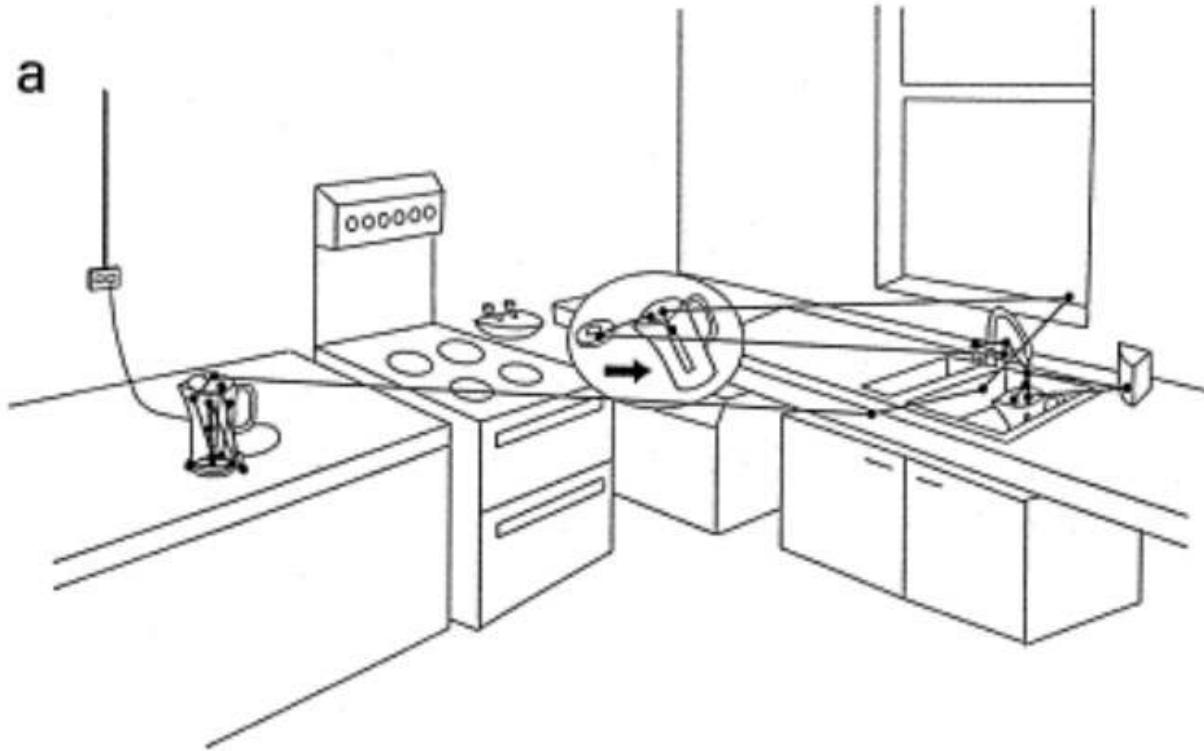
Hand Tracking



Conclusion



# Gaze and Fixations

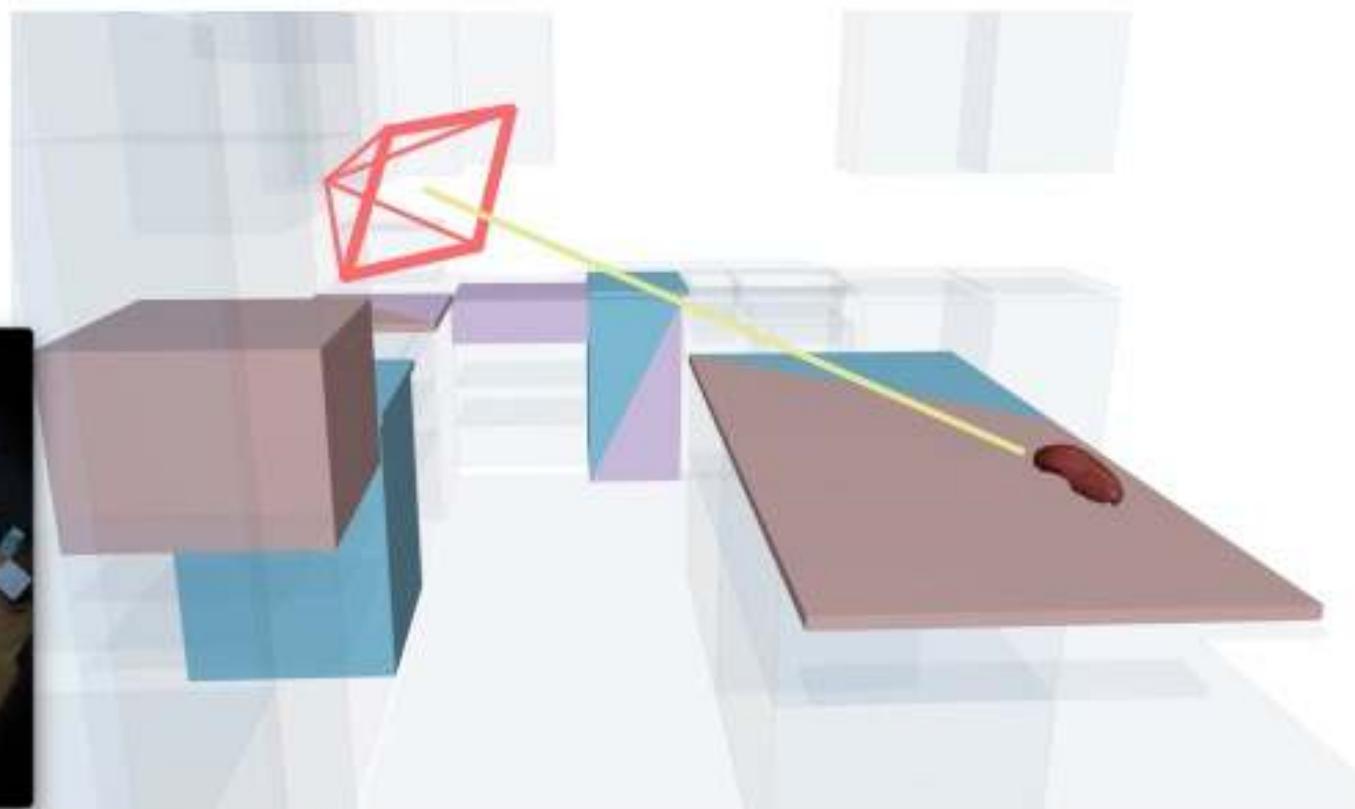
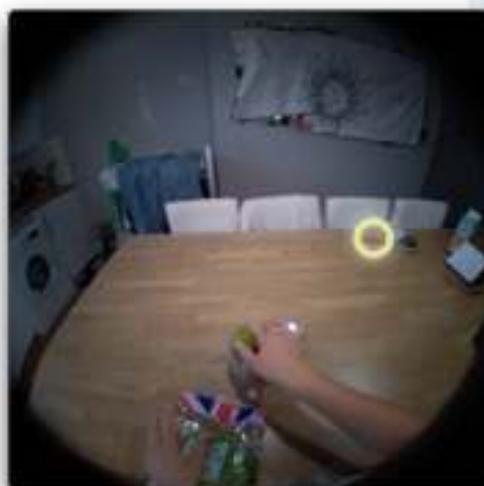




# HD-EPIC

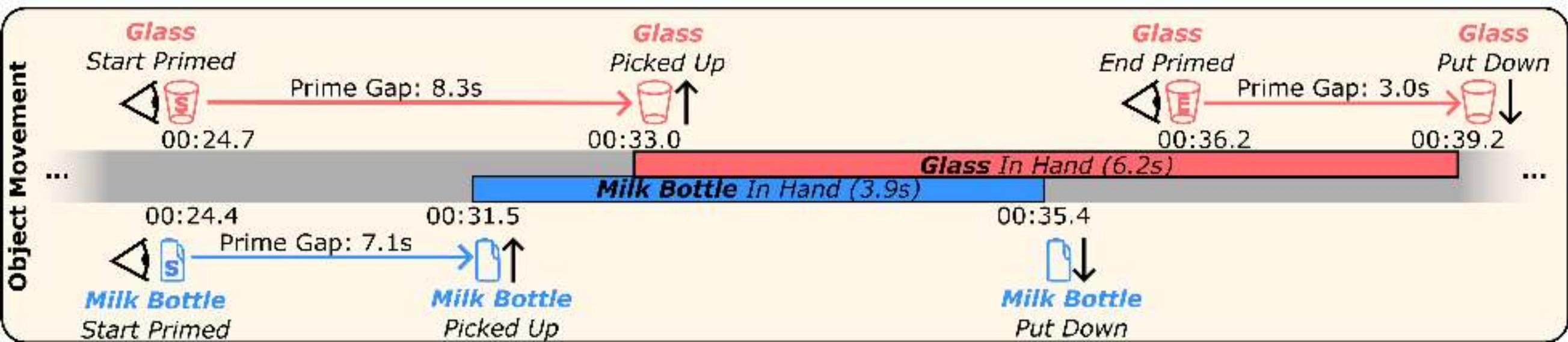


## Gaze priming





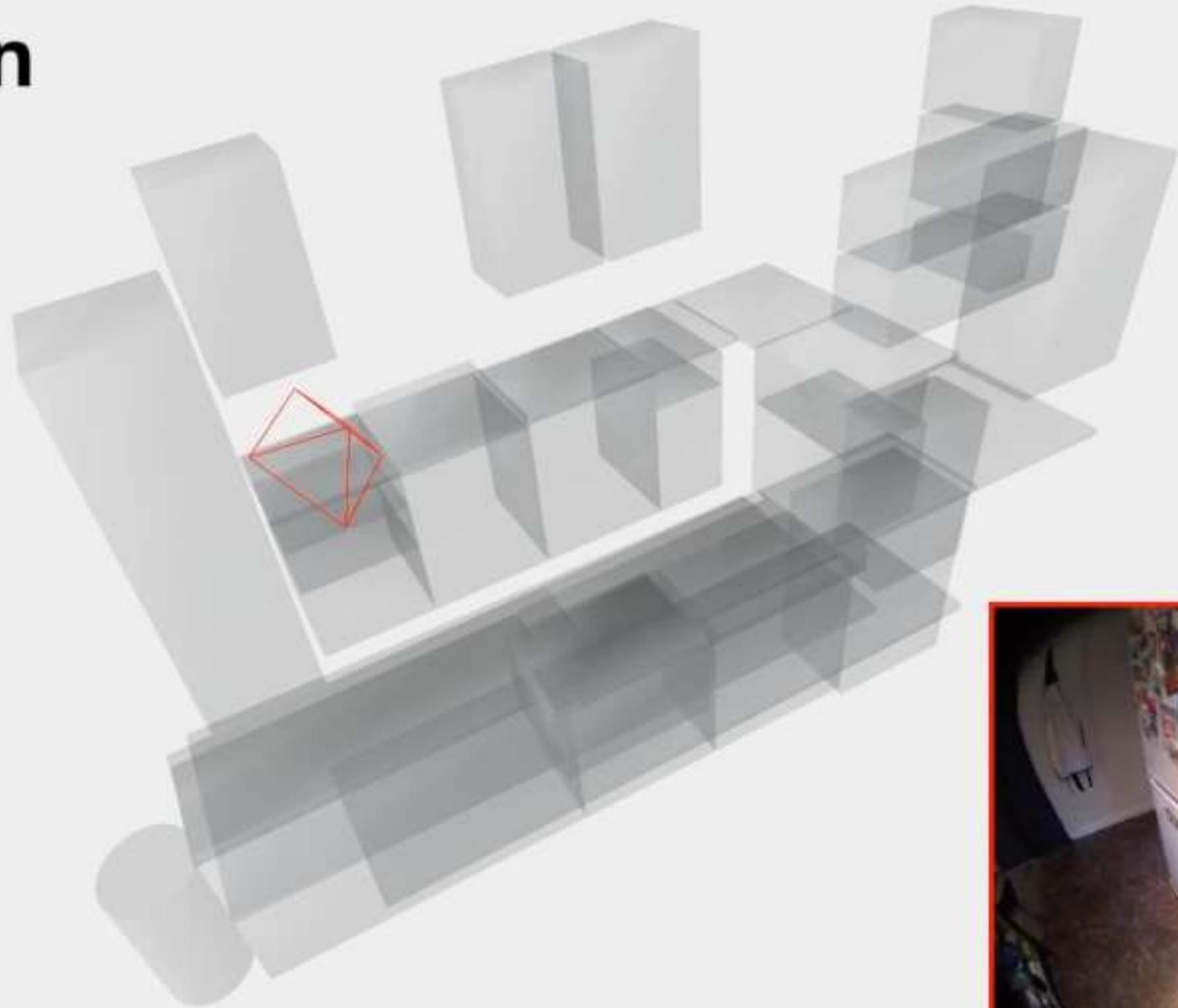
# HD-EPIC



# Digital Twin

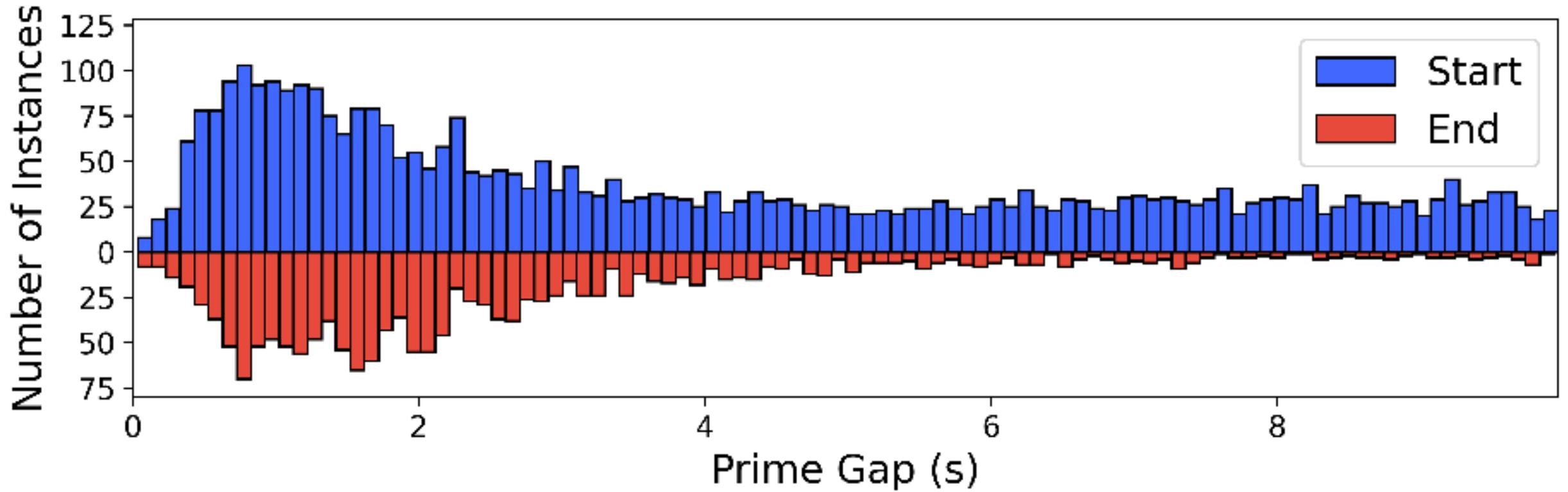
Fixtures

Open drawer



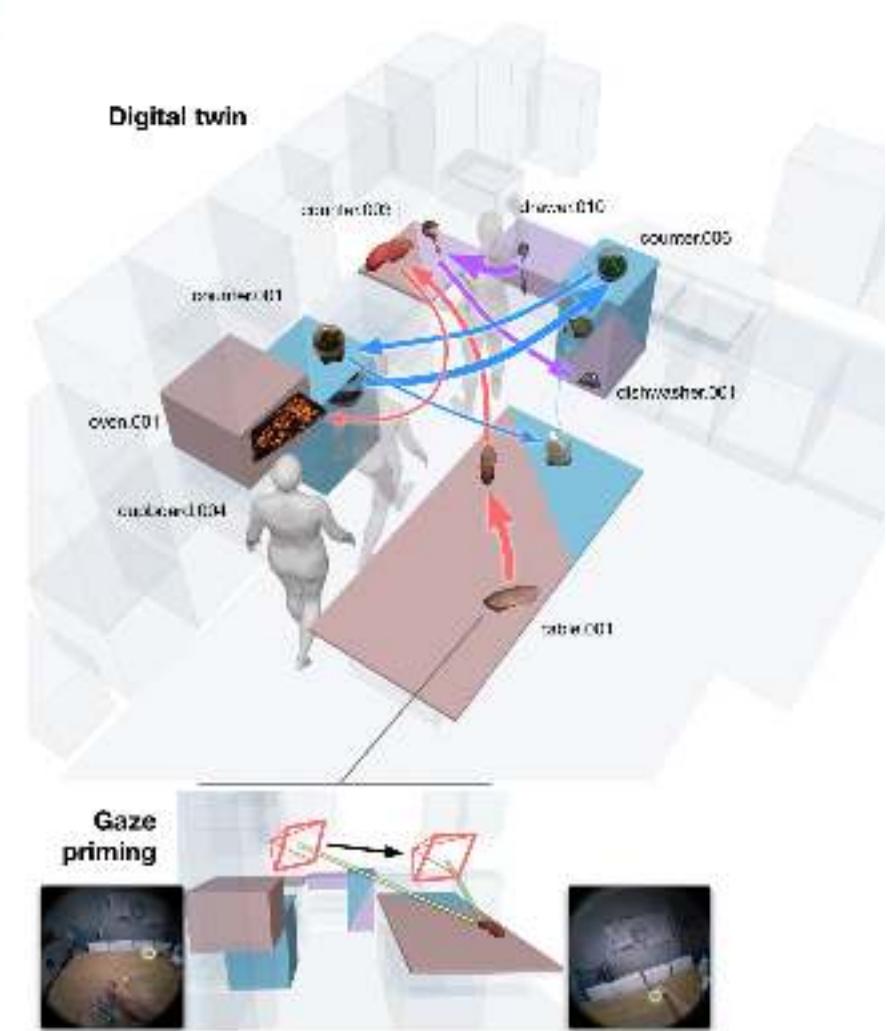
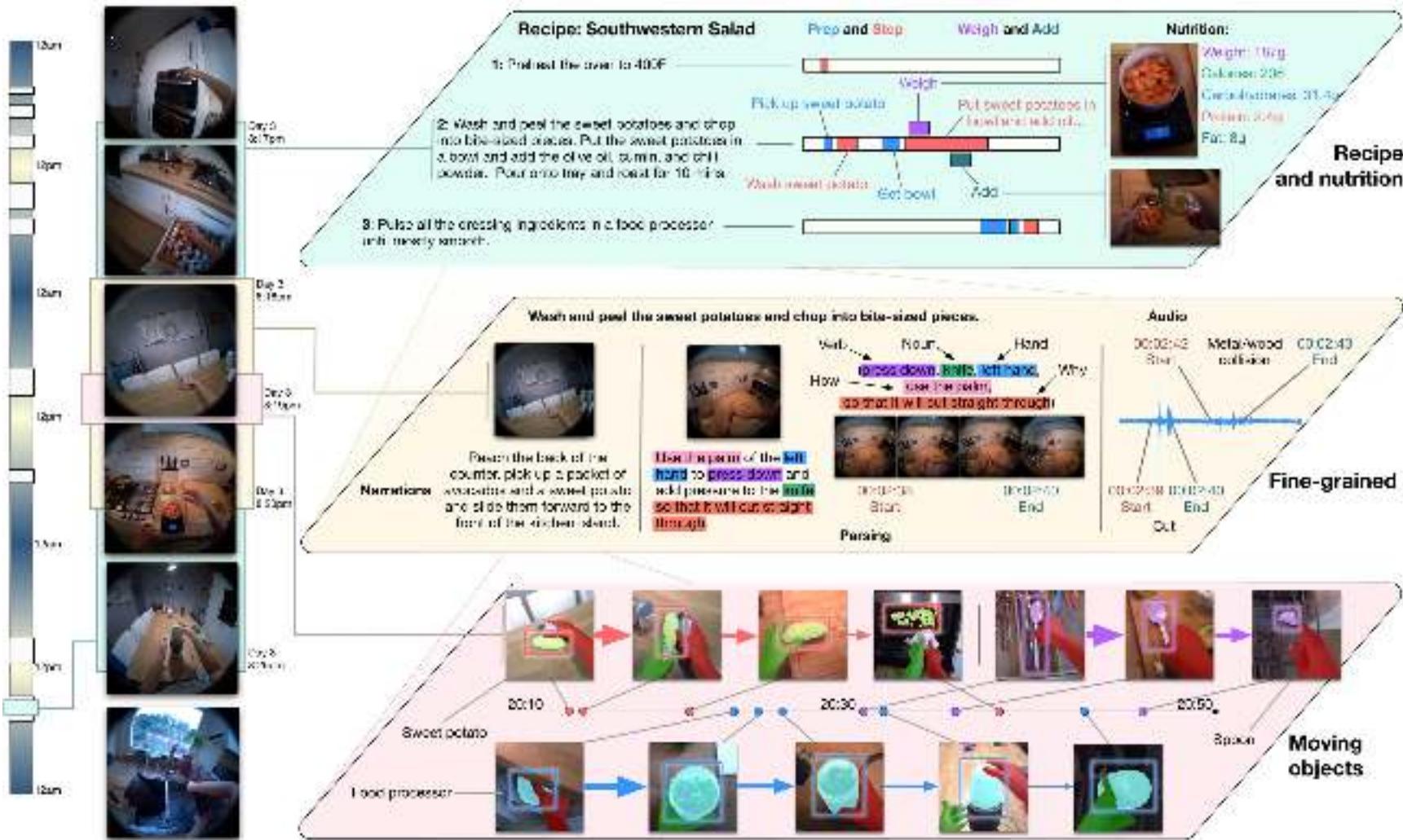


# HD-EPIC





# 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



Try it Yourself

## Use Wise to Search through HD-EPIC



SCAN ME

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

# In today's talk...



Motivation and Datasets in  
Egocentric Video Understanding



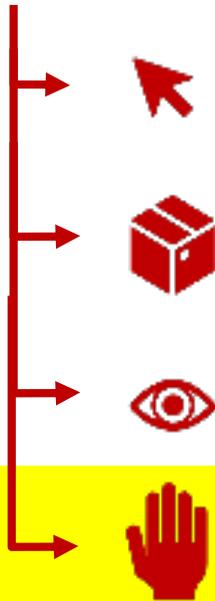
Teaser: The Wizard of Oz  
& Genie 3



Video Understanding  
Out of the Frame



Outlook into the Future of  
Egocentric Vision



Point Tracking



Object Tracking



Gaze Priming



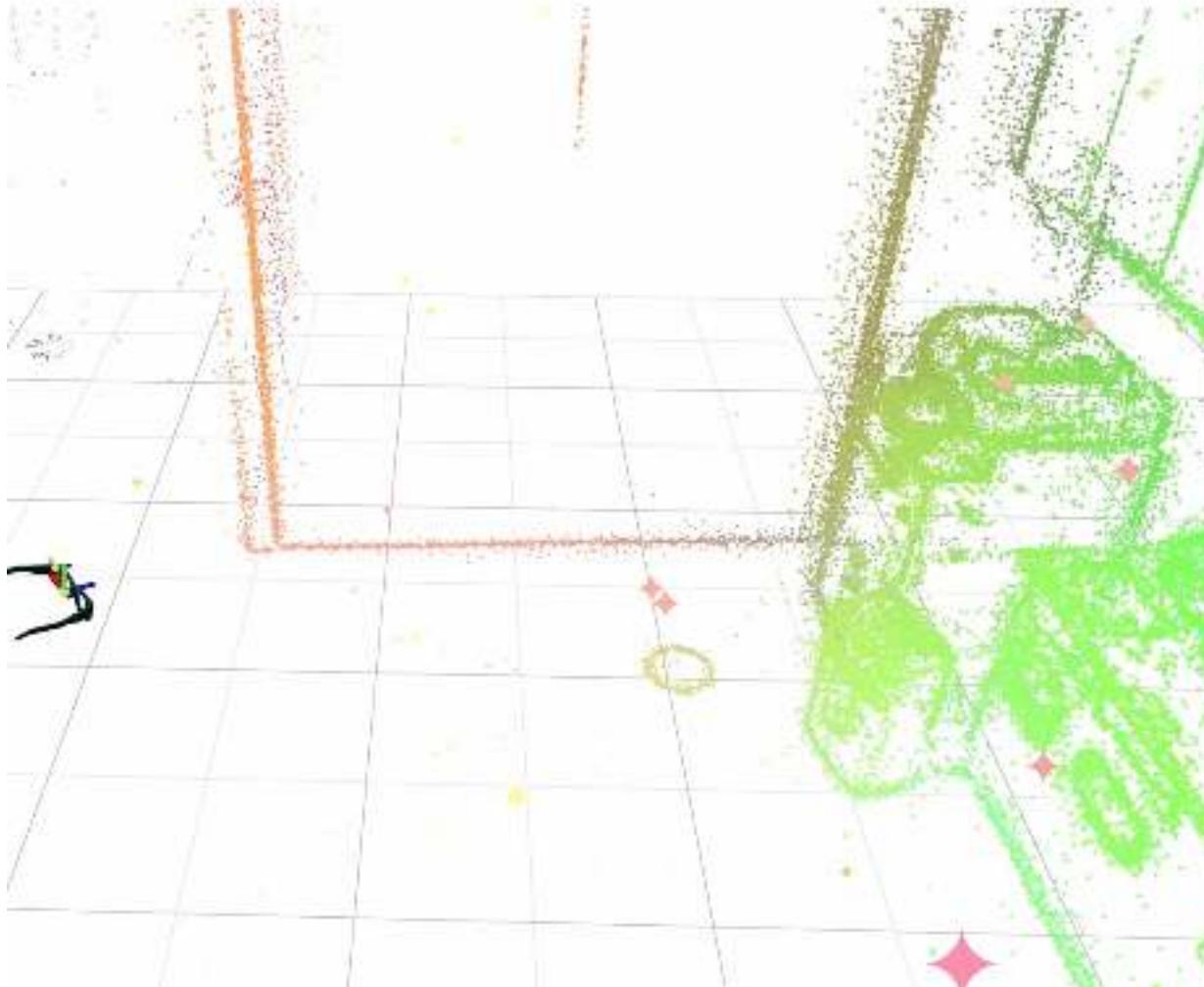
Hand Tracking



Conclusion

# EgoBody

EgoAllo uses egocentric (  $\epsilon$  ) SLAM poses and images



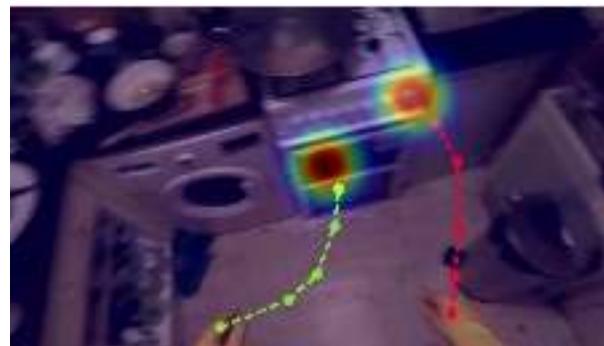
# EgoHand Forecasting – Previous Works

with: Masashi Hatano  
Zhifan Zhu  
Hideo Saito

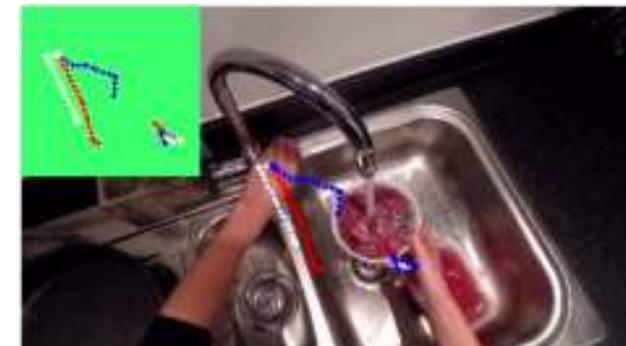
## 2D Hand Forecasting

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

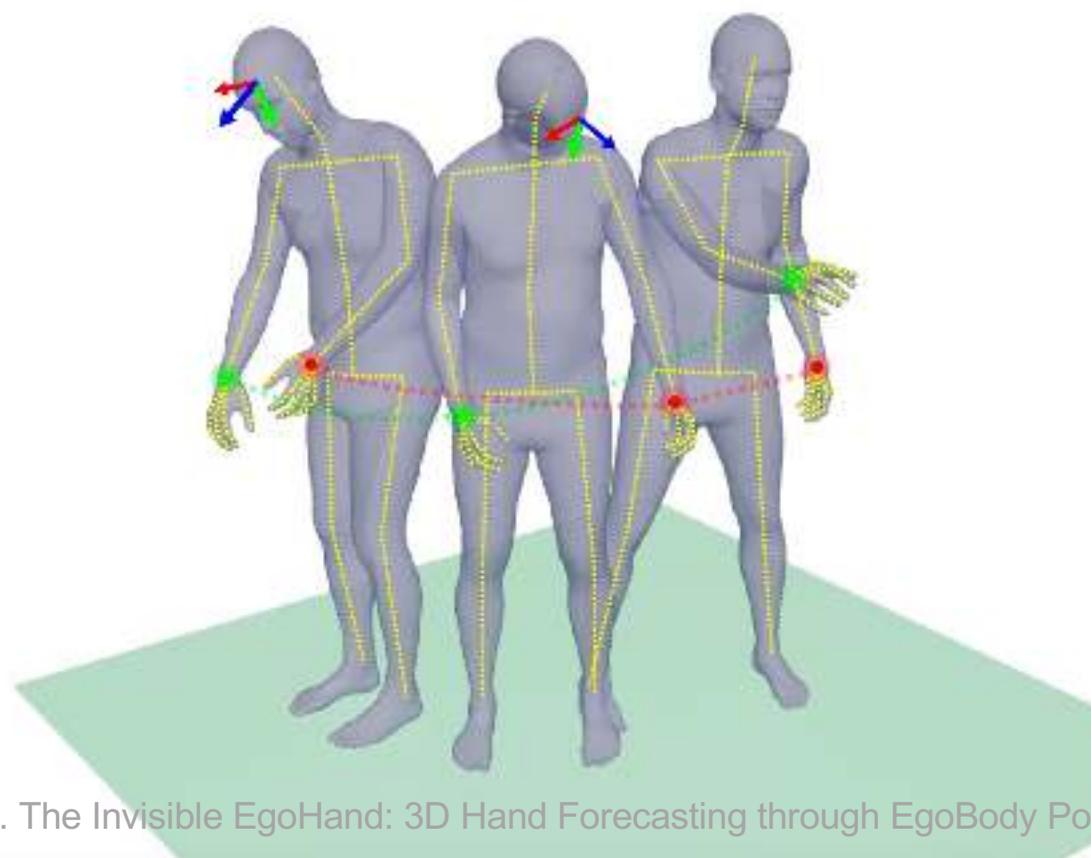
Given an egocentric video & 3D hand trajectory,  
forecast 3D hand positions of one hand



USST [ICCV'23]

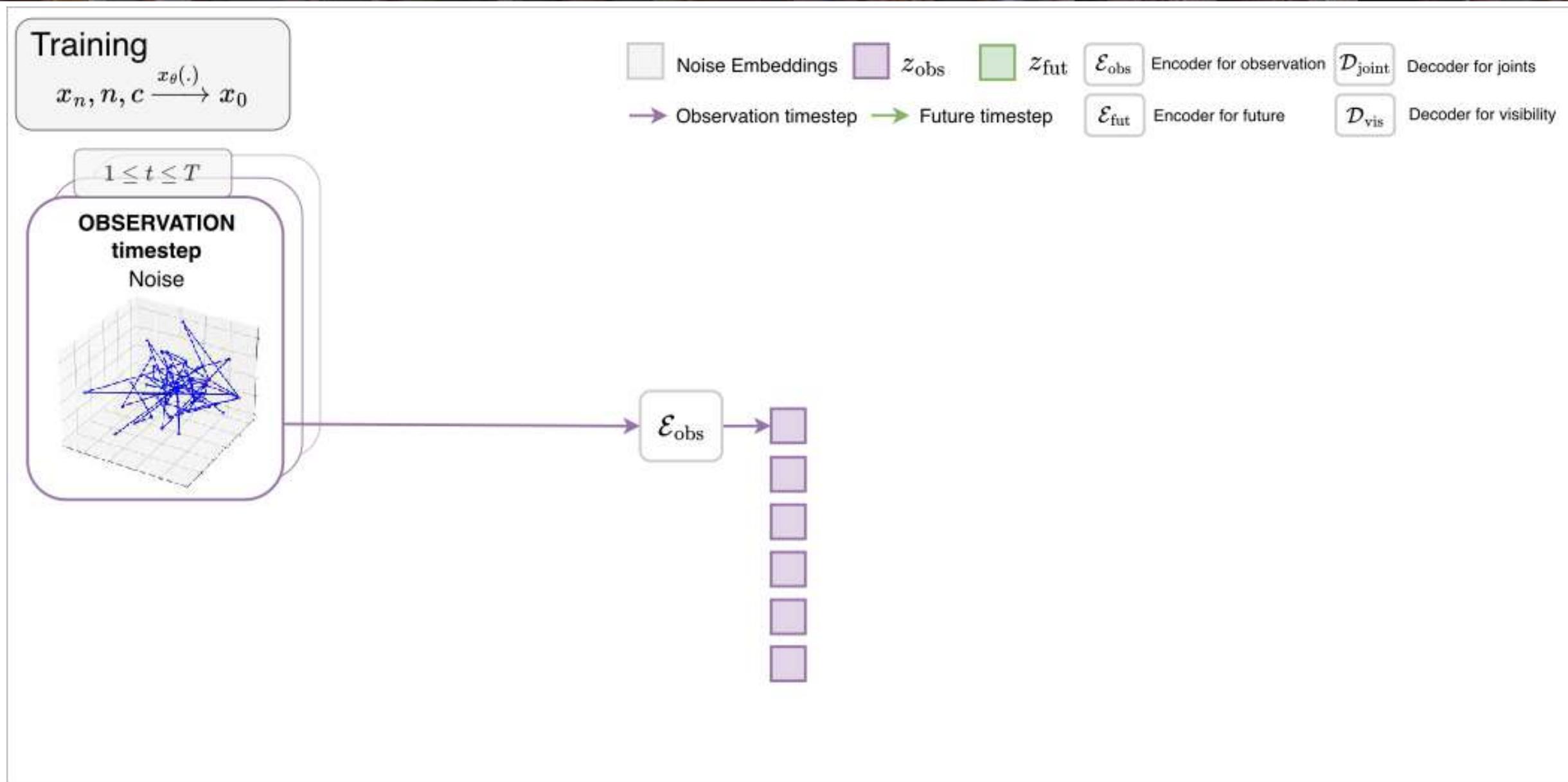
# The Invisible EgoHand

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	-	-
<b>Ours</b>	<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

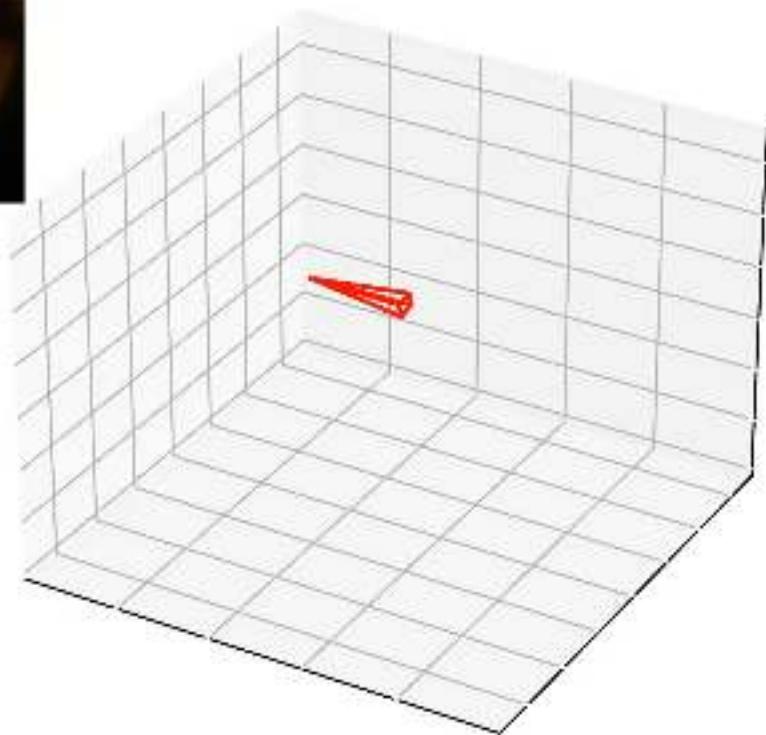
# The Invisible EgoHand

with: Masashi Hatano  
Zhifan Zhu  
Hideo Saito



Observation

 Camera Pose



Observation



In today's talk...



Motivation and Datasets in  
Egocentric Video Understanding



Teaser: The Wizard of Oz  
& Genie 3



Video Understanding  
Out of the Frame



Point Tracking



Object Tracking



Gaze Priming



Hand Tracking



Outlook into the Future of  
Egocentric Vision



Conclusion

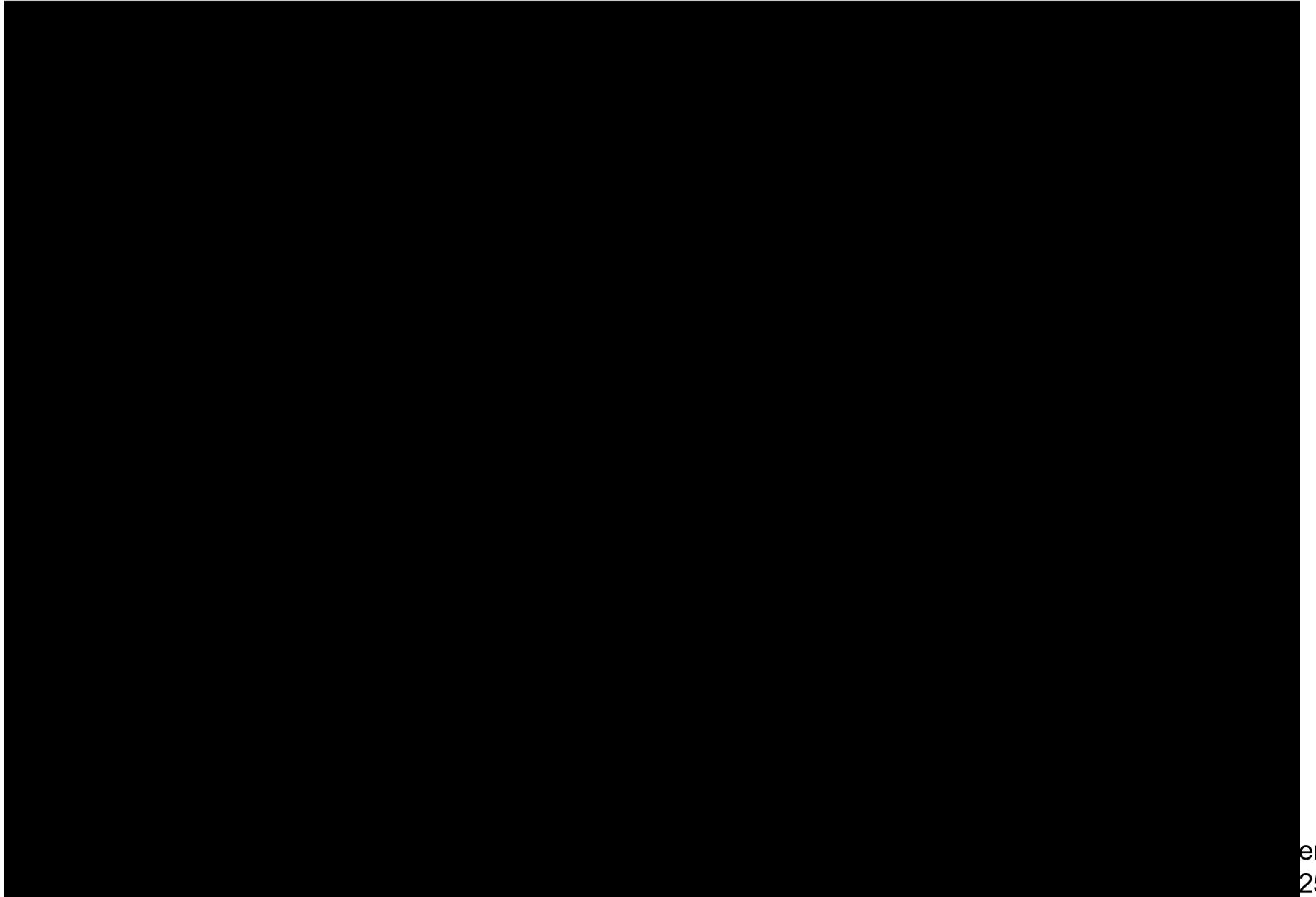
# The Wizard of Oz at the Sphere

Sold out tickets – August 2025



# The Wizard of Oz @ The Sphere

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



# The Wizard of Oz @ The Sphere



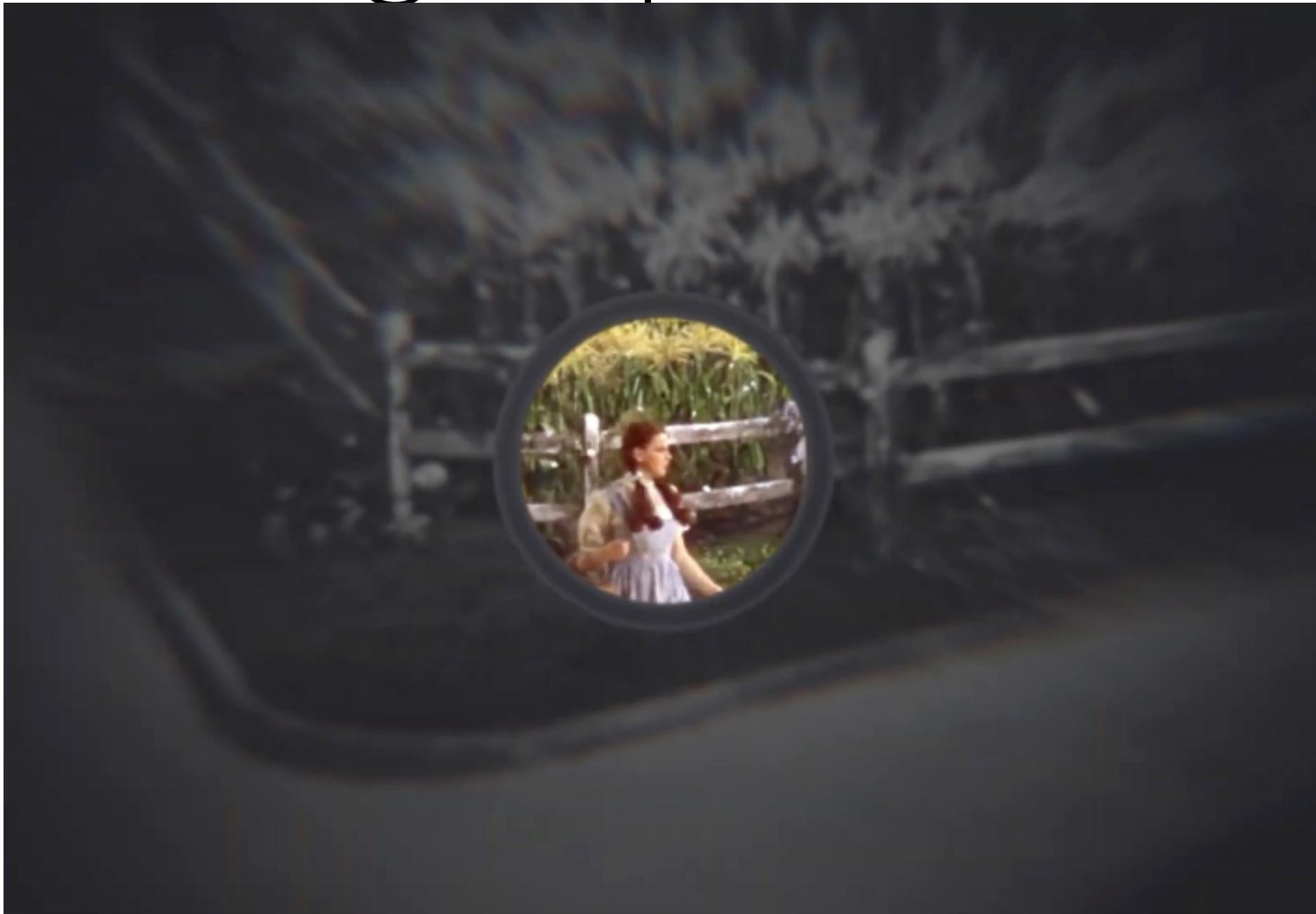
**Ralph Winte**

Head of Physical Production -  sphere



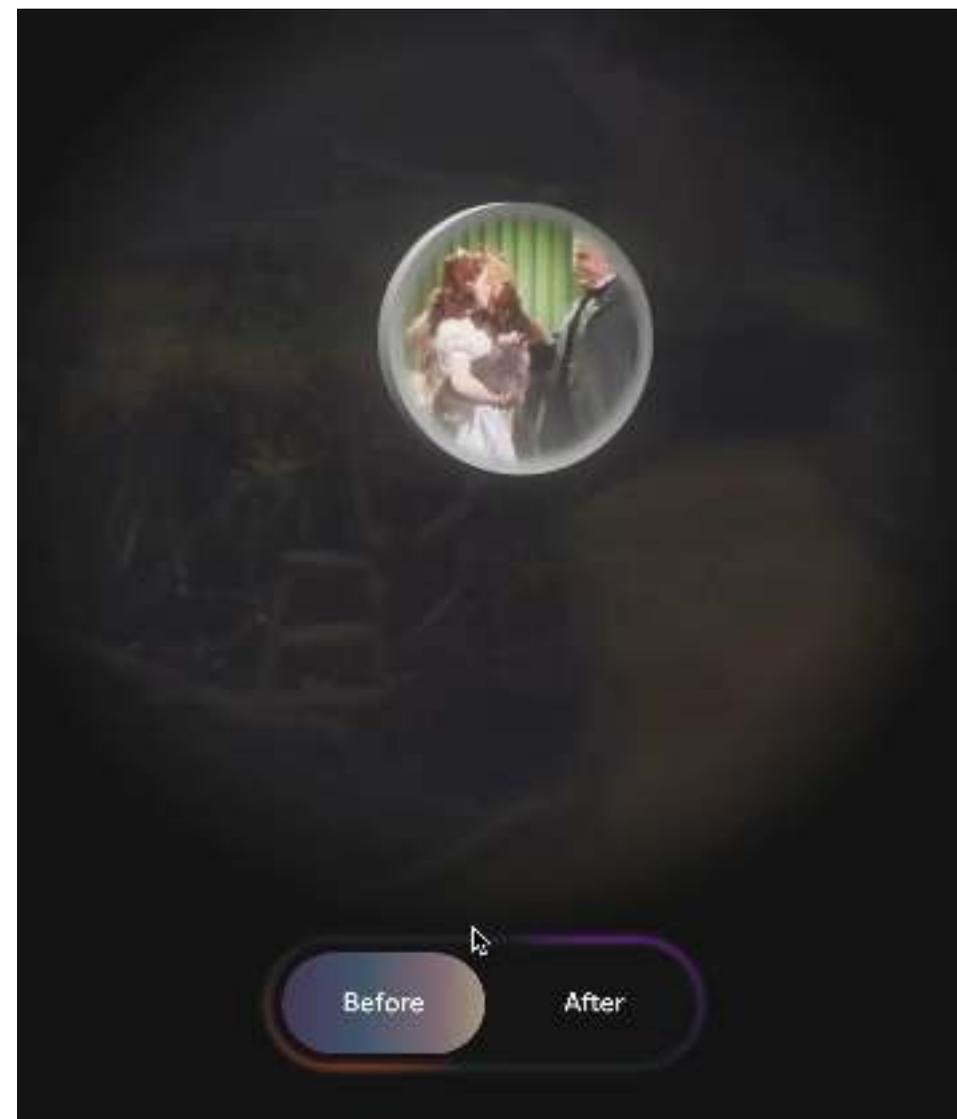
Dima Damen  
Workshop @ICCV2025

# The Wizard of Oz @ The Sphere



# The Wizard of Oz @ The Sphere

- Super-resolution,
- Outpainting...

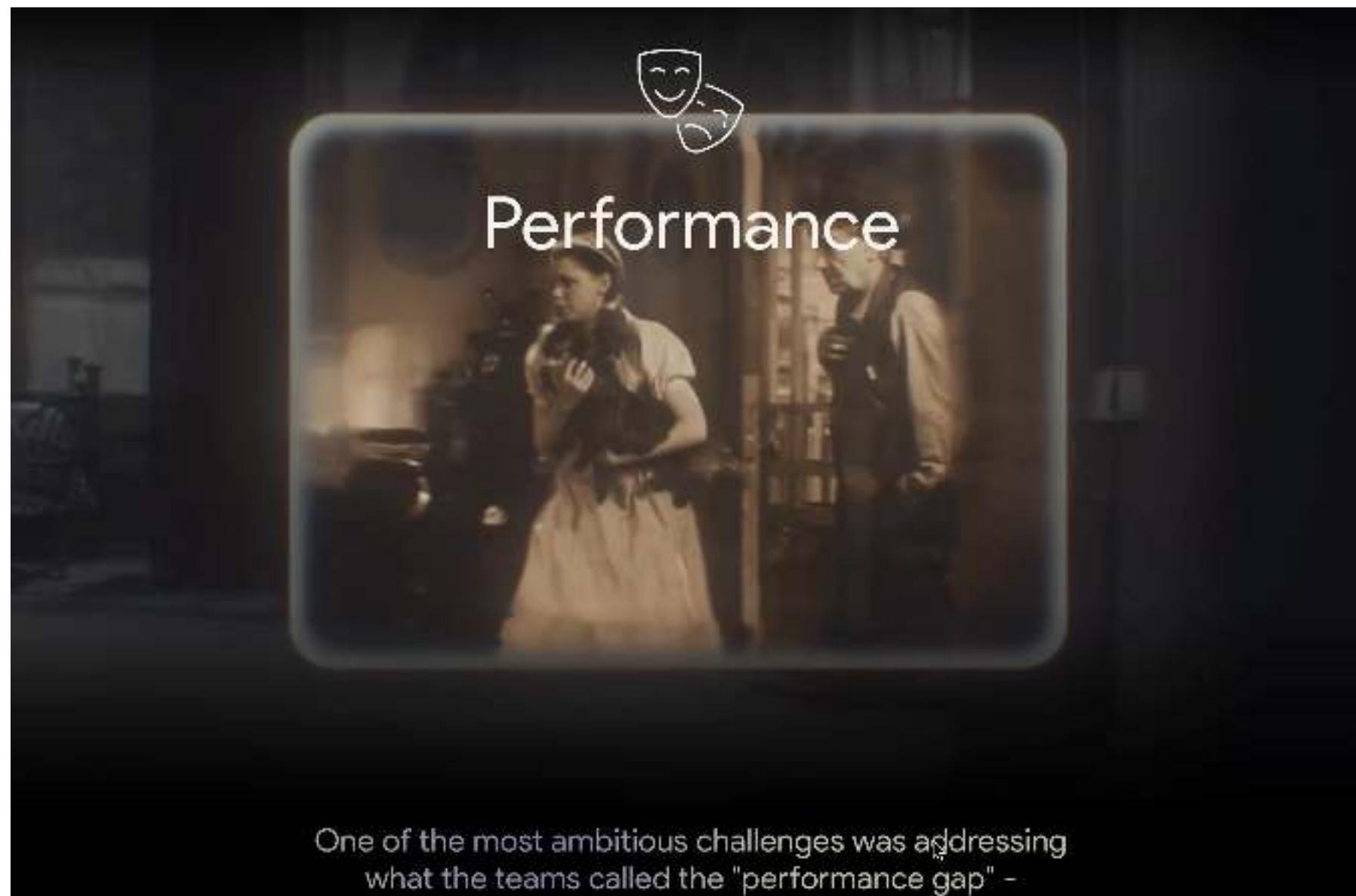


<https://behindthecurtain.withgoogle.com>

Dima Damen  
BinEgo-360 Workshop @ICCV2025

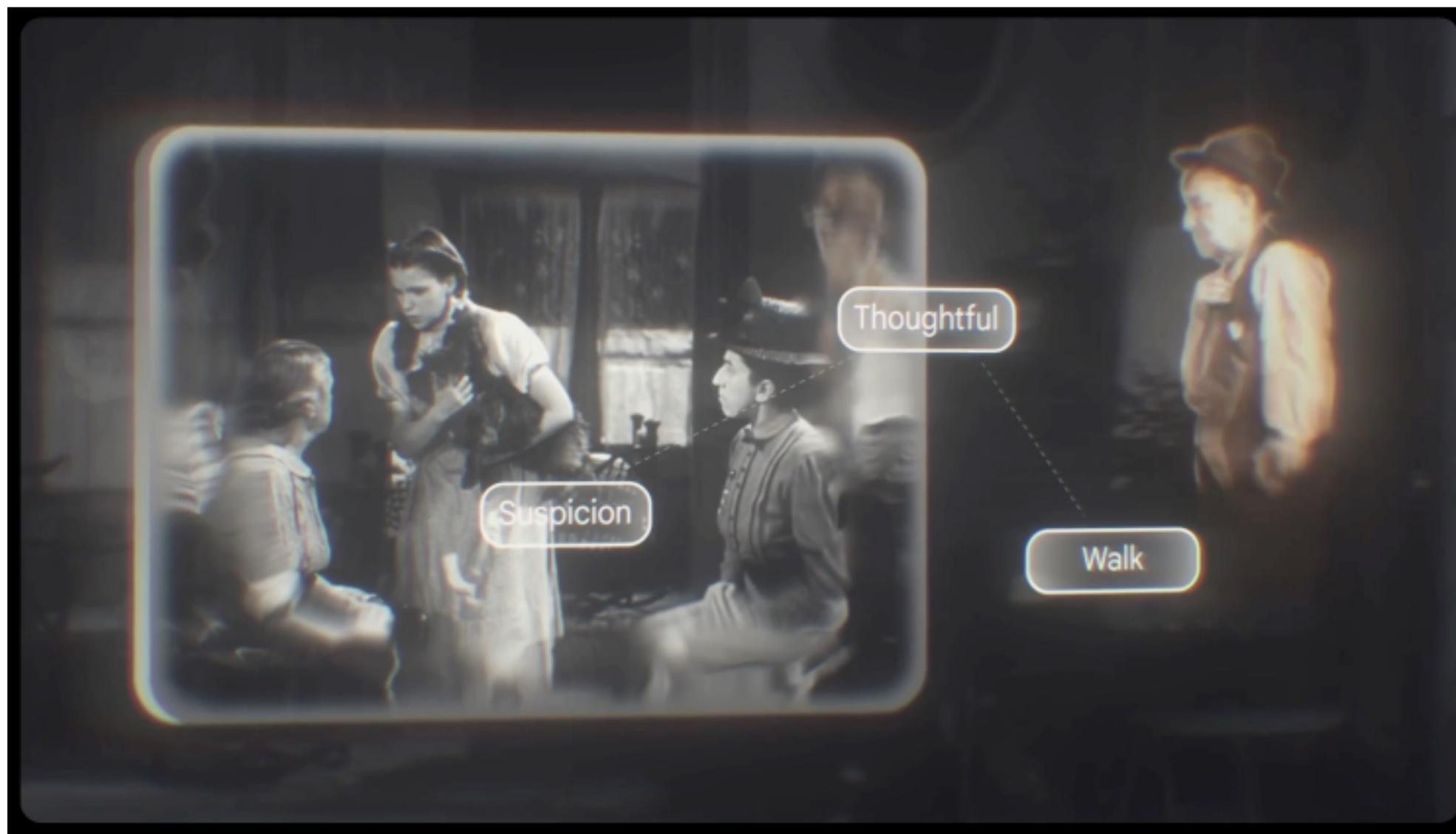
# The Wizard of Oz @ The Sphere

- Performance Interpolation,



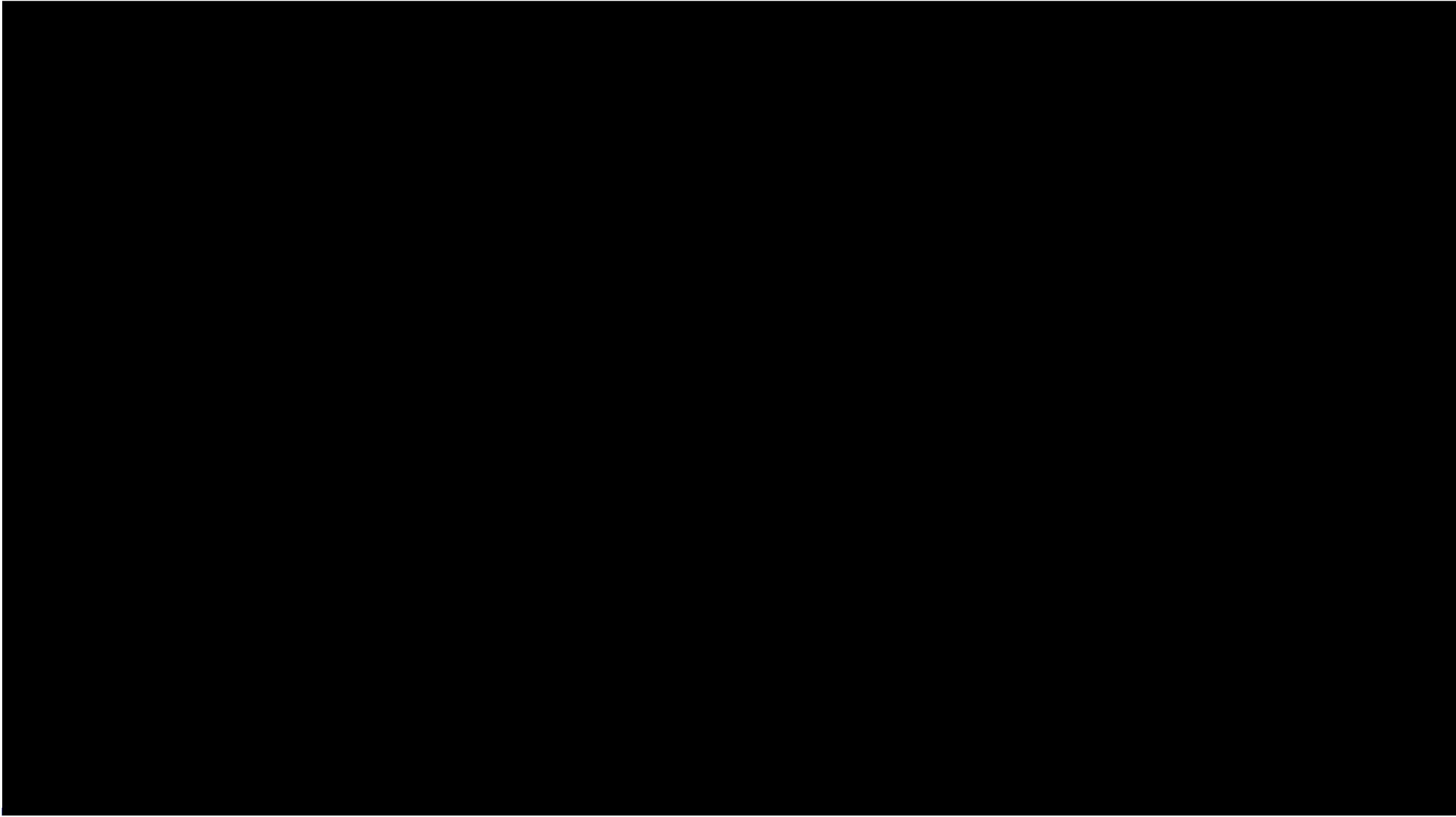
# The Wizard of Oz @ The Sphere

- Auto-Director



<https://behindthecurtain.withgoogle.com>

Dima Damen  
BinEgo-360 Workshop @ICCV2025





# 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  
BinEgo-360 Workshop @ICCV2025

# Genie 3

Released Aug 2025



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

Dima Damen  
BinEgo-360 Workshop @ICCV2025



# 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**



In today's talk...



Motivation and Datasets in  
Egocentric Video Understanding



Teaser: The Wizard of Oz  
& Genie 3



Video Understanding  
Out of the Frame



Point Tracking



Object Tracking



Gaze Priming



Hand Tracking



Outlook into the Future of  
Egocentric Vision



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>



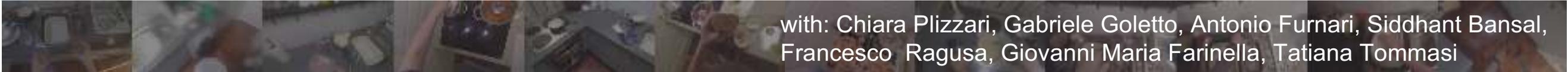
Politecnico  
di Torino



University of  
BRISTOL



UNIVERSITÀ  
degli STUDI  
di CATANIA



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**



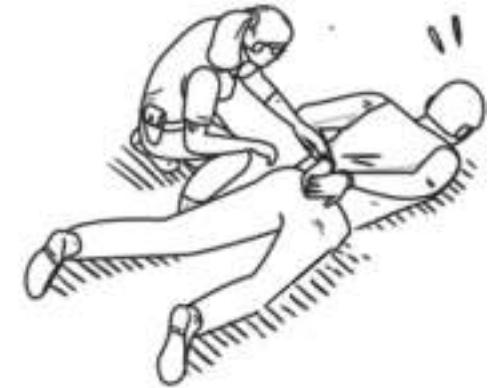
**EGO-Tourist**



**EGO-Worker**



**EGO-Home**

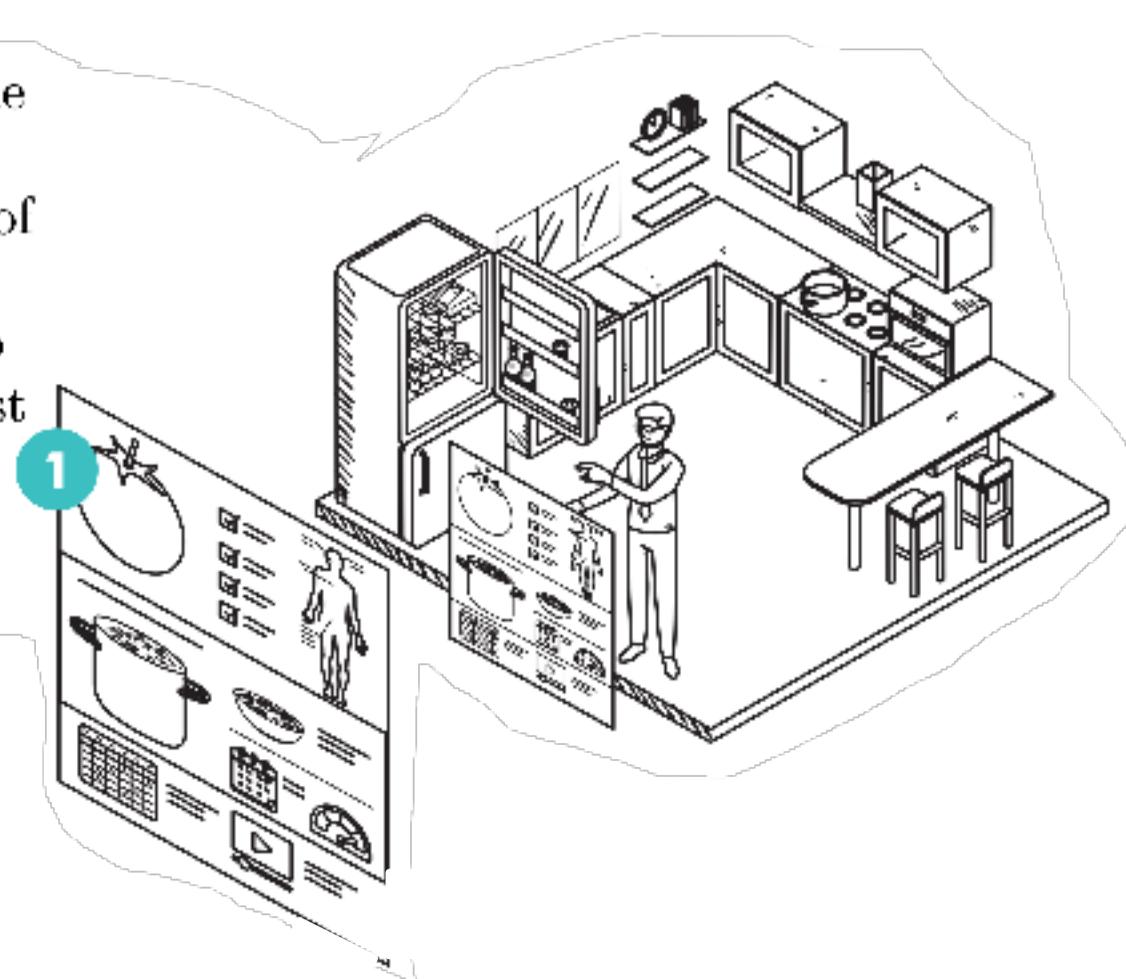


**Ego-Police**

# EGO-Home

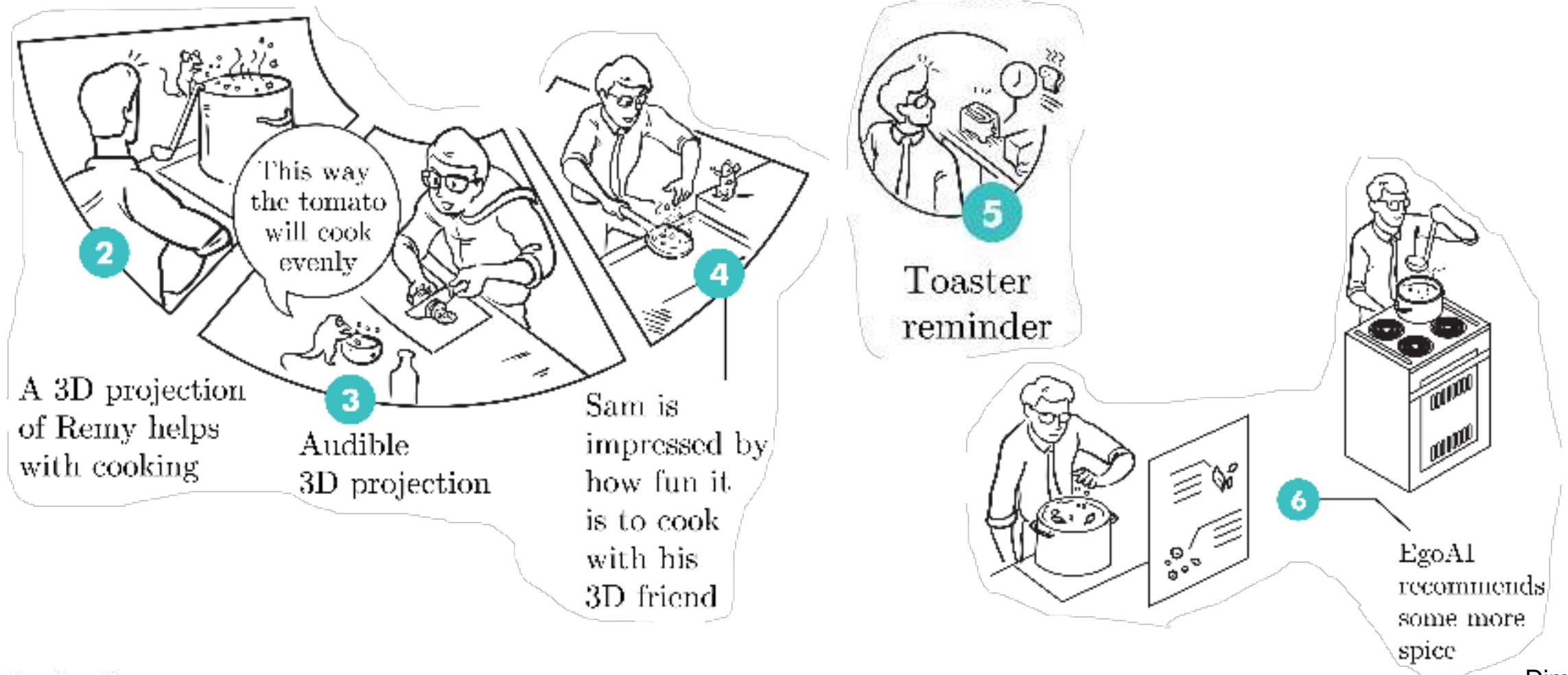
with: Chiara Plizzari, Gabriele Goletto, Antonio Furnari, Siddhant Bansal, Francesco Ragusa, Giovanni Maria Farinella, Tatiana Tommasi

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



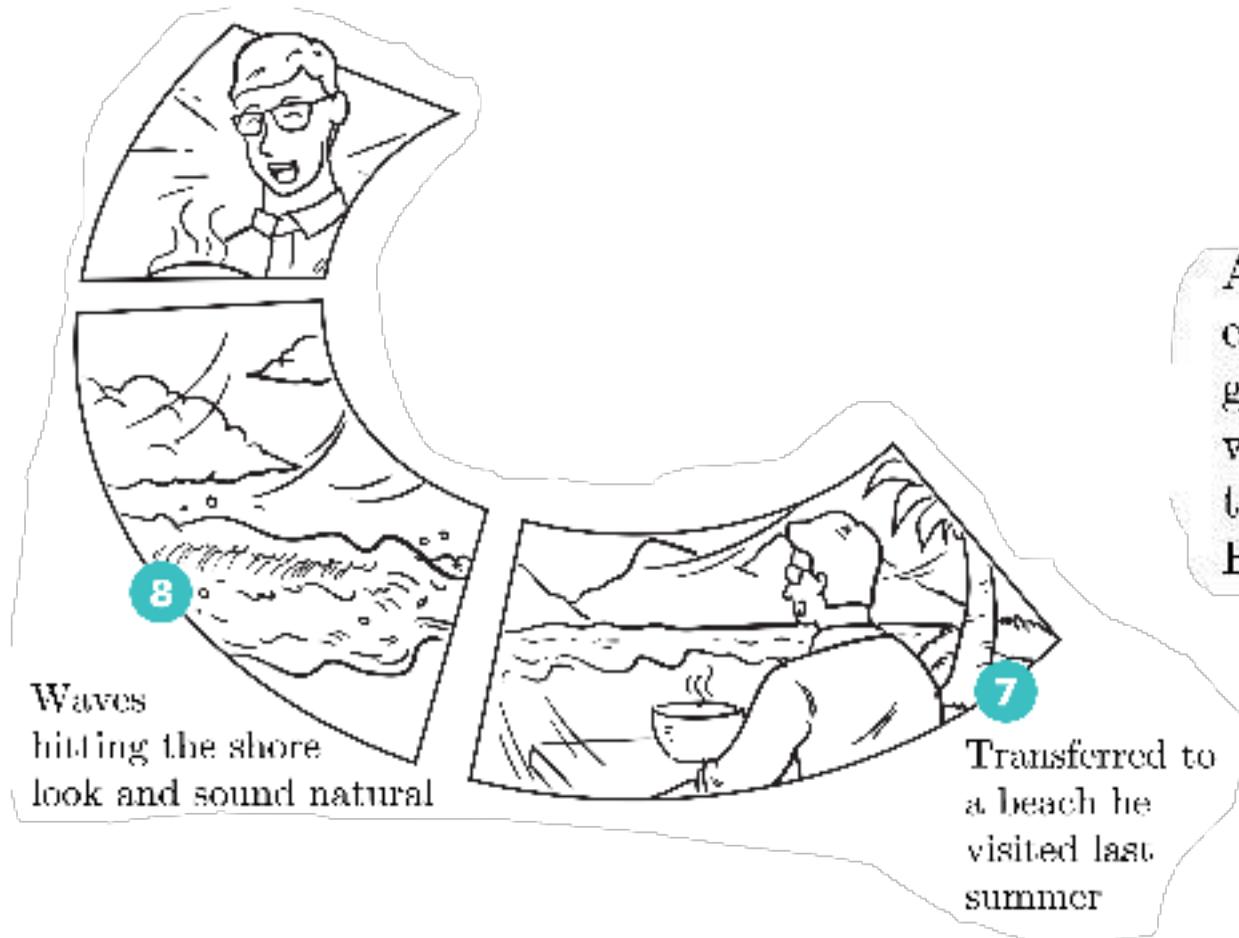
# EGO-Home

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



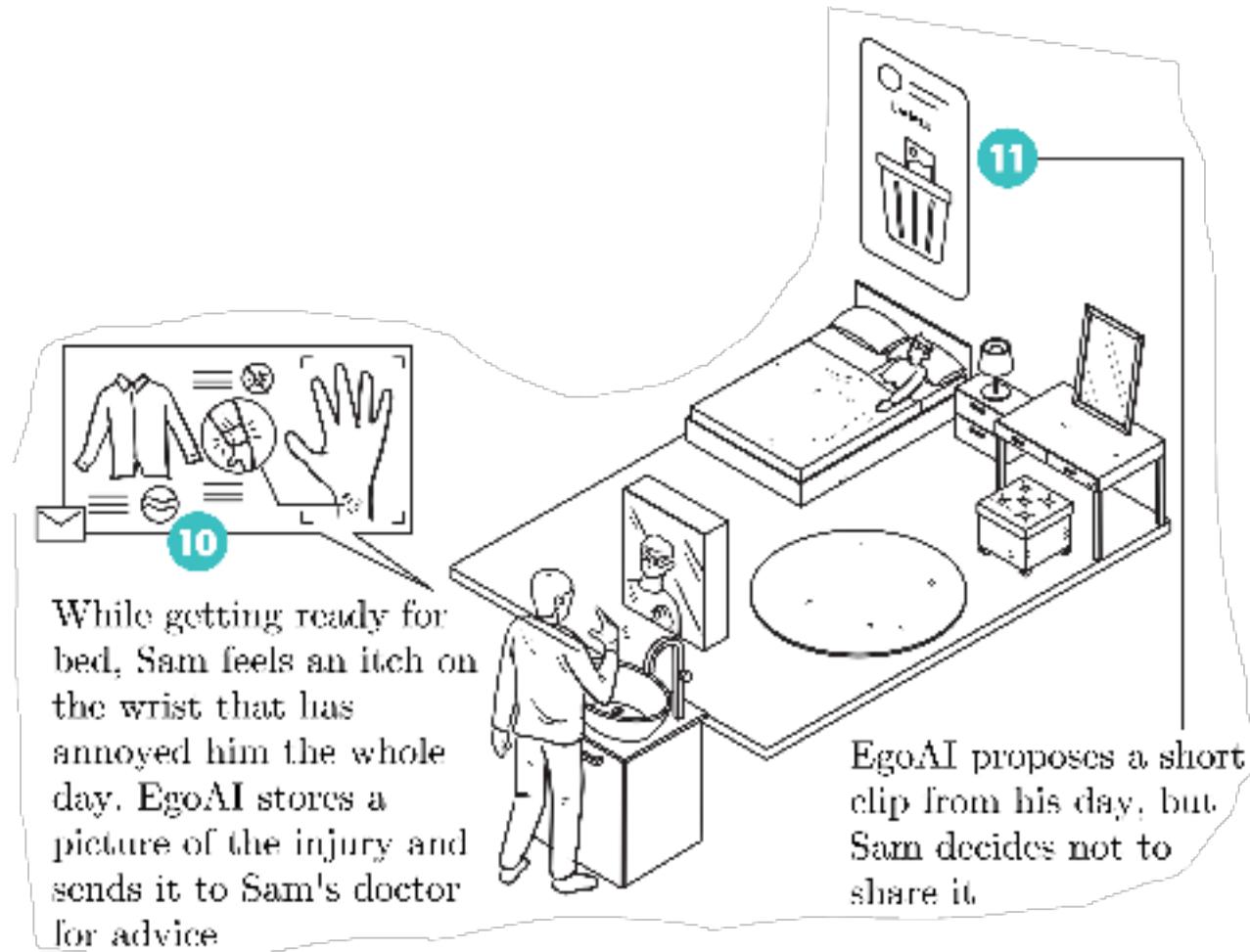
# EGO-Home

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

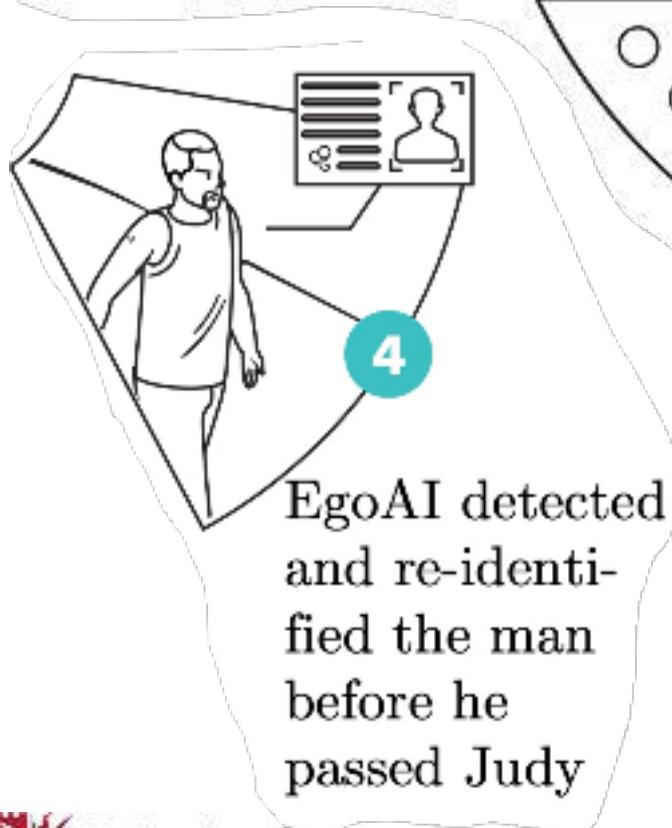




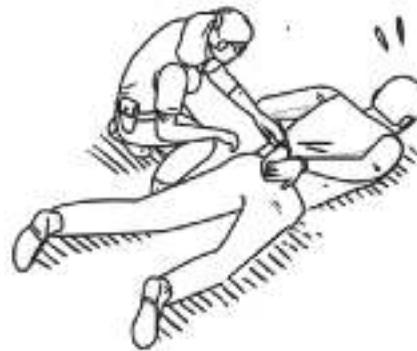
# From Stories to Tasks

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

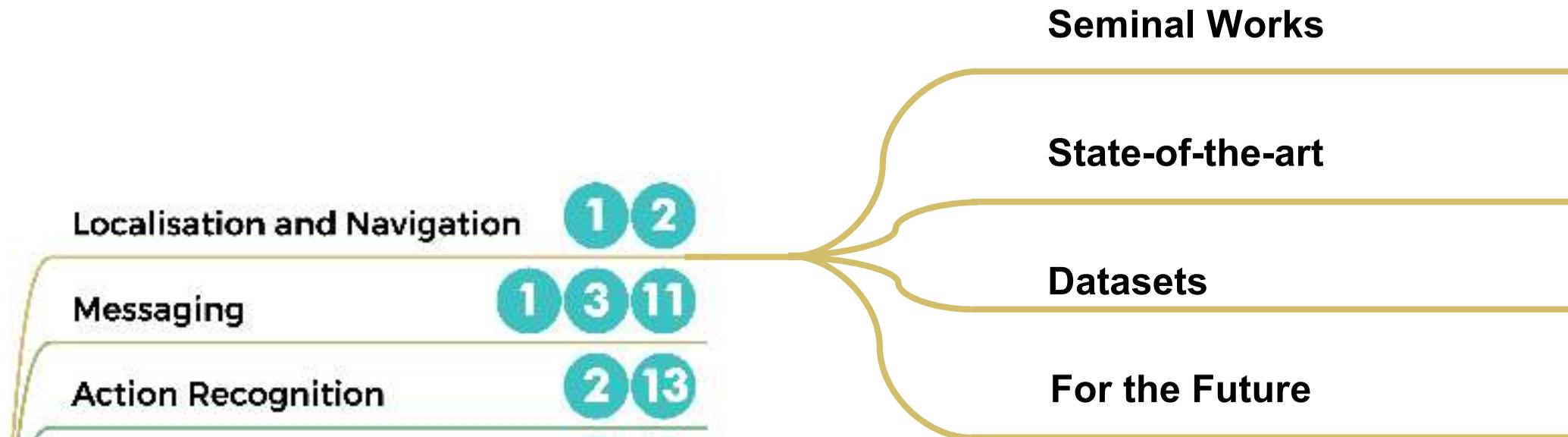
13

Privacy

14

# The Survey Part

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



In today's talk...



Motivation and Datasets in Egocentric Video Understanding



Teaser: The Wizard of Oz & Genie 3



Video Understanding Out of the Frame



Outlook into the Future of Egocentric Vision



Point Tracking



Object Tracking



Gaze Priming



Hand Tracking



Conclusion

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