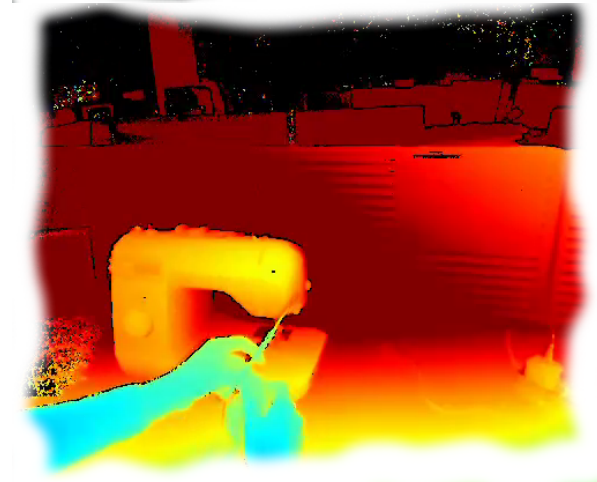


Egocentric Vision

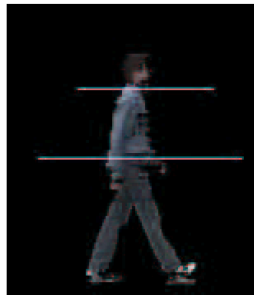
Dr Dima Damen

Department of Computer Science

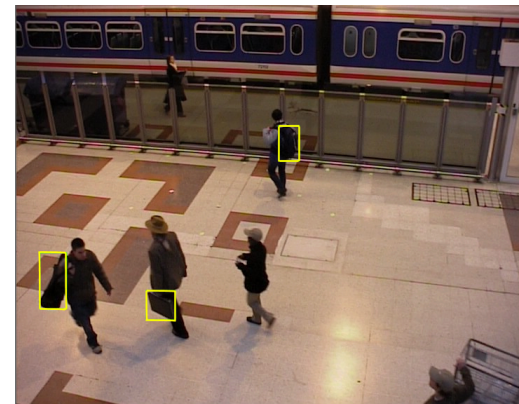


Short Bio

- 1998-2002 BSC in Computer Science
- 2002-2003 MSc in Distributed Multimedia Sys.

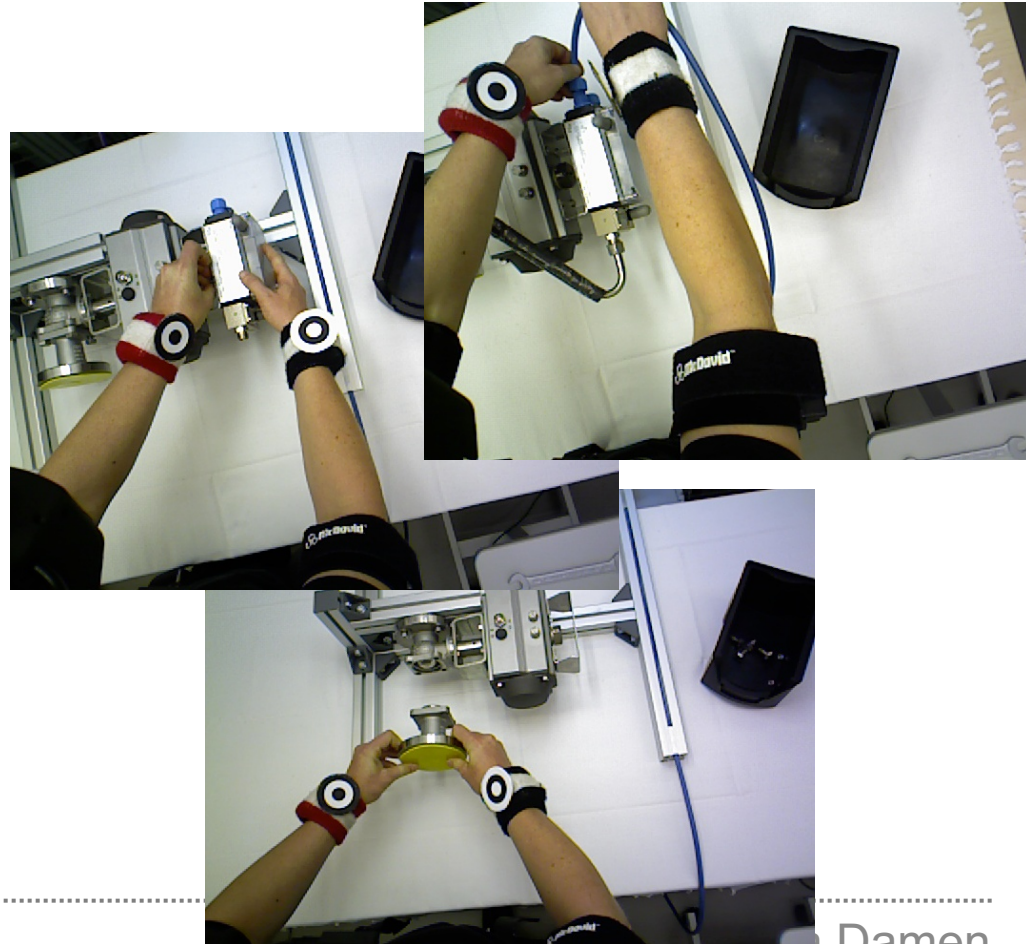


- 2006-2009 PhD in Computer Vision



Short Bio

- 2010-2012 Postdoc on EU-FP7 project



Short Bio

- 2013- Assistant Prof in Computer Vision

Egocentric Vision?

- Research interests: action and activity recognition
- Particularly centred around the viewpoint or the perspective

Ego...

*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*

Ego...



Visual Sensing – the landscape



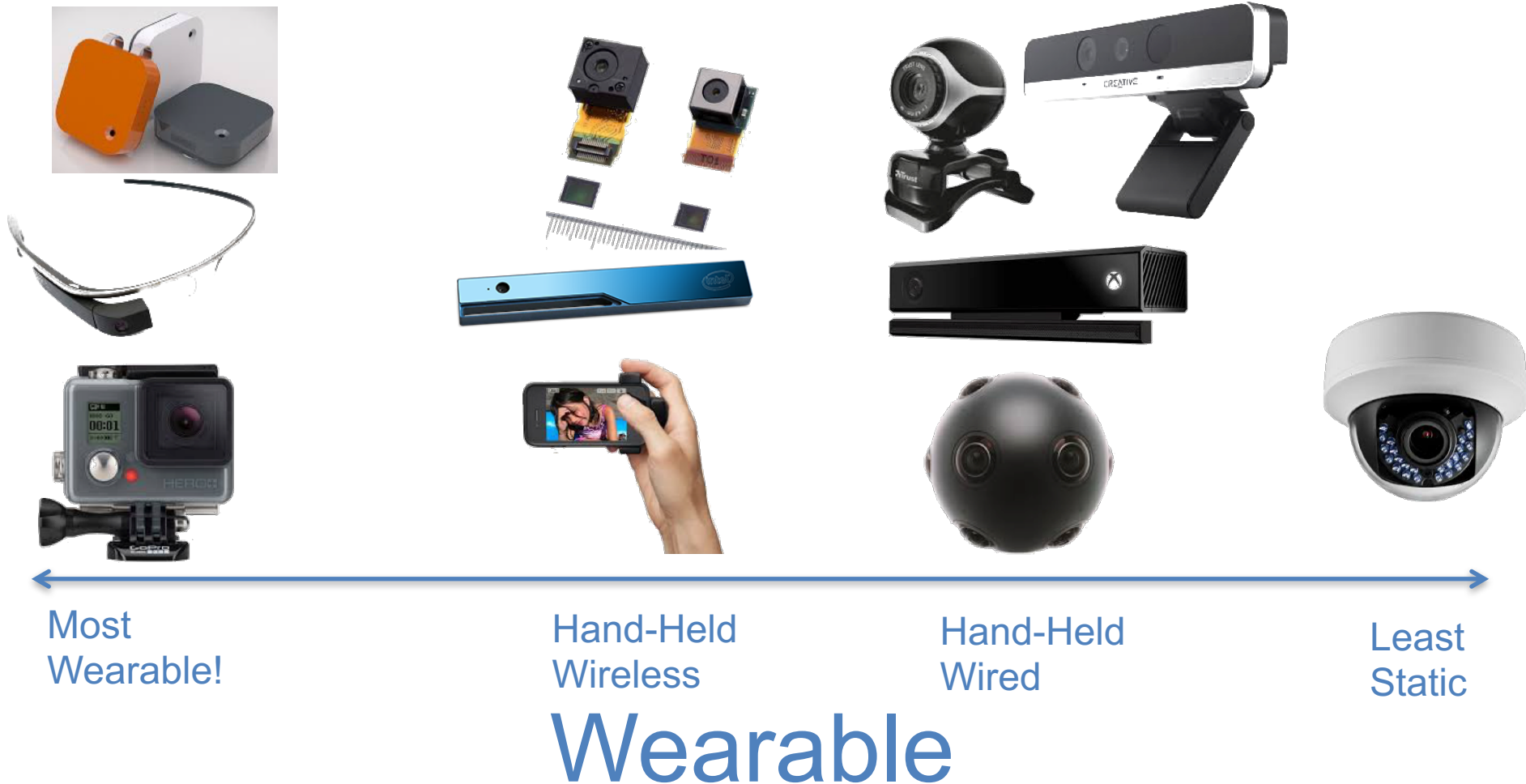
Visual Sensing – the landscape



Visual Sensing – the landscape



Visual Sensing – the landscape



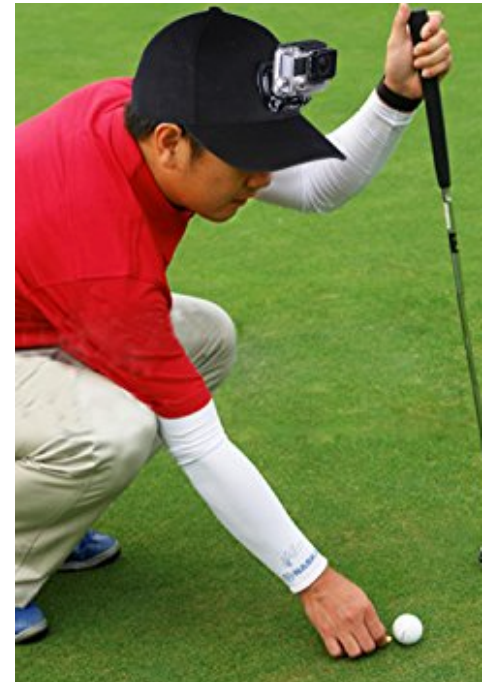
Wearable?



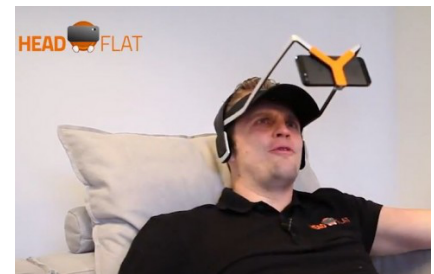
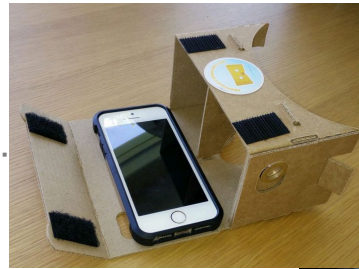
Wearable?



Wearable?



Wearable?

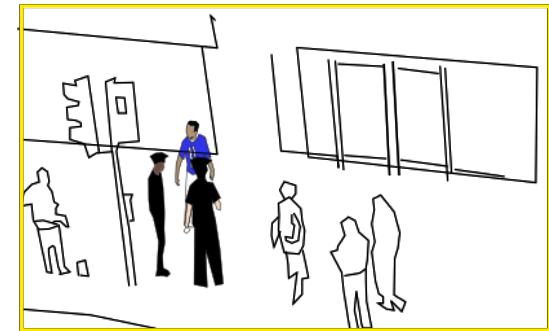
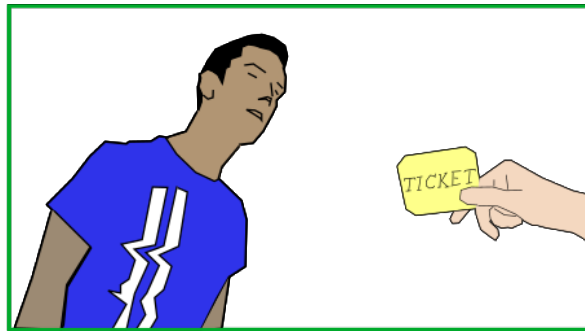


Wearable?

- Hat-Mounted
- Head-Mounted 
- Glass-Mounted 
- Shoulder-Mounted
- Chest-Mounted 
- Wrist-Mounted
- Belt-Mounted
- Ankle-Mounted

But why do we care about... hardware???

- OPV (Ordinal-Person Views)
 - FPV (First-Person View)
 - SPV (Second-Person View)
 - TPV (Third-Person View)



See for yourself!

- [Videos...](#)

Conclusions?

- Just another camera?
- Just a shaking camera?

Egocentric Vision

- The Unique Problems
 1. Camera Motion
 2. Mapping and Localisation (ref tomorrow's talk)
 3. Attention and Task-Relevance
 4. Object Interactions
 5. Multi-view Solutions
- The Unique Applications
 1. Video Summarisation
 2. Skill Determination
 3. Real-time solutions

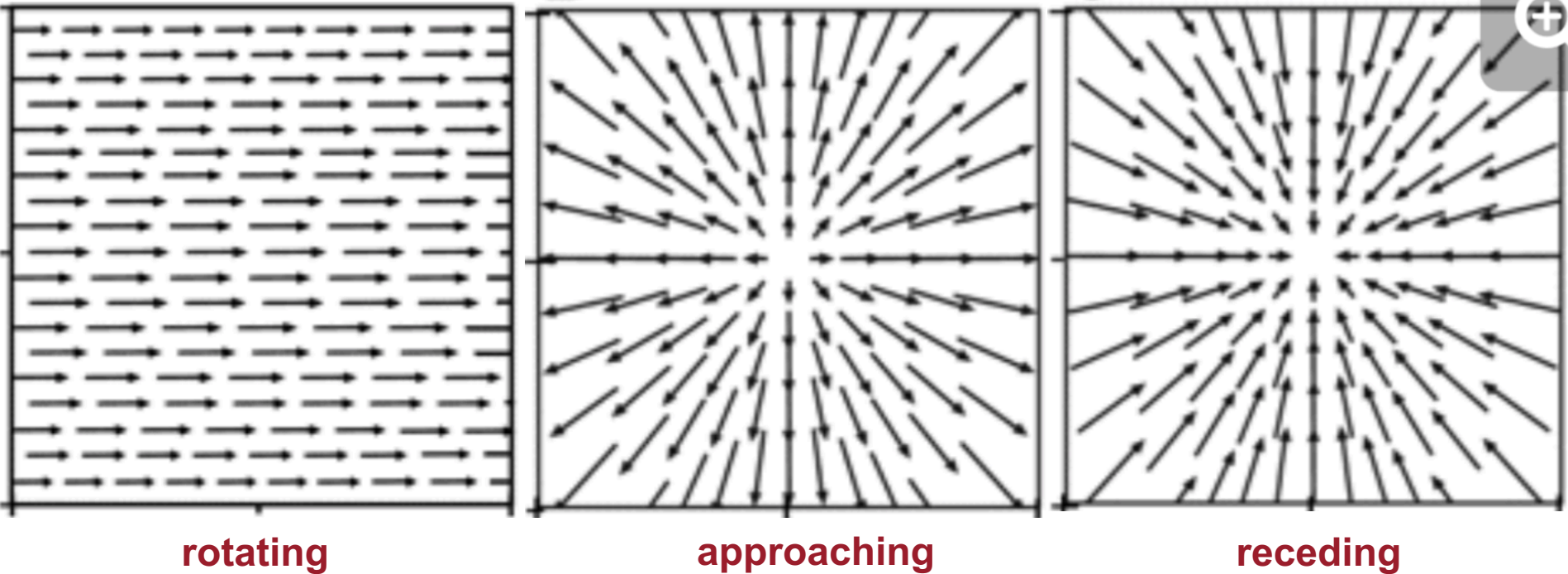
The Unique Problems

1. Camera Motion

1. Camera Motion

- Two types of motion
 - Egomotion
 - Foreground motion

Ego-motion



Ego-motion

- Detect to:
 - Use?
 - Remove?

Hyperlapse

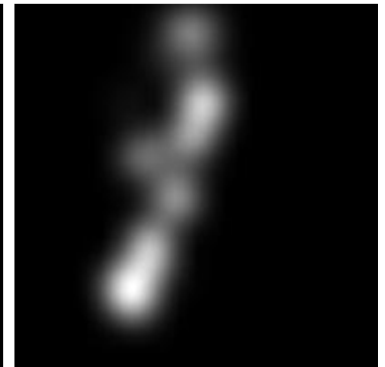
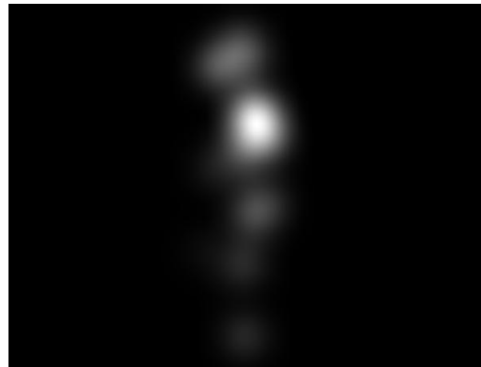
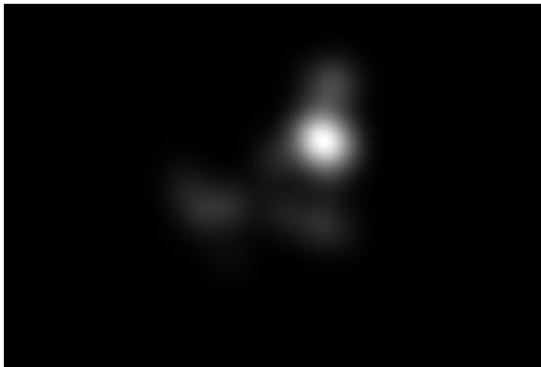
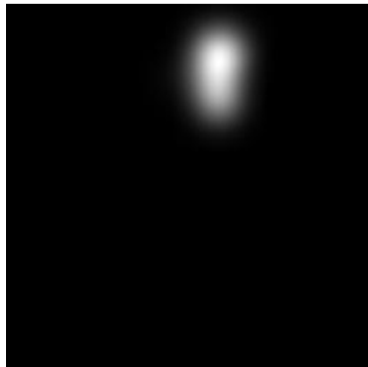
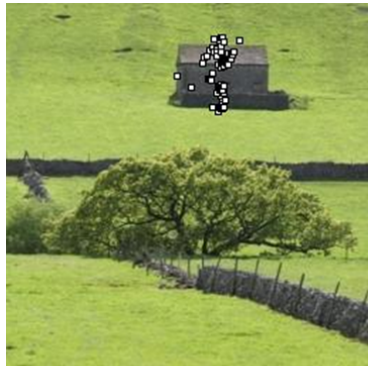
- <https://youtu.be/sA4Za3Hv6ng>

The Unique Problems

3. Attention and Task Relevance

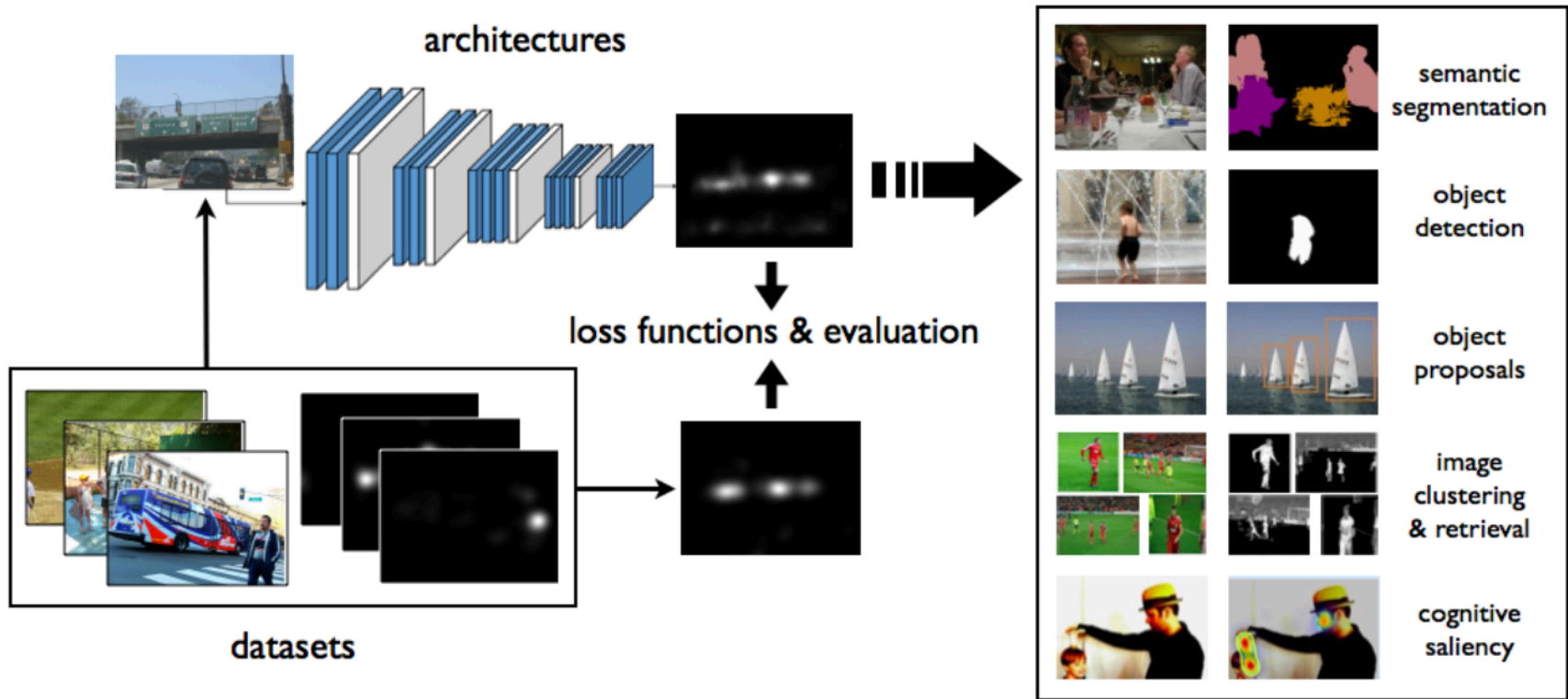
Attention and Task Relevance

- What is attention?
 - Non-Egocentric Attention Models (→ Saliency)



Attention and Task Relevance

- What is attention?
 - Non-Egocentric Attention Models (→ Saliency) applications



Attention and Task Relevance



Attention and Task Relevance

- Attention in egocentric vision
 - Foreground segmentation
 - Hand-region segmentation
 - Gaze tracking



Quick introduction to human gaze

- Humans iterate between “fixations” and “saccades”
 - Fixation: short stops
 - Saccade: quick movements between fixations
- <https://youtu.be/pknohrs4Qs>

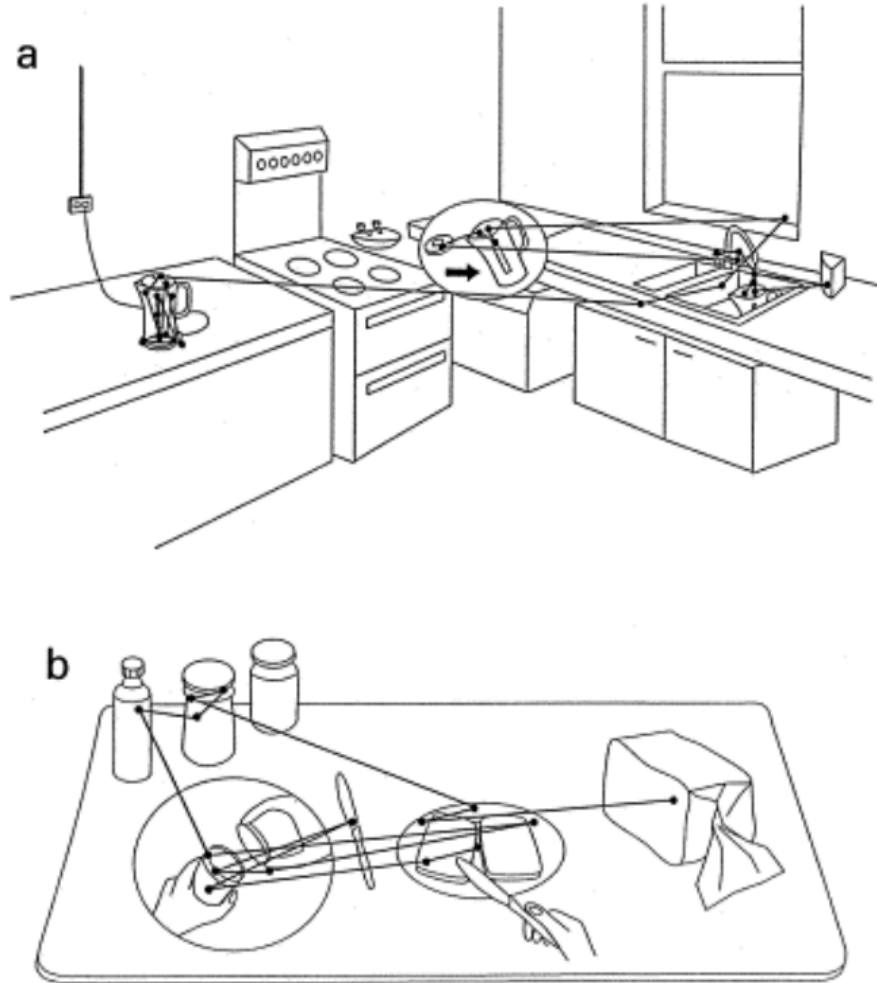
Quick introduction to human gaze



Quick introduction to human gaze

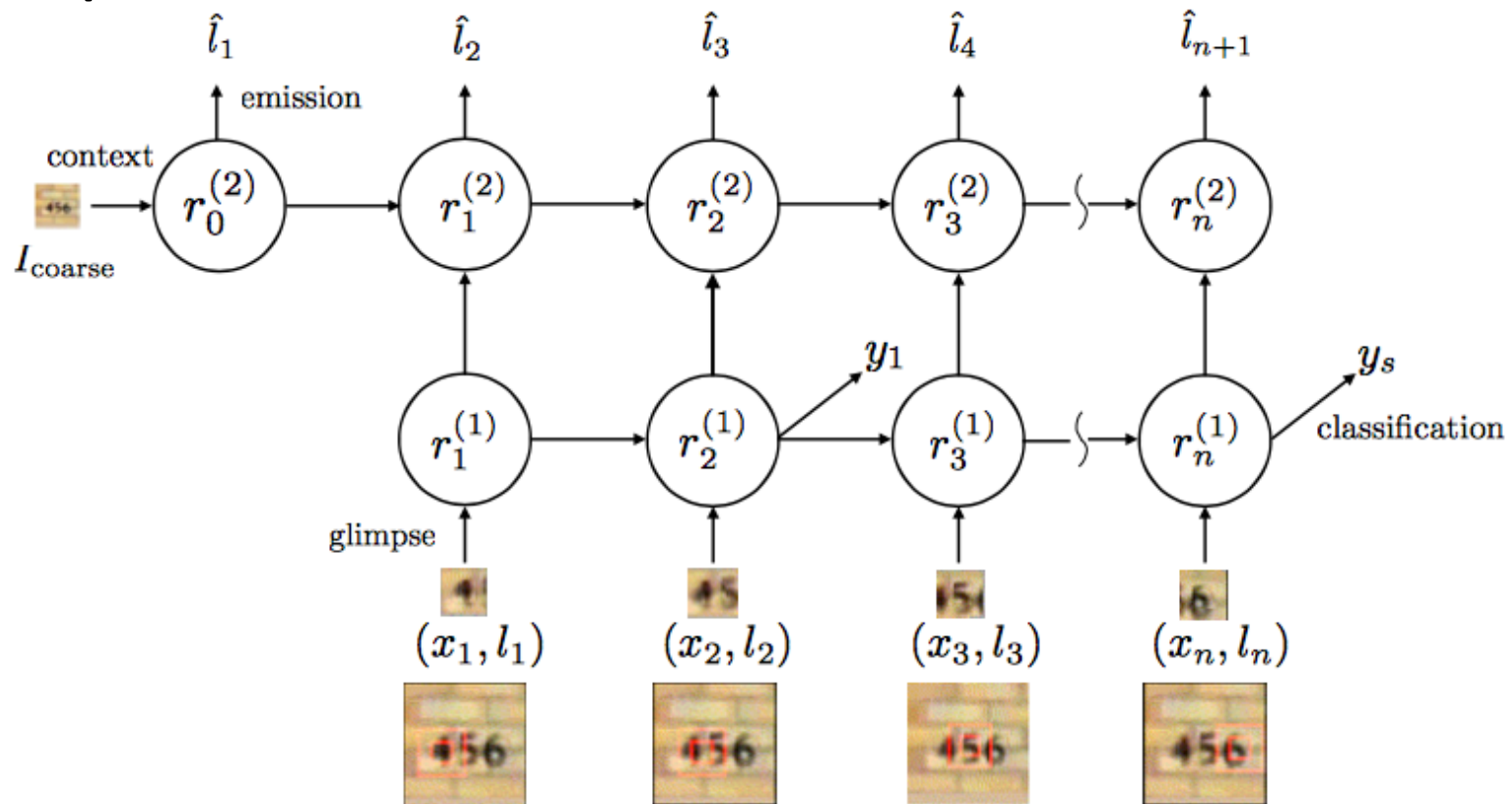


Quick introduction to human gaze



Quick introduction to human gaze

- The notion of fixation/saccade has recently inspired attention models in vision



Quick introduction to human gaze

Gaze Information to prime Object Detection

Dima Damen, Osian Haines, Andrew Calway and Walterio Mayol-Cuevas.

Object detection is based on the paper:

Dima Damen, Pished Bunnun, Andrew Calway and Walterio Mayol-Cuevas.

Real-time Learning and Detection of 3D Texture-less Objects: A Scalable Approach.

British Machine Vision Conference (BMVC), 2012. **[Best Poster Paper]**

Jan. 2013



The Unique Problems

3. Attention and Task Relevance

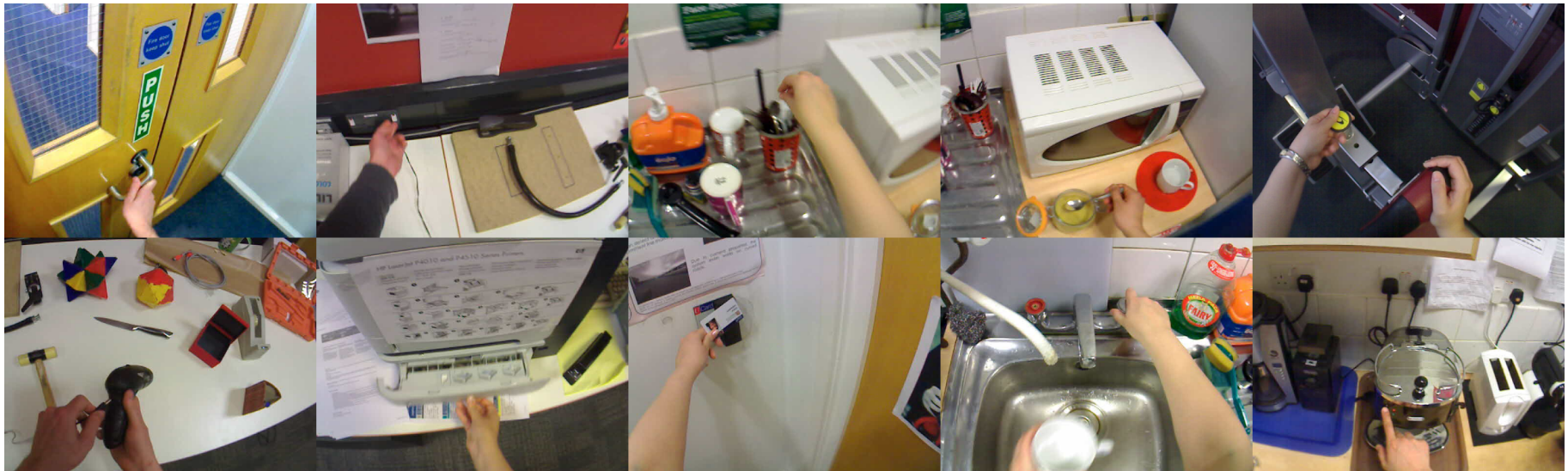
Case Study: You-Do, I-Learn

You-Do, I-Learn

- First-person view
- Offers a unique insight into ‘used’ or ‘attended-to’ objects
- How these objects have been used

BEOID Dataset

- Released July 2014
- Wearable gaze tracker (ASL Mobile Eye XG)
- 6 locations: kitchen, workspace, printer, corridor with locked door, cardiac gym and weight-lifting machine
- 5 operators (2 sequences each)



BEOID Dataset

- Q. How to 'ground-truth' objects that have been used?
- Q. How to 'ground-truth' how these objects have been used?

Try it yourself



BEOID

- Ground-truth by written narration
- Released with dataset

pick the charger and plug it into the socket. Check that the screwdriver is powered by looking at the button. Pick the tape and place it in the box. Walk to the printer. Open the drawer to check the paper, and press keys on the printer pad. Use the card to unlock the door

You Do, I Learn

- Discover used objects
- Discover how objects have been used
- Extract guidance videos
- Fully unsupervised
 - No prior knowledge of objects (number, size)
 - Static and moveable objects

Definition

Task-Relevant Object (TRO)

an object, or part of an object, with which a person interacts during task performance

Which Objects?



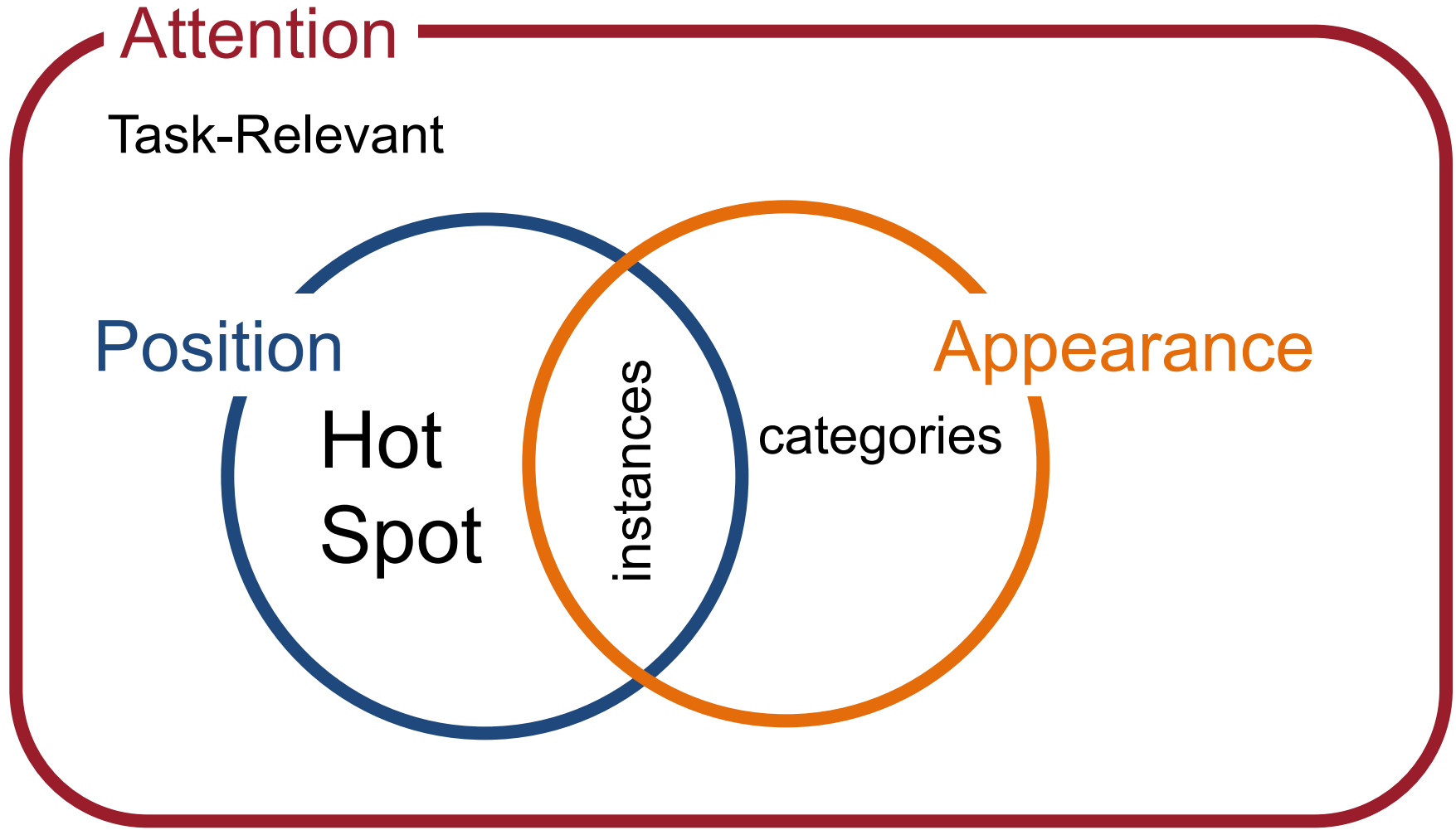
Discovering Task-Relevant Objects



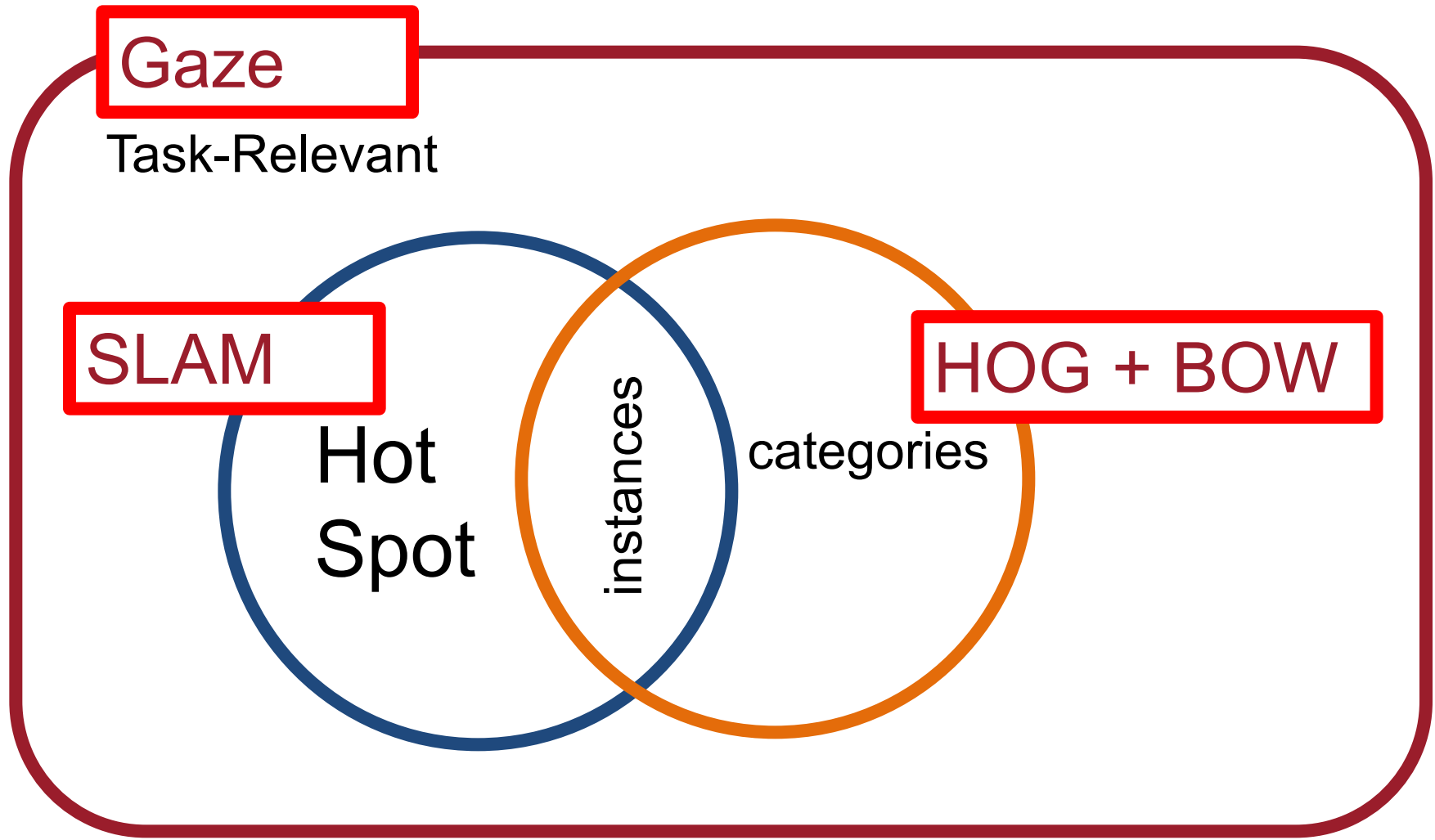
Discovering Task-Relevant Objects

- *Suggested* Problem Formulation...
 - Given a sequence of egocentric images $\{I_1, \dots, I_T\}$
 - Collected from multiple operators around a common environment
 - Automatically discover all task-relevant objects
 $\{O_k; 1 \leq k \leq K\}$
 $O_k = \{\Omega(I_t); 1 \leq t \leq T\}$
 - *Assumption:* at most one task-relevant image part is present within each image

Discovering Task-Relevant Objects



Discovering Task-Relevant Objects

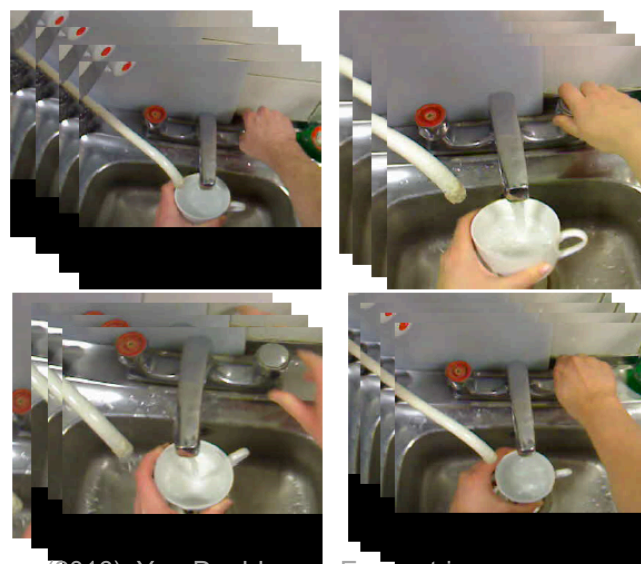
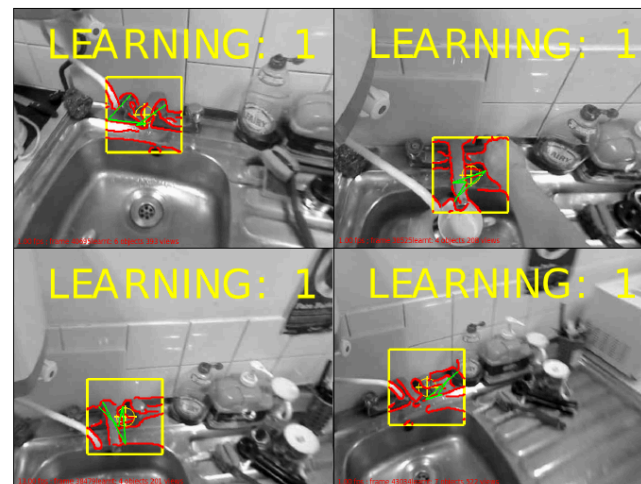
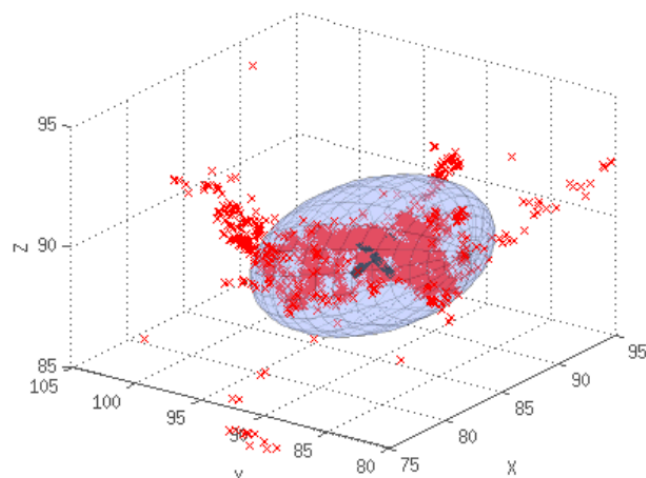


Discovering TROs

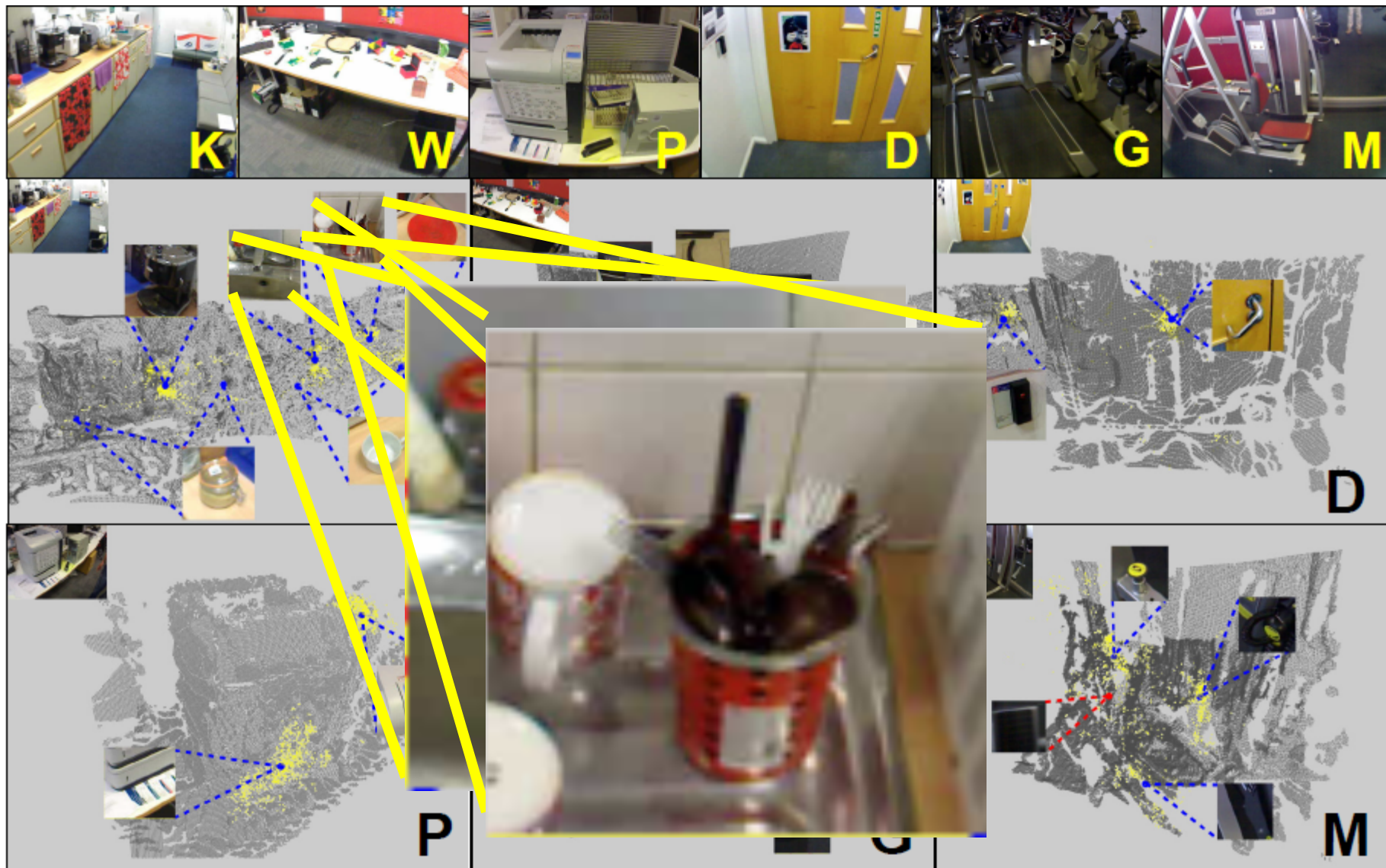
Discovering becomes a clustering task...

- Considers attention, position and appearance
- Unknown number of objects
- Davies-Bouldin (DB) index
- K-Means vs Spectral

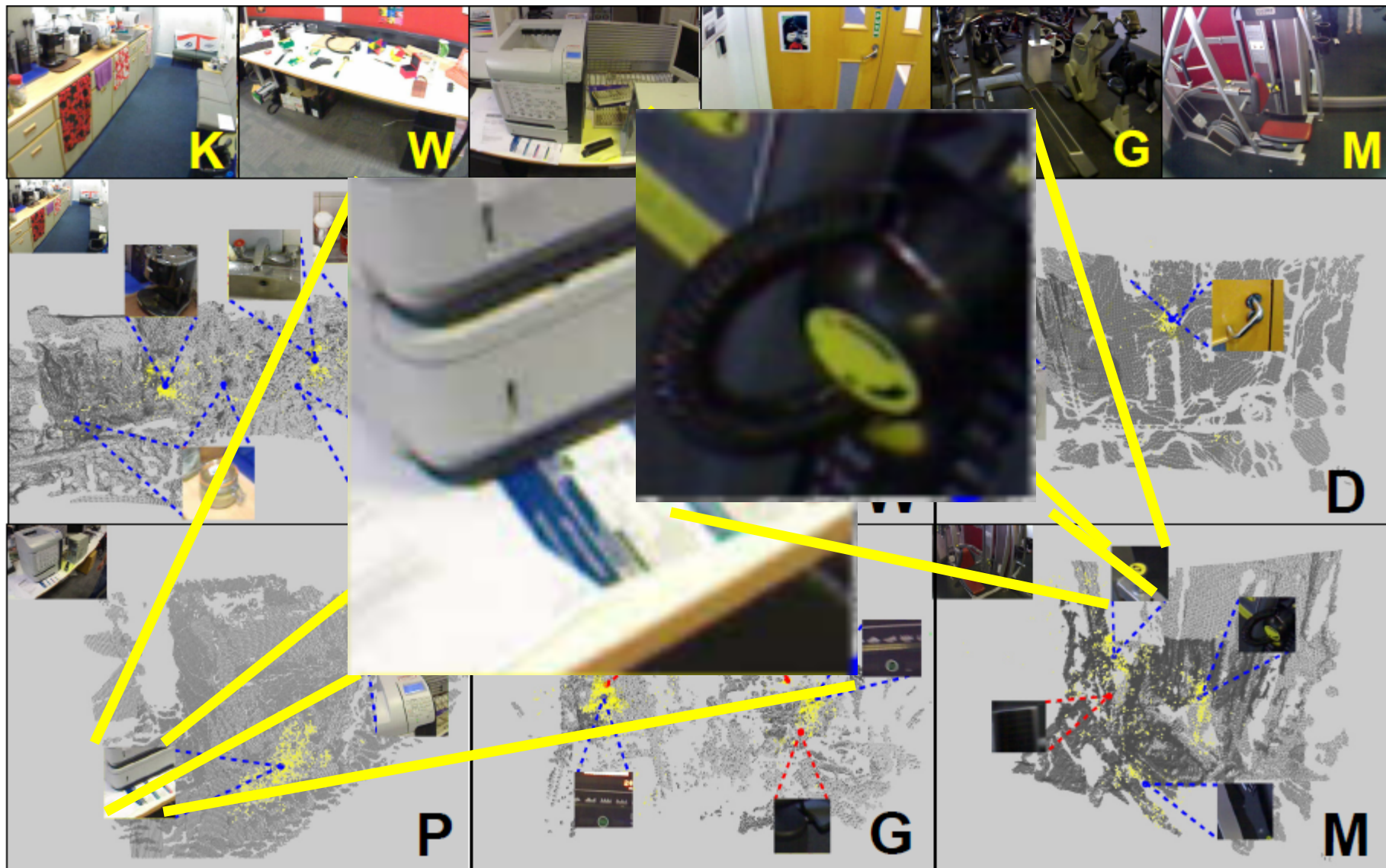
Discovering Task-Relevant Objects



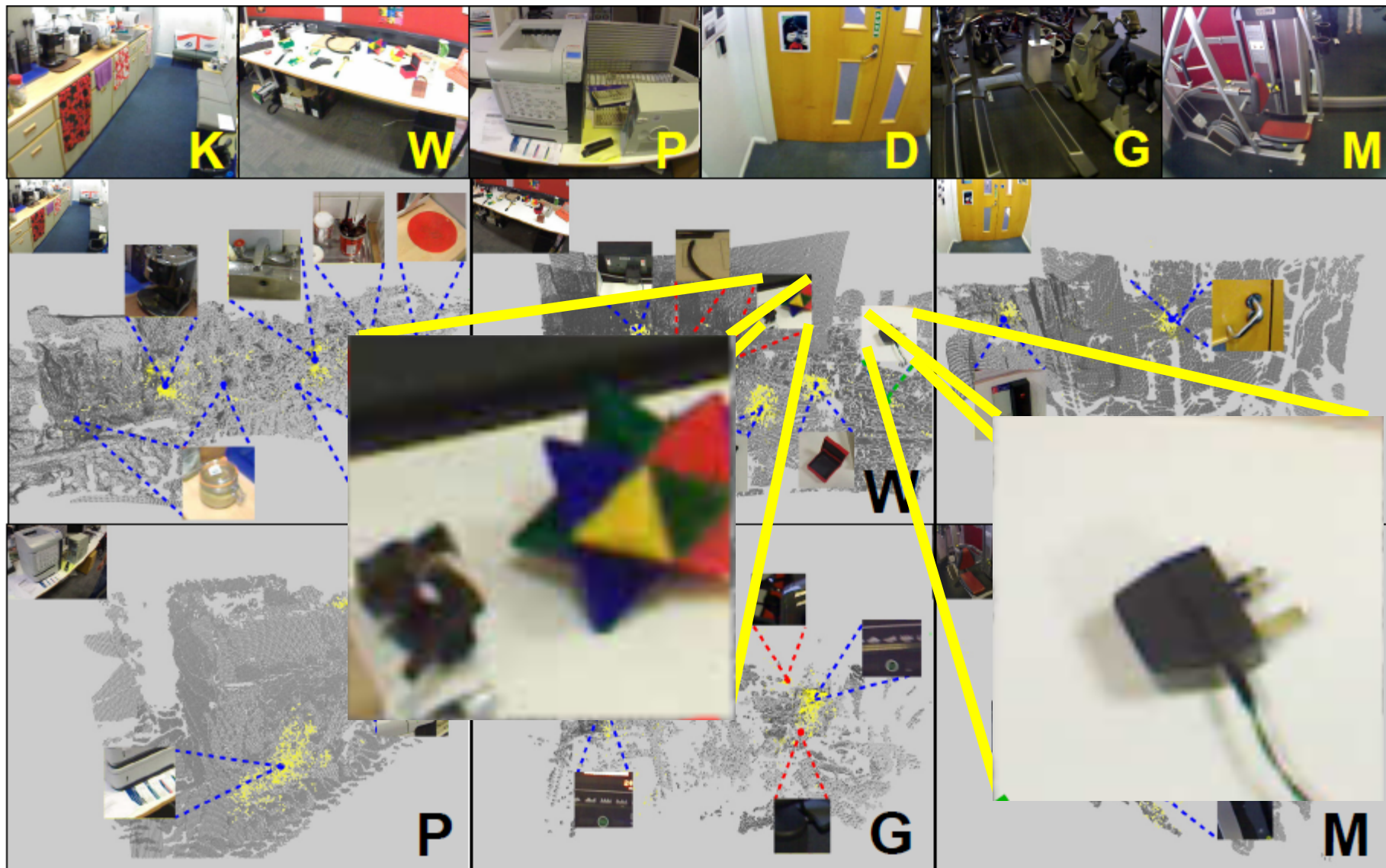
Discovering Task-Relevant Objects



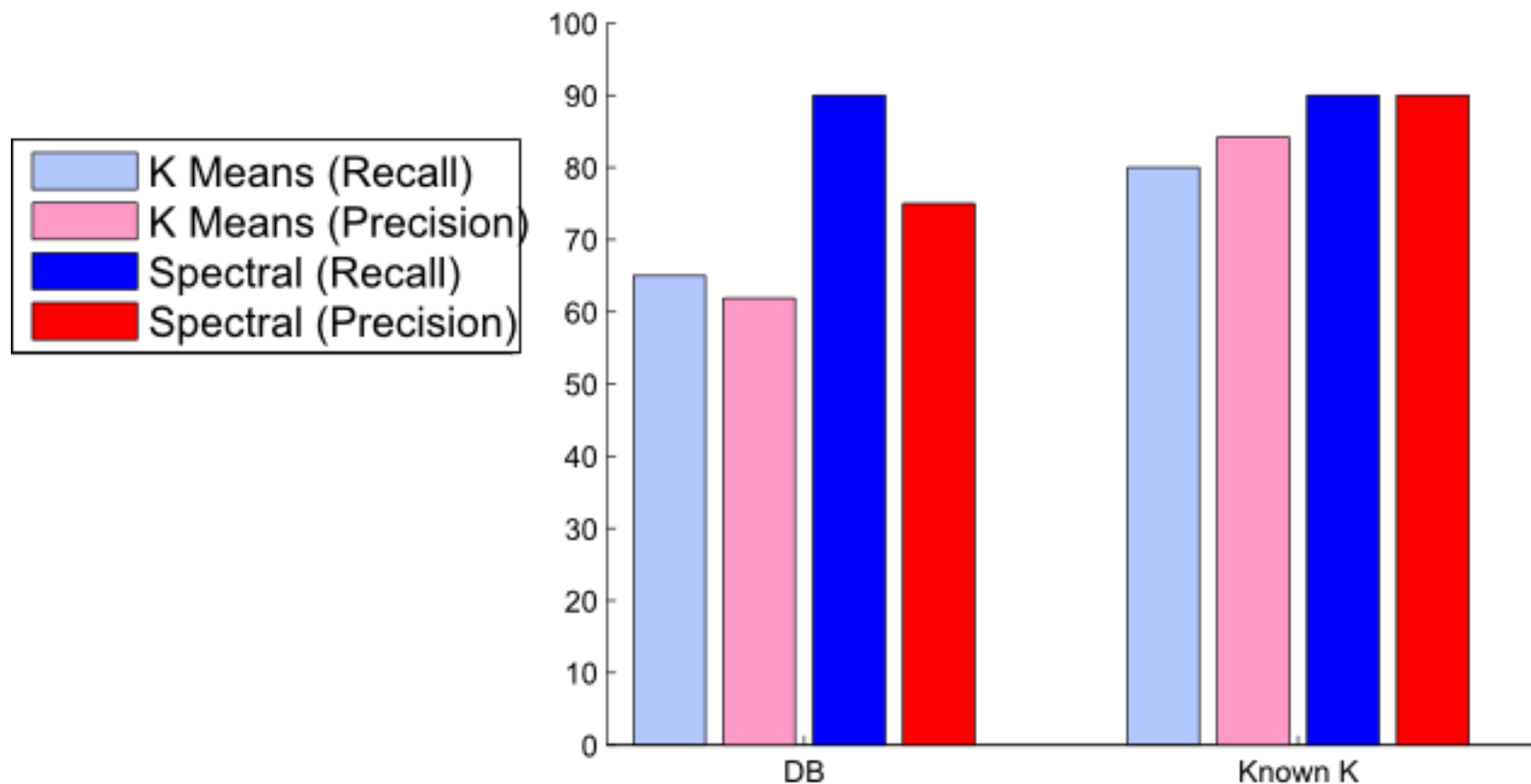
Discovering Task-Relevant Objects



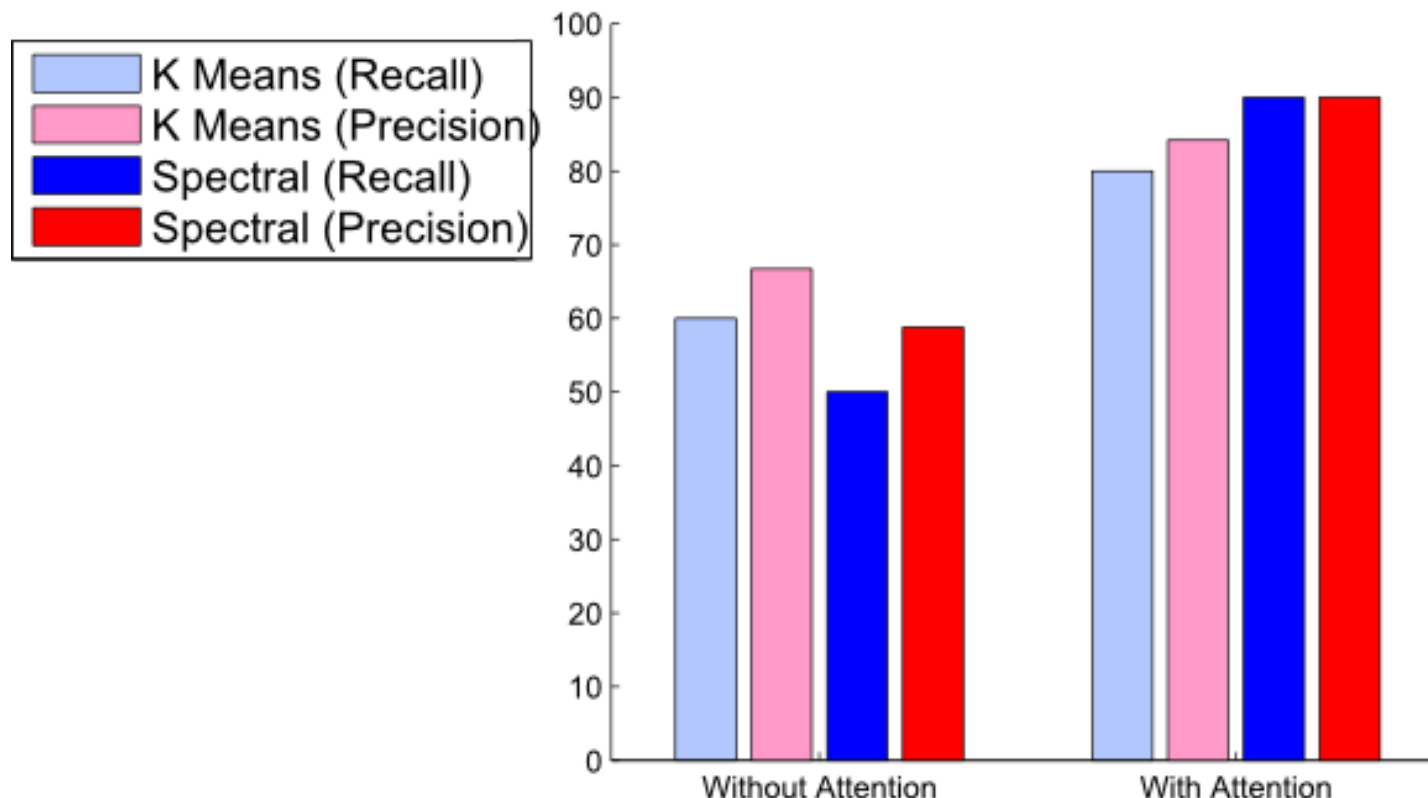
Discovering Task-Relevant Objects



Discovering Task-Relevant Objects



Discovering Task-Relevant Objects



Discovering Modes of Interaction

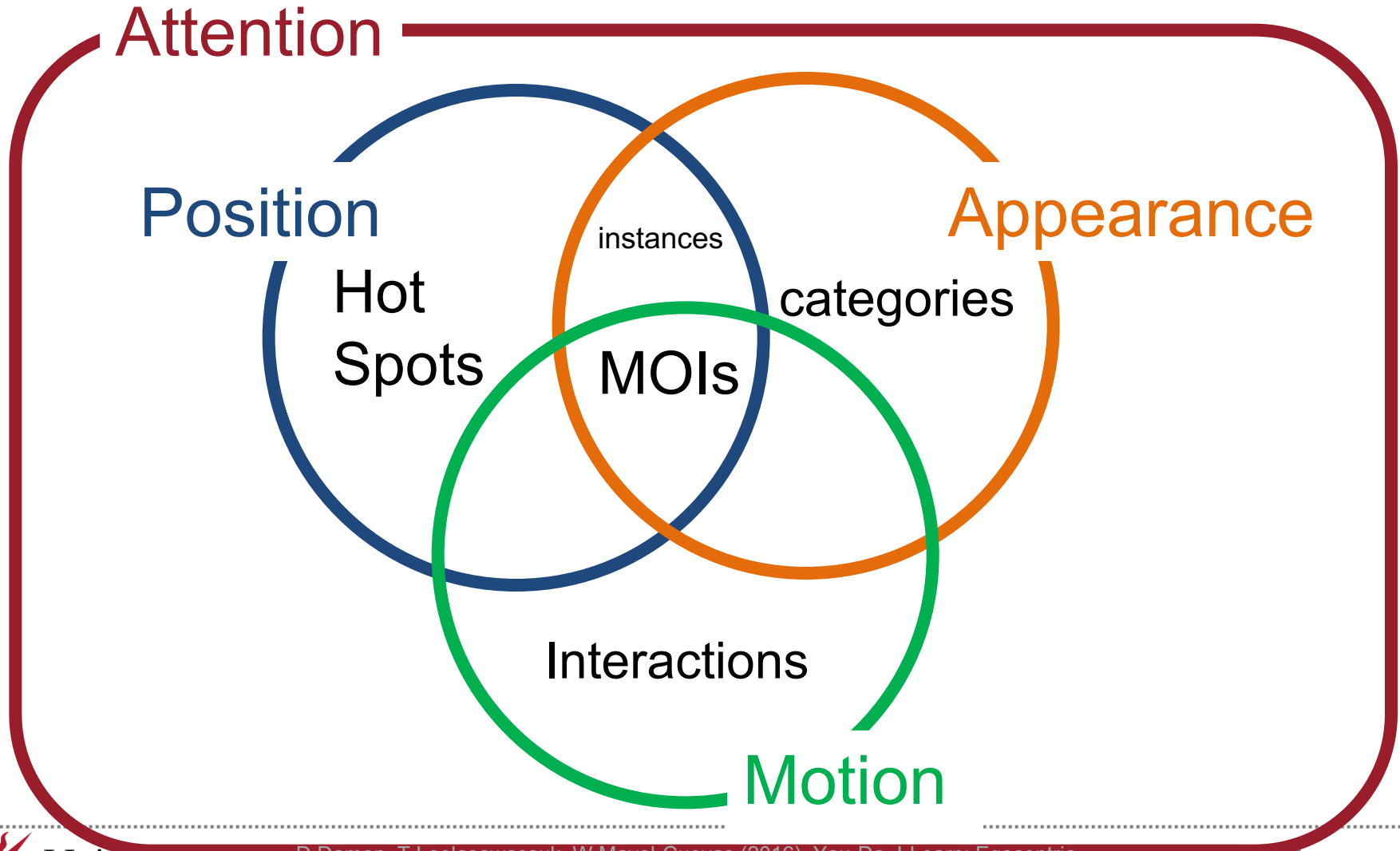


Definition

Modes of Interaction (MOI)

the different ways in which TROs are used

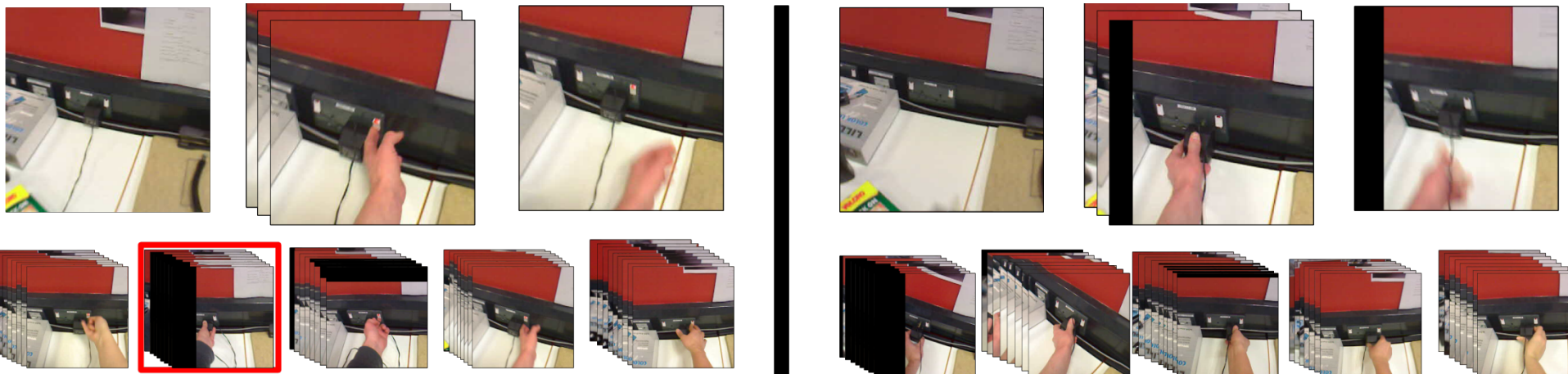
Discovering Modes of Interaction



Discovering Modes of Interaction

- Motion
 - Video snippets for each discovered object
 - Descriptor per snippet
 - Clustering using DB-index

Discovering Modes of Interaction



Discovering Modes of Interaction

Open & get sugar



Put



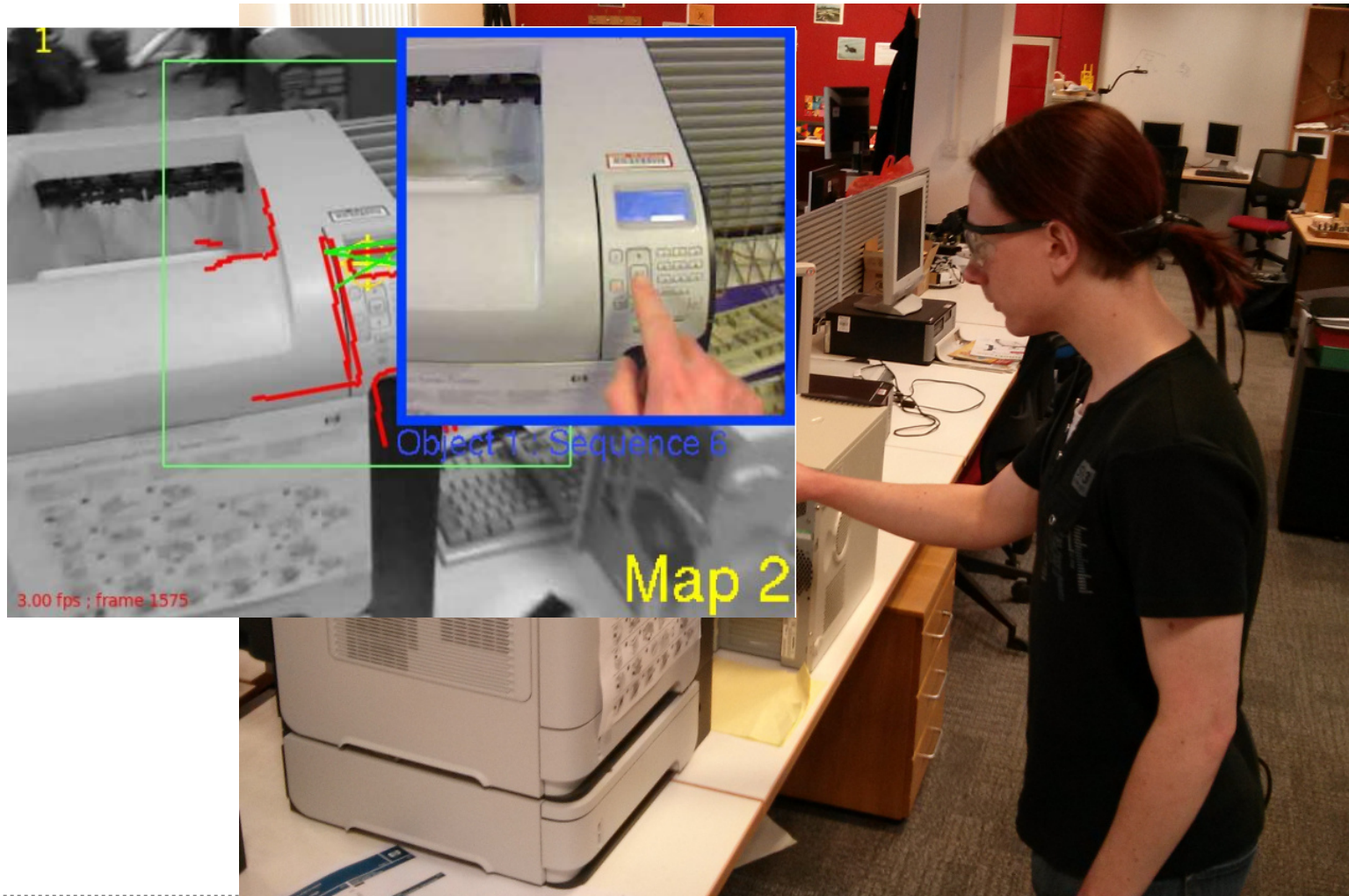
Pick



Open door



Back to.... the goal...



You Do, I Learn - Demonstration

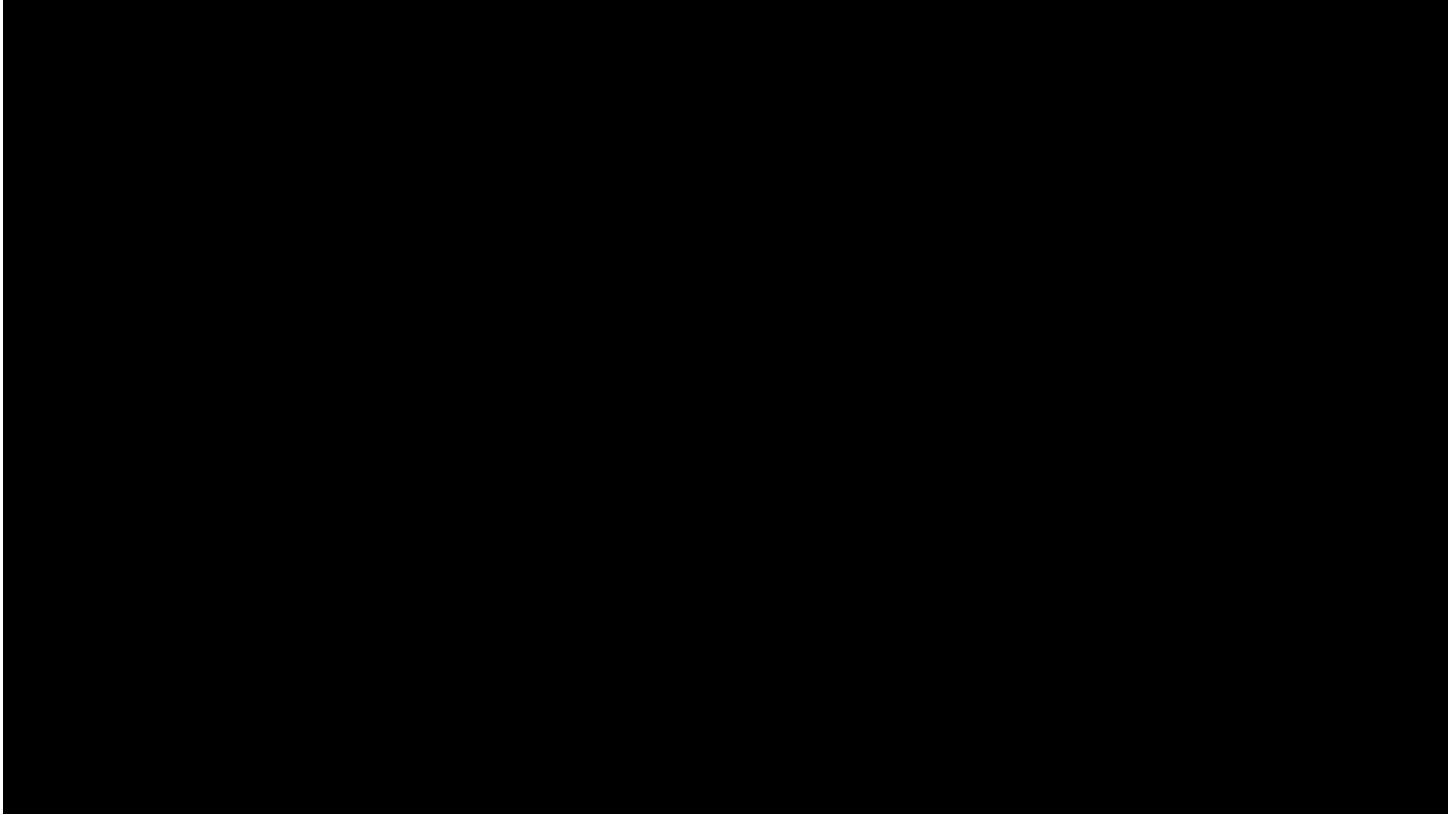


You Do, I Learn – Google Glass Prototype



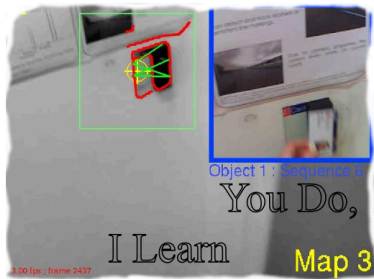
Task Monitoring - 2017

with: Longfei Chen
Kazuaki Kondo
Yuichi Nakamura
Walterio Mayol-Cuevas



More info...

Project You-Do, I-Learn



[Video1 \(2014\)](#), [Video2 \(2017\)](#)

Automated capture and delivery of assistive task guidance with an eyewear computer: The GlaciAR system. T Leelasawassuk, D Damen, W Mayol-Cuevas. Augmented Human, Mar 2017 [pdf](#)

You-Do, I-Learn: Discovering Task Relevant Objects and their Modes of Interaction from Multi-User Egocentric Video. D Damen, T Leelasawassuk, O Haines, A Calway, W Mayol-Cuevas. British Machine Vision Conference (BMVC), Sep 2014. [PDF](#) | [Abstract](#) | [Dataset](#)

Multi-user egocentric Online System for Unsupervised Assistance on Object Usage. D Damen, O Haines, T Leelasawassuk, A Calway, W Mayol-Cuevas. ECCV Workshop on Assistive Computer Vision and Robotics (ACVR), Sep 2014. [PDF Preprint](#)

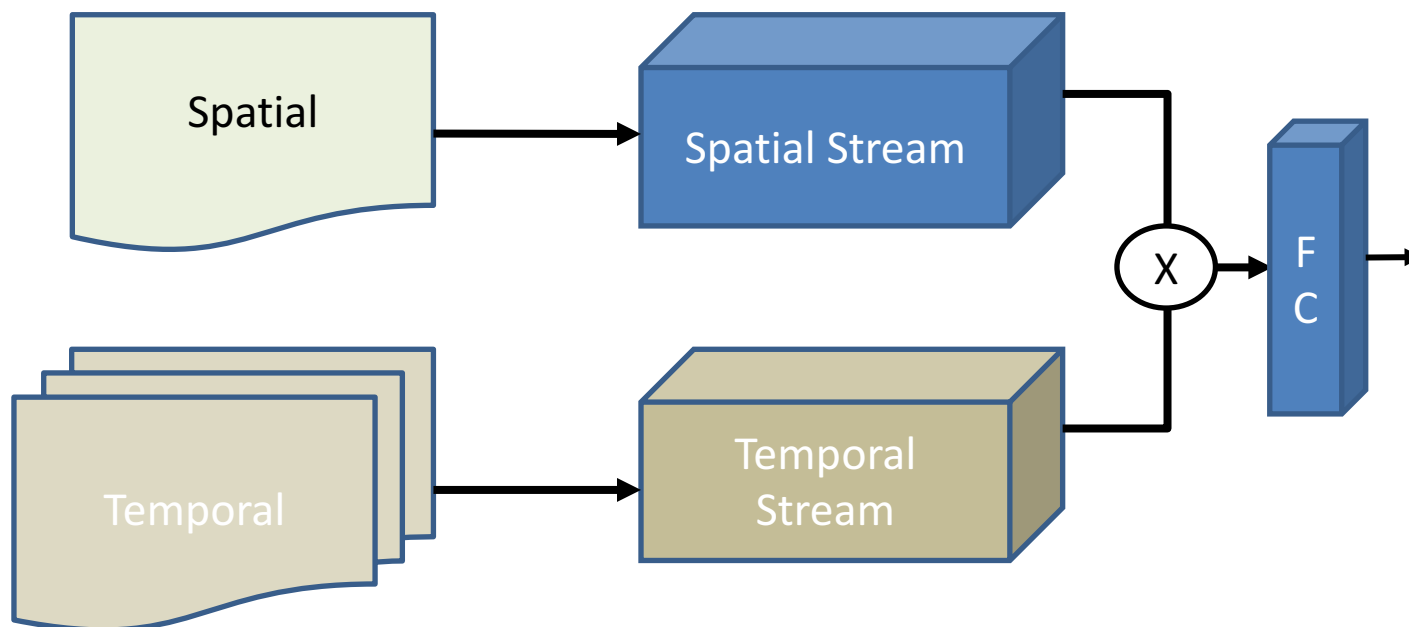
Estimating Visual Attention from a Head Mounted IMU. T Leelasawassuk, D Damen, W Mayol-Cuevas. International Symposium on Wearable Computers (ISWC), Sep 2015. [PDF](#)

The Unique Problems

4. Object Interactions

Action Recognition – an Introduction

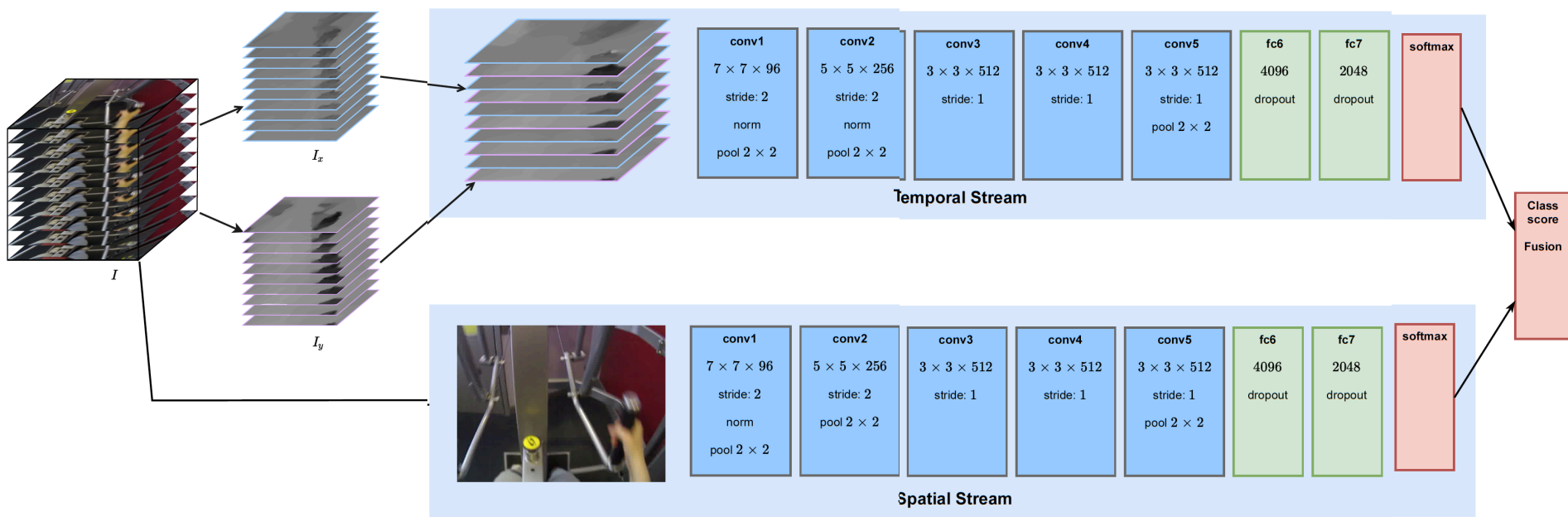
- CNNs for Action Recognition
 1. Dual-Stream Neural Networks



Action Recognition – an Introduction

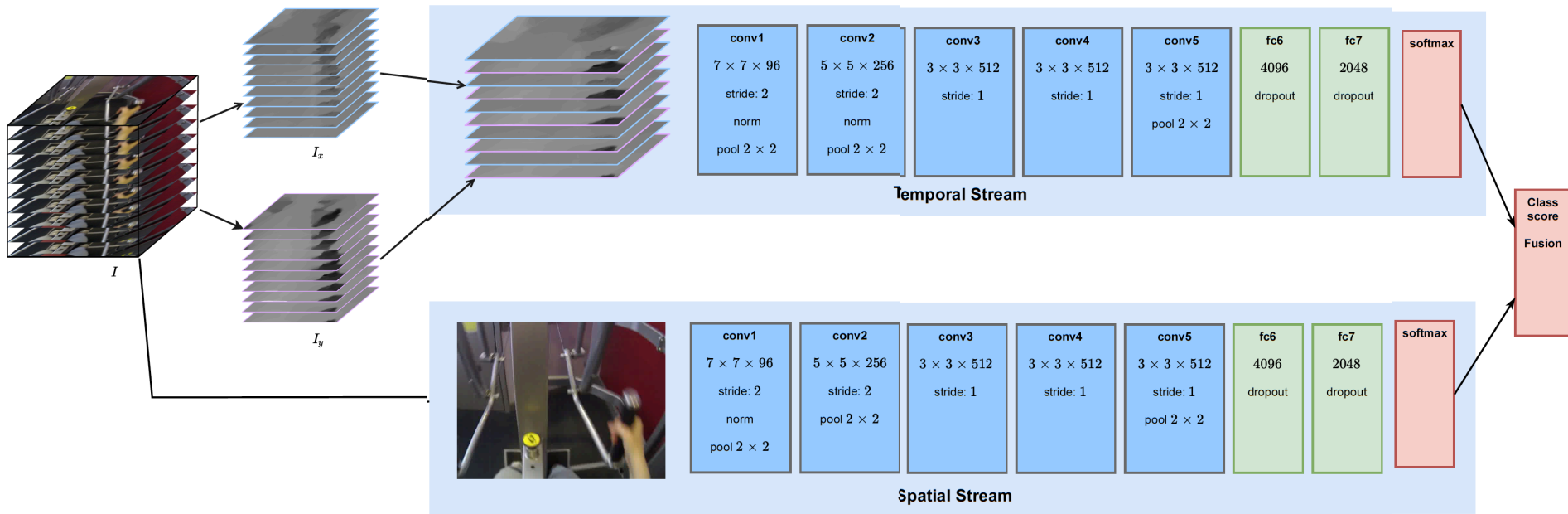
- CNNs for Action Recognition

1. Dual-Stream Neural Networks



Action Recognition – an Introduction

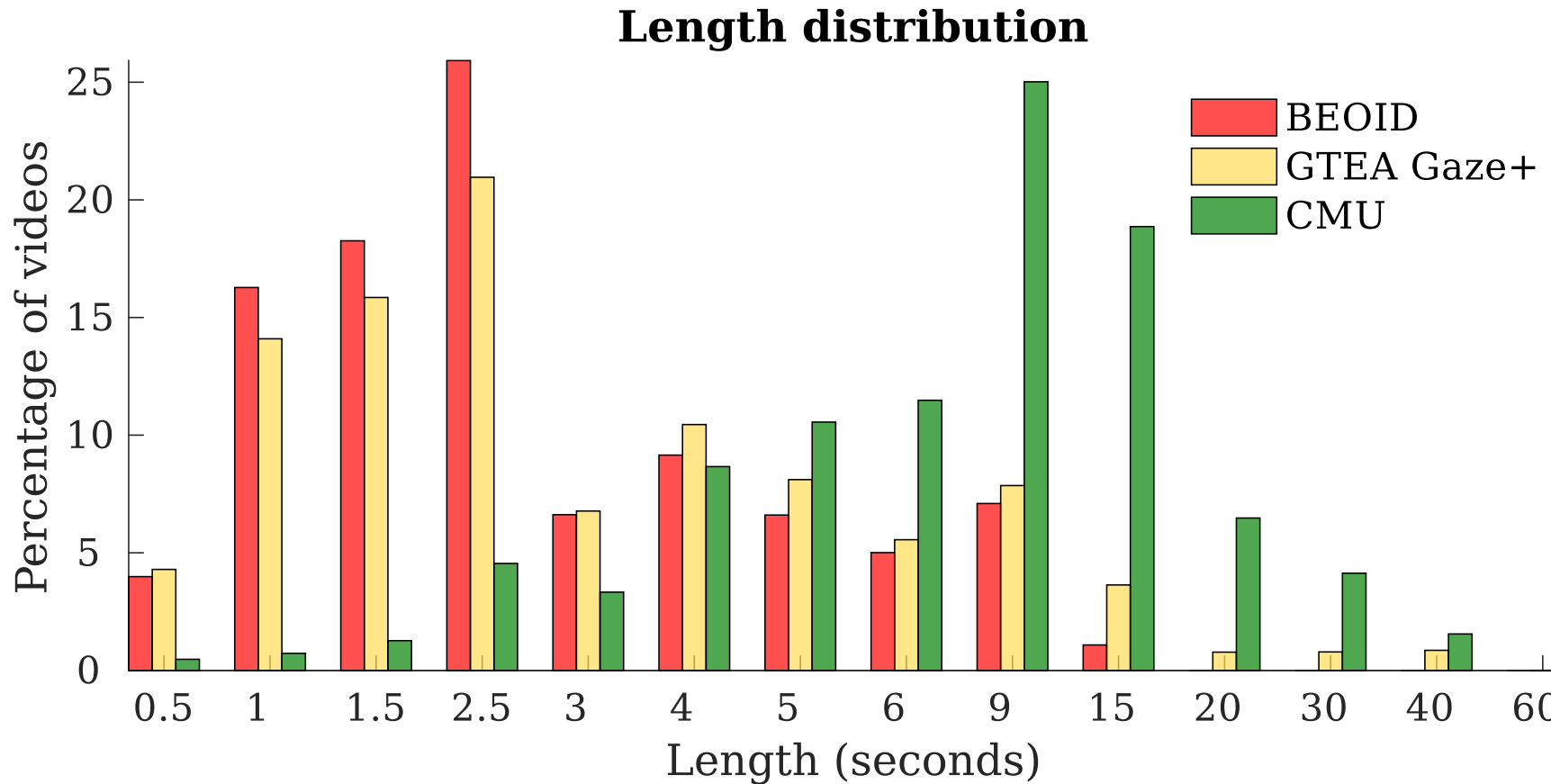
- CNNs for Action Recognition
 1. Dual-Stream Neural Networks



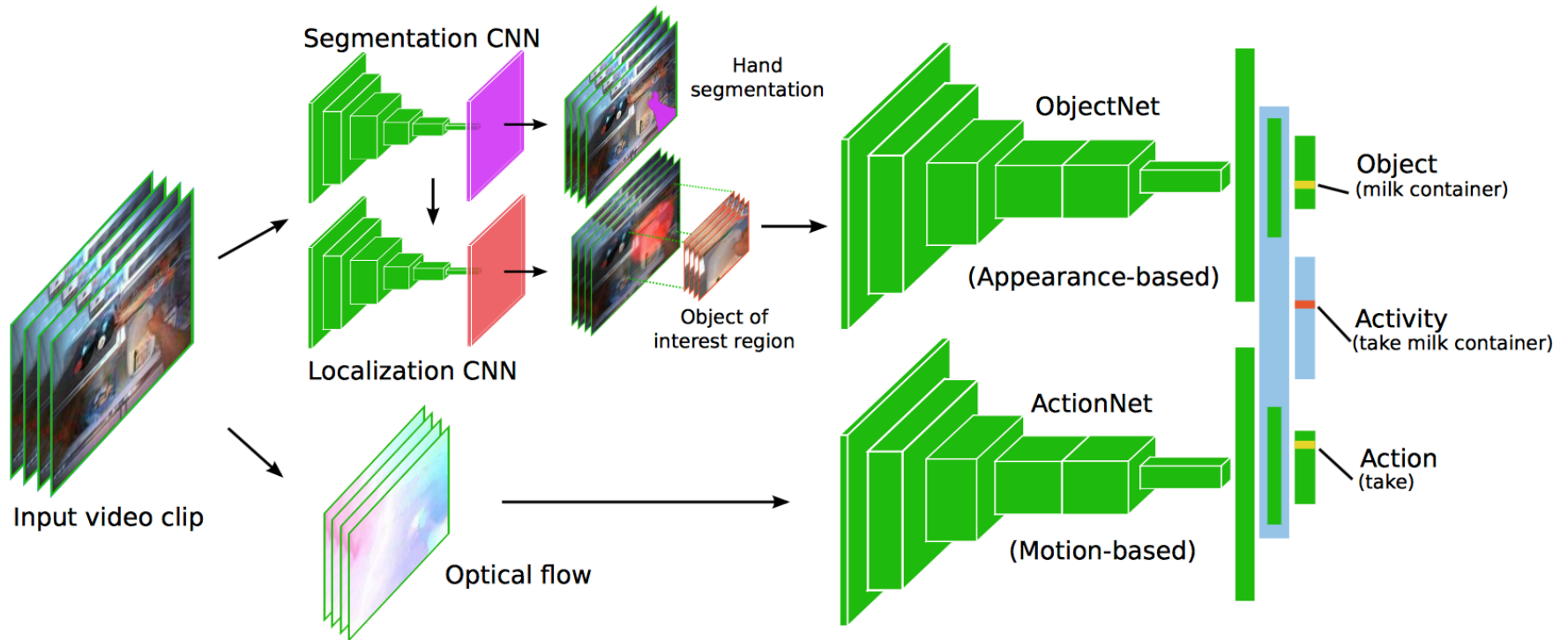
Action Recognition – an Introduction

Dataset	N. of <i>gt</i> segments	N. of <i>gen</i> segments	Classes
BEOID [1]	742	16691	34
GTEA Gaze+ [6]	1141	22221	42
CMU [2]	450	26160	31

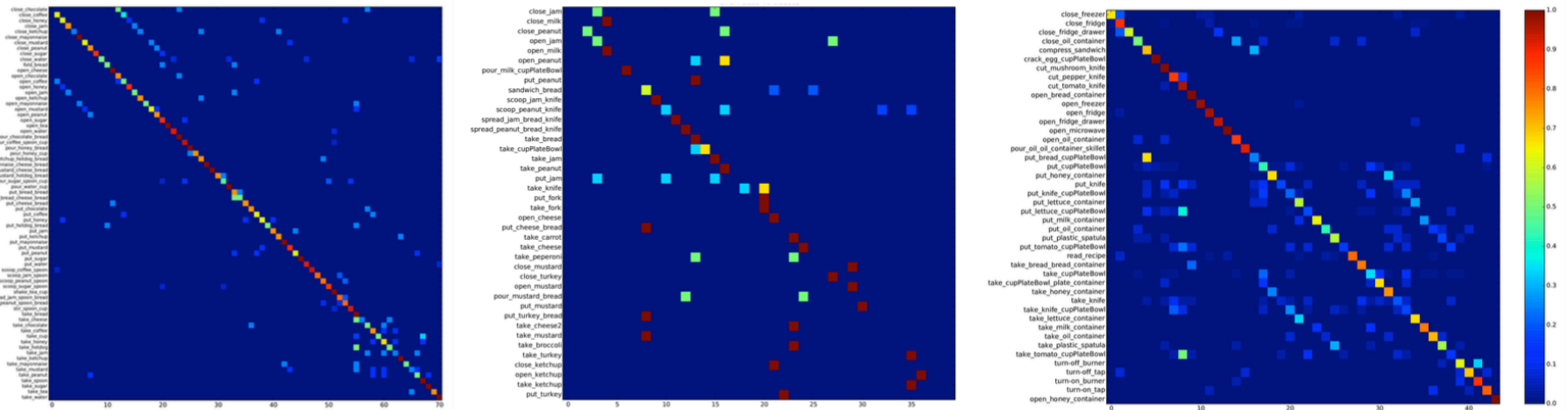
Action Recognition – an Introduction



Egocentric Action Recognition



Egocentric Action Recognition

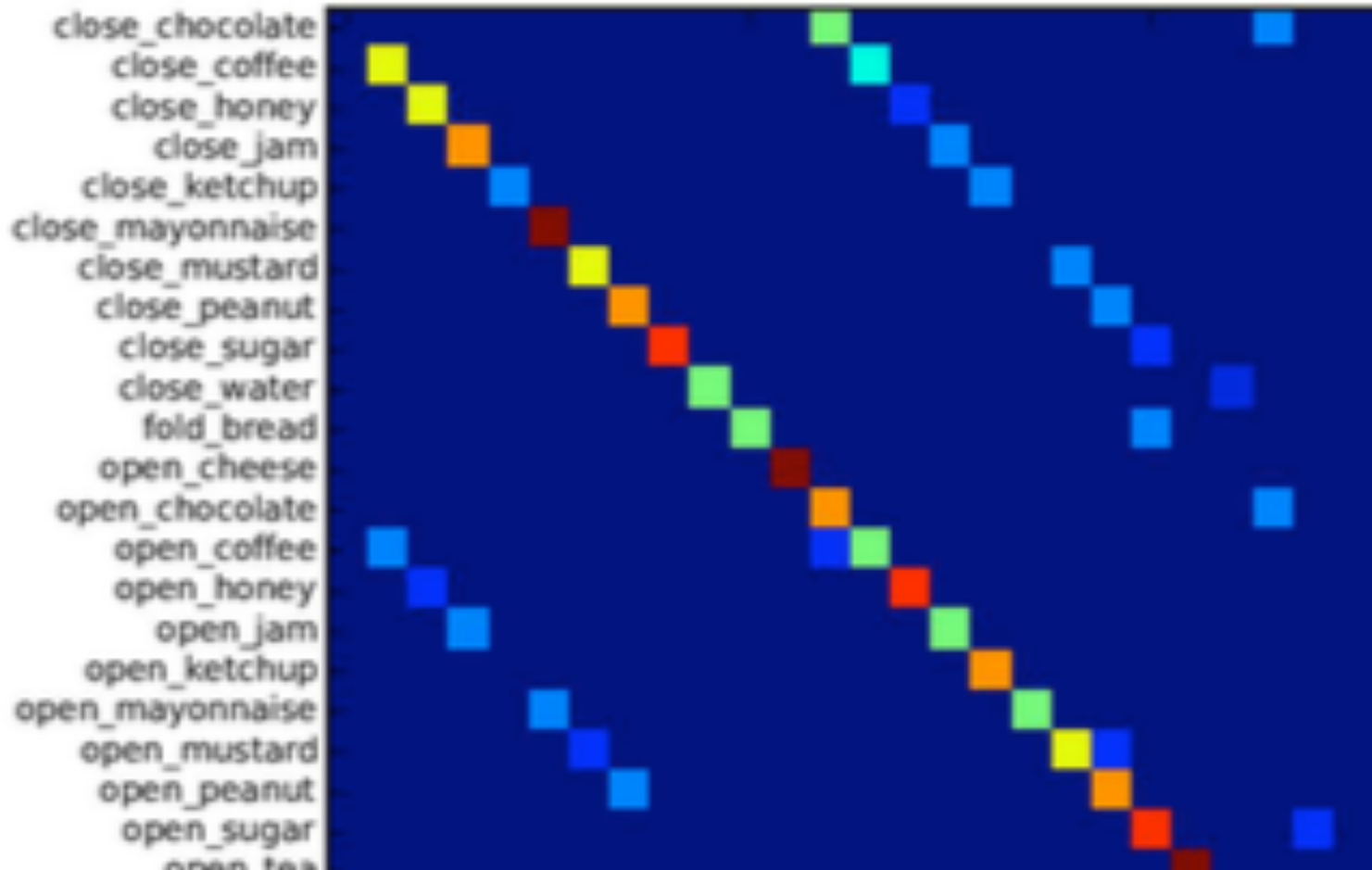


(a) GTEA 71 classes

(b) Gaze 40 classes

(c) Gaze+ 44 classes

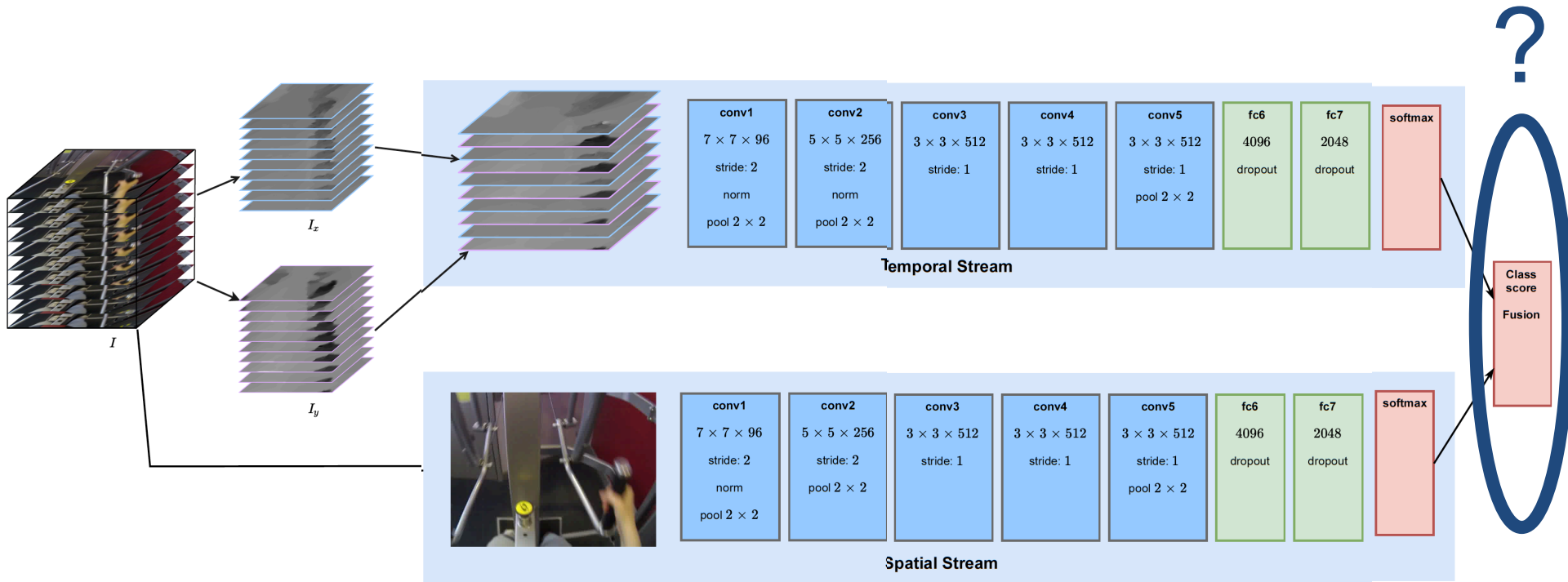
Egocentric Action Recognition



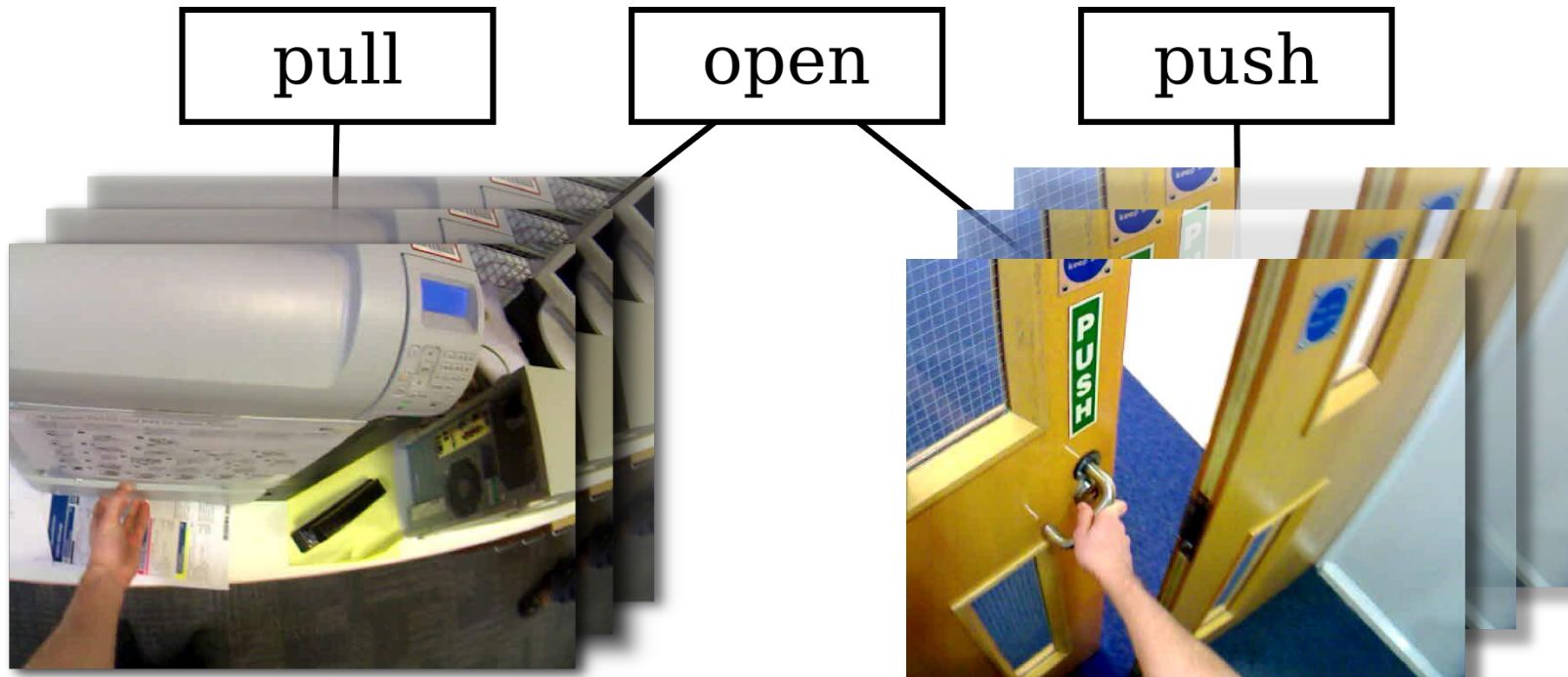
Action Recognition – an Introduction

- CNNs for Action Recognition

1. Dual-Stream Neural Networks

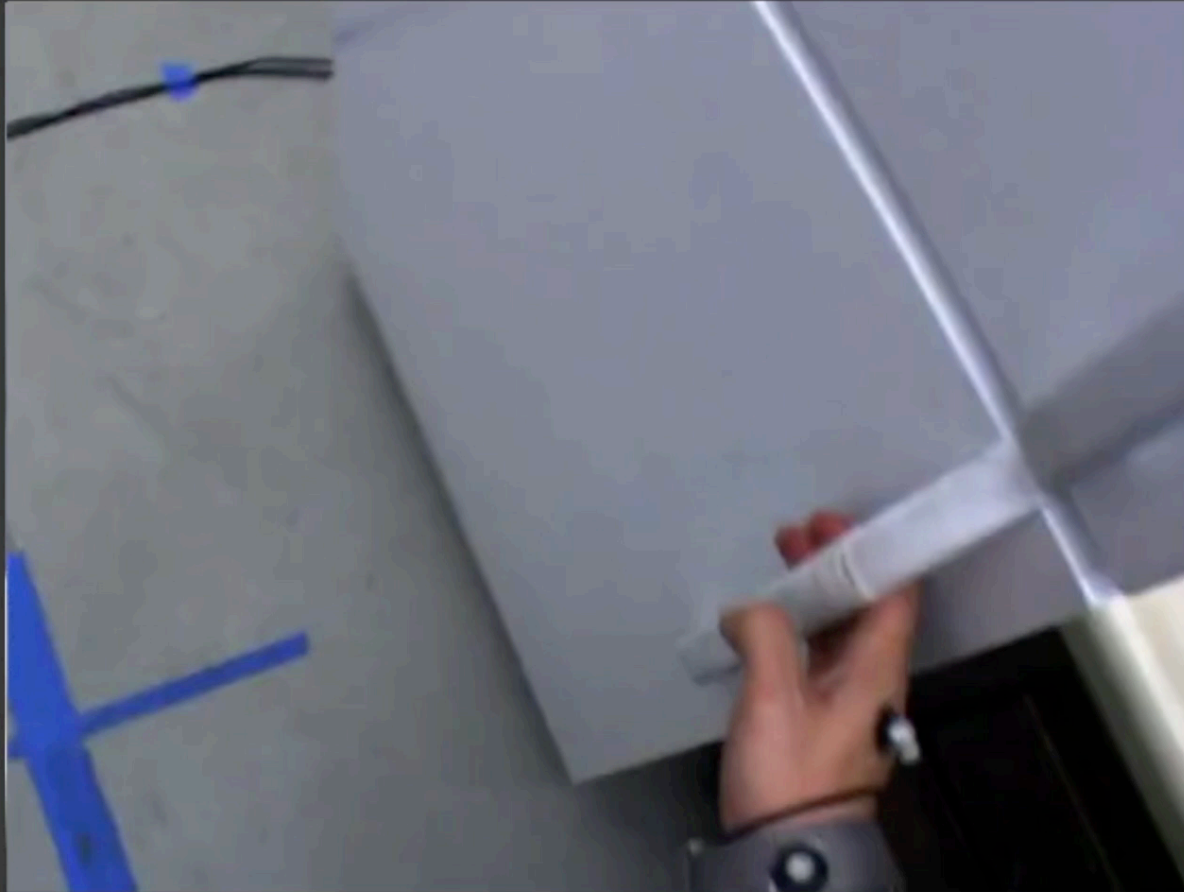


Object Interactions – the Dilemma



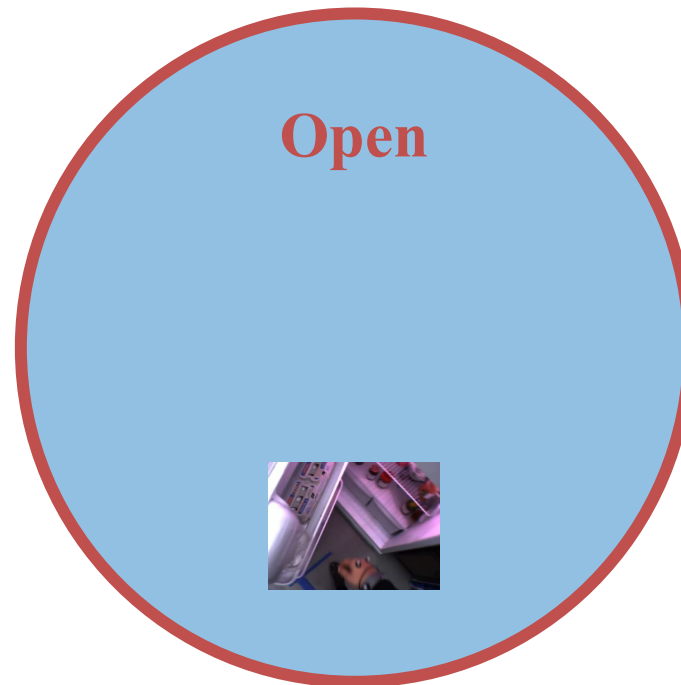
Object Interactions – the Dilemma

with: Michael Wray
Davide Moltisanti
Walterio Mayol-Cuevas



Object Interactions – the Dilemma

with: Michael Wray
Davide Moltisanti
Walterio Mayol-Cuevas



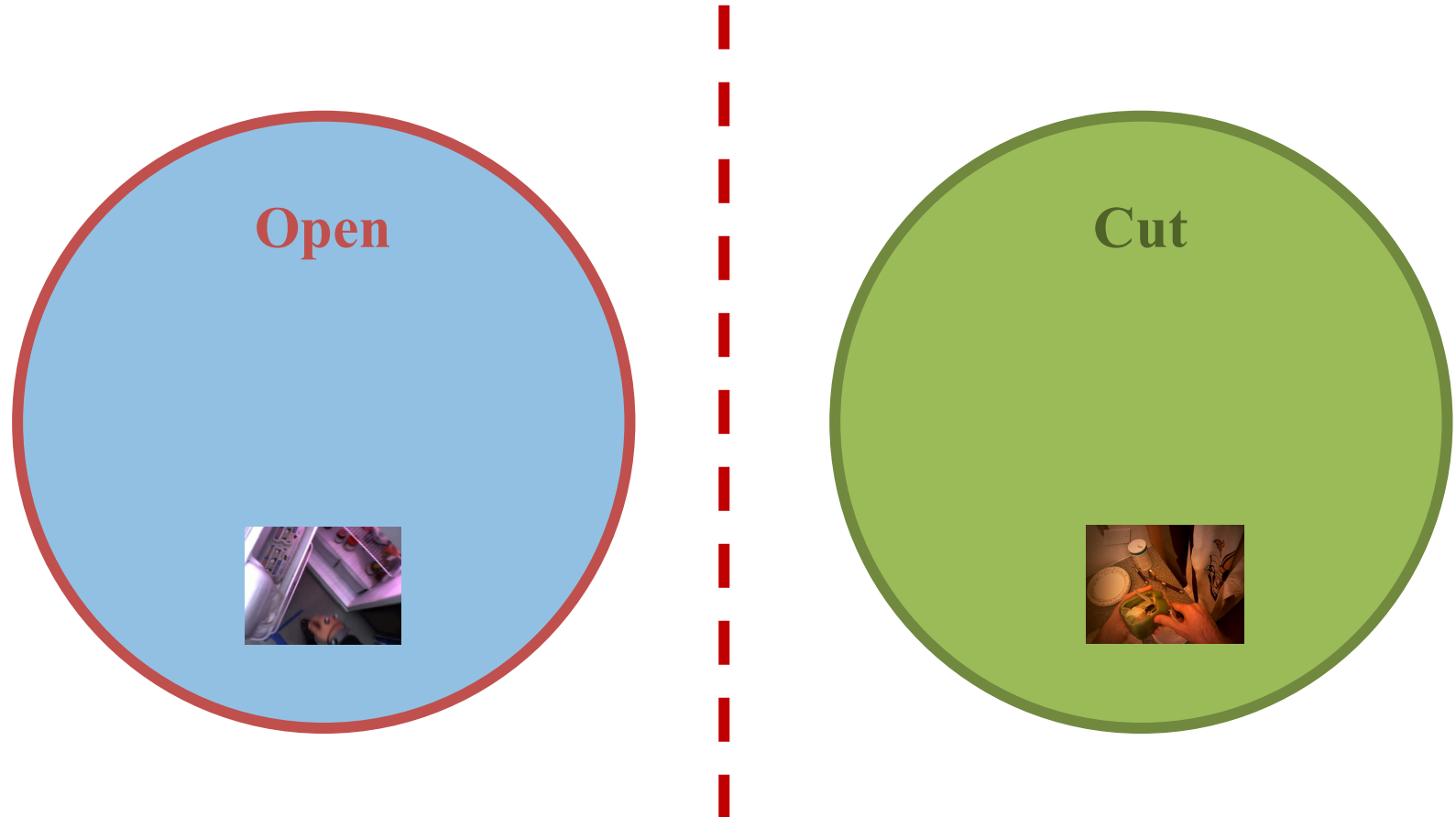
Object Interactions – the Dilemma

with: Michael Wray
Davide Moltisanti
Walterio Mayol-Cuevas



Object Interactions – the Dilemma

with: Michael Wray
Davide Moltisanti
Walterio Mayol-Cuevas



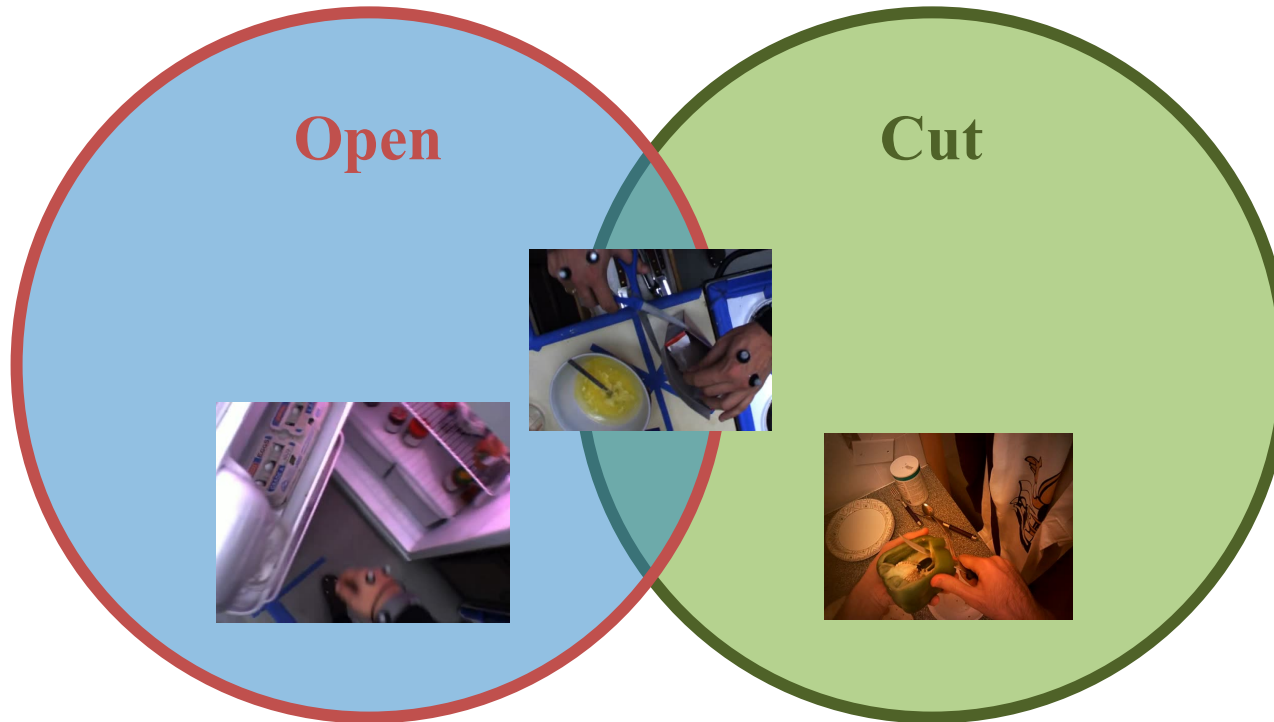
Object Interactions – the Dilemma

with: Michael Wray
Davide Moltisanti
Walterio Mayol-Cuevas



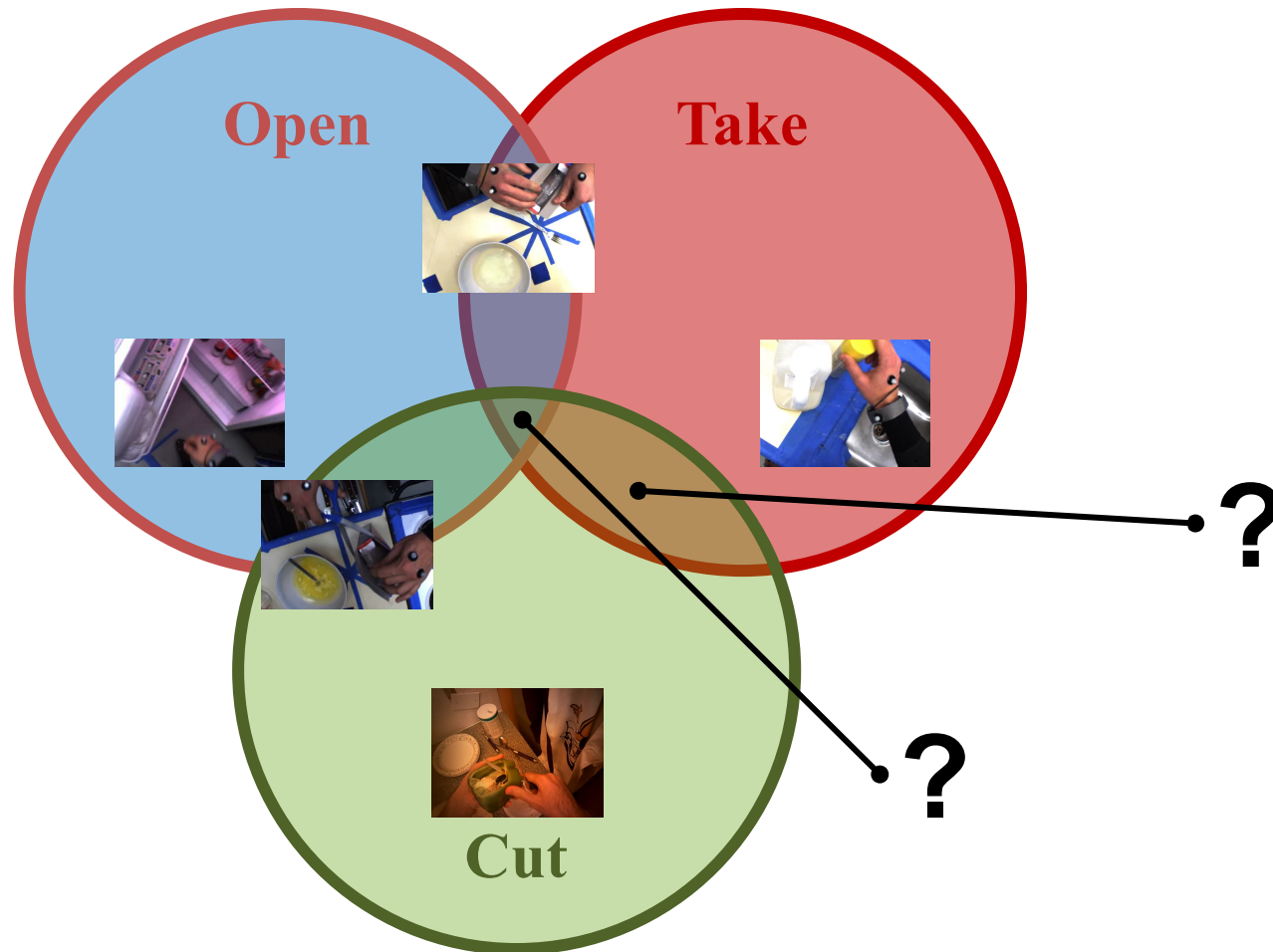
Object Interactions – the Dilemma

with: Michael Wray
Davide Moltisanti
Walterio Mayol-Cuevas



Object Interactions – the Dilemma

with: Michael Wray
Davide Moltisanti
Walterio Mayol-Cuevas



Object Interactions – the Dilemma

- Verbs cannot be separated into classes with hard boundaries.
- Rather the boundaries are more nuanced – what is correct in one video is incorrect for another.
- Singular classes are not enough.

Visualising Learnt Models

- BEOID EBP videos:

https://youtu.be/Fu7Db7Pau_A

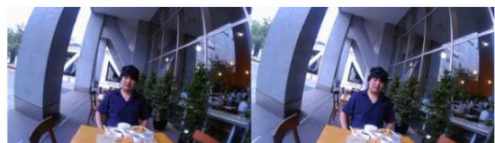
<https://youtu.be/4aDwQ-a3M68>

The Unique Problems

5. Multi-View Action Recognition

FPV with SPV

Input: paired egocentric videos



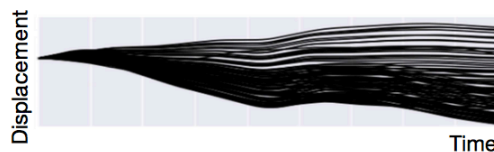
Egocentric video of person A



Egocentric video of person B



Multiple POV features of A



f_A : First-person POV feature of A

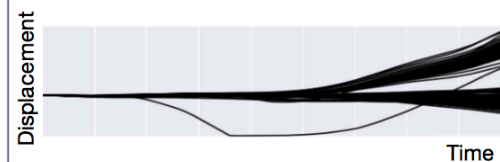


$f_{A \leftarrow B}$: Second-person POV feature of A

Multiple POV features of B



$f_{B \leftarrow A}$: Second-person POV feature of B

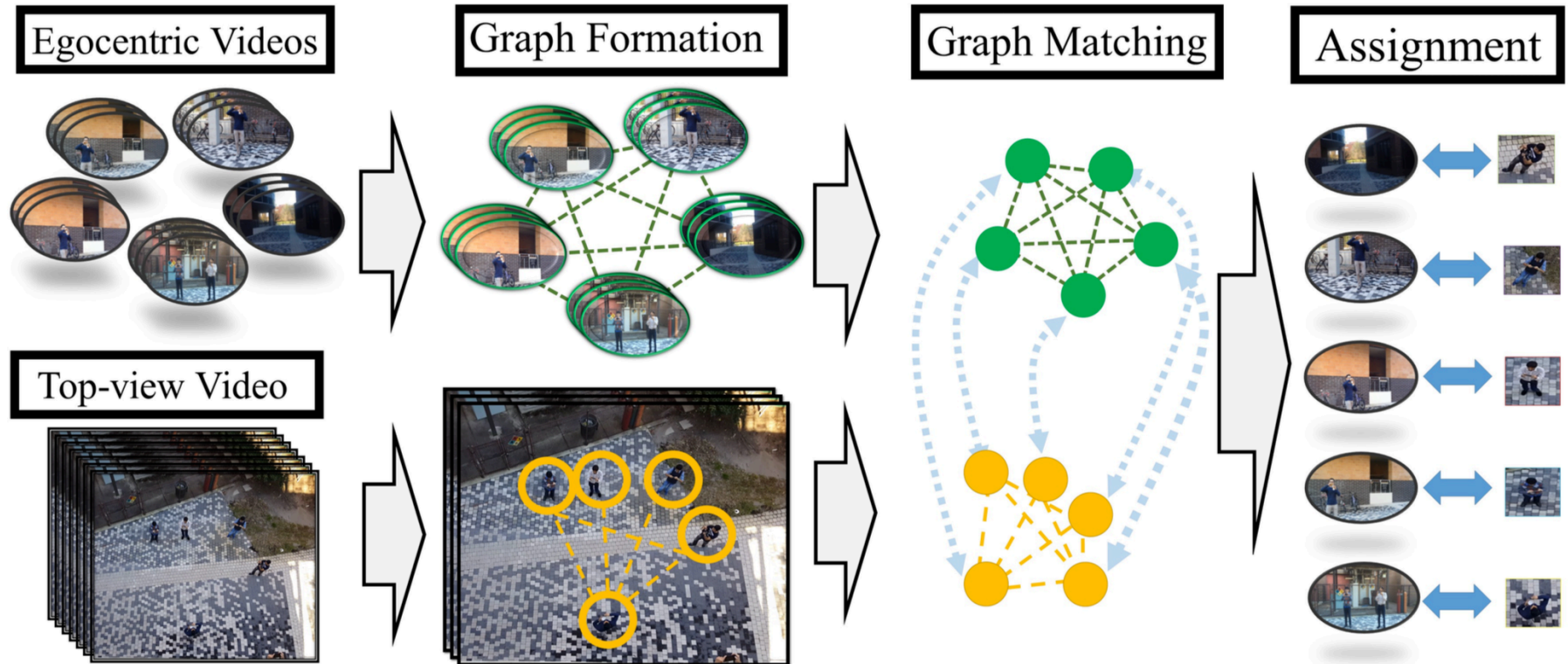


f_B : First-person POV feature of B

FPV with TPV (top-view)



FPV with TPV (top-view)



Egocentric Vision

- The Unique Problems
 1. Camera Motion
 2. Mapping and Localisation (ref tomorrow's talk)
 3. Attention and Task-Relevance
 4. Object Interactions
 5. Multi-view Solutions
- The Unique Applications
 1. Video Summarisation
 2. Skill Determination
 3. Real-time solutions

The Unique Applications

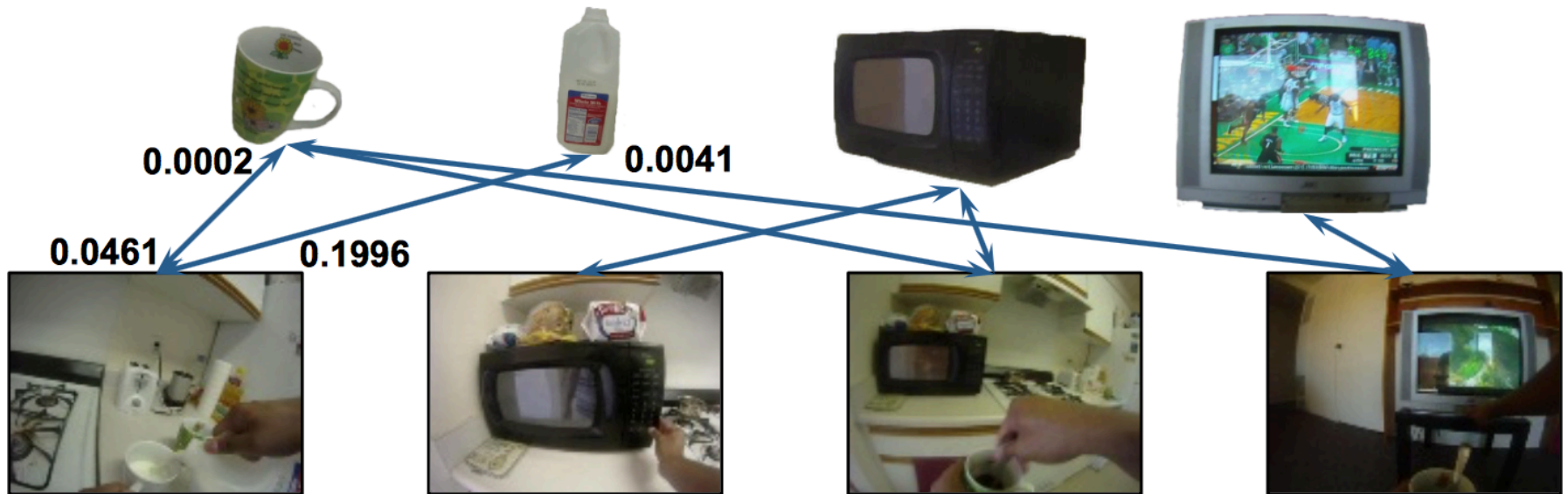
1. Video Summarisation

Video Summarisation

- Fixations
- Highlight Detection

Egocentric Video Summarisation

- Object-Driven



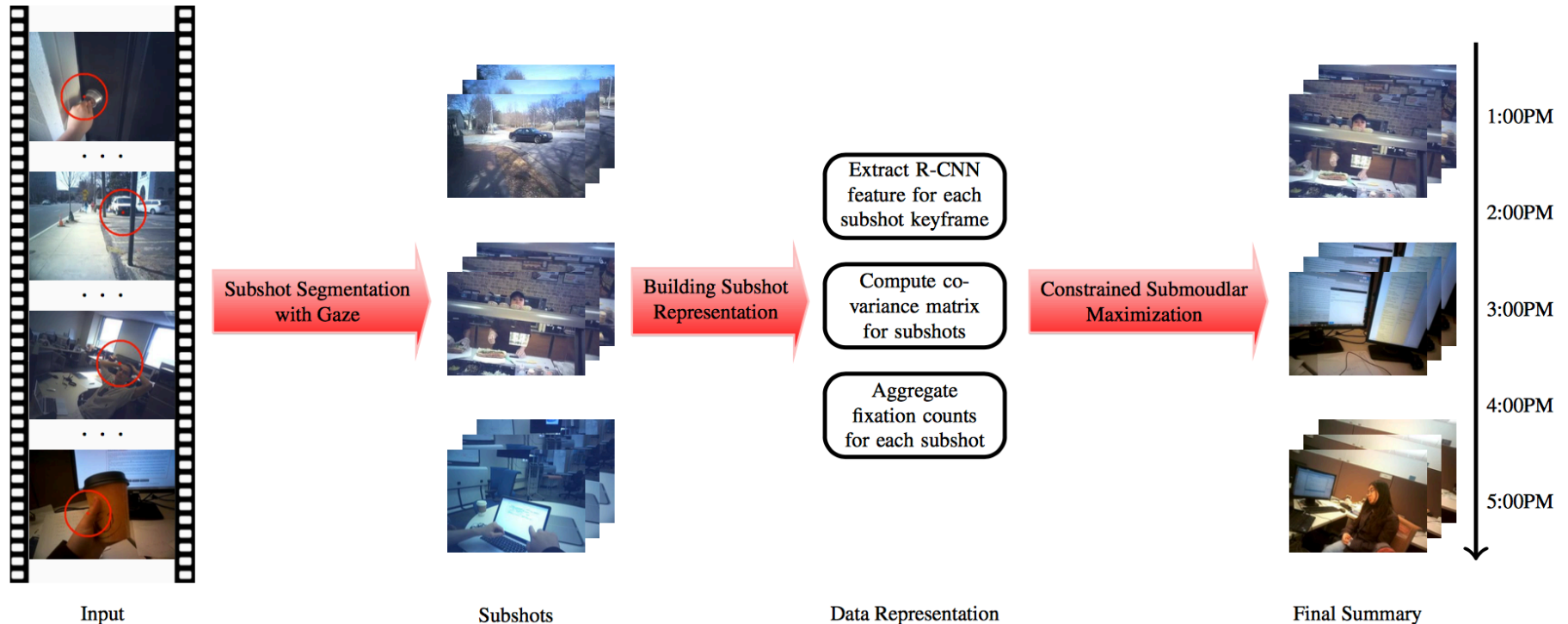
Egocentric Video Summarisation

- Object-Driven



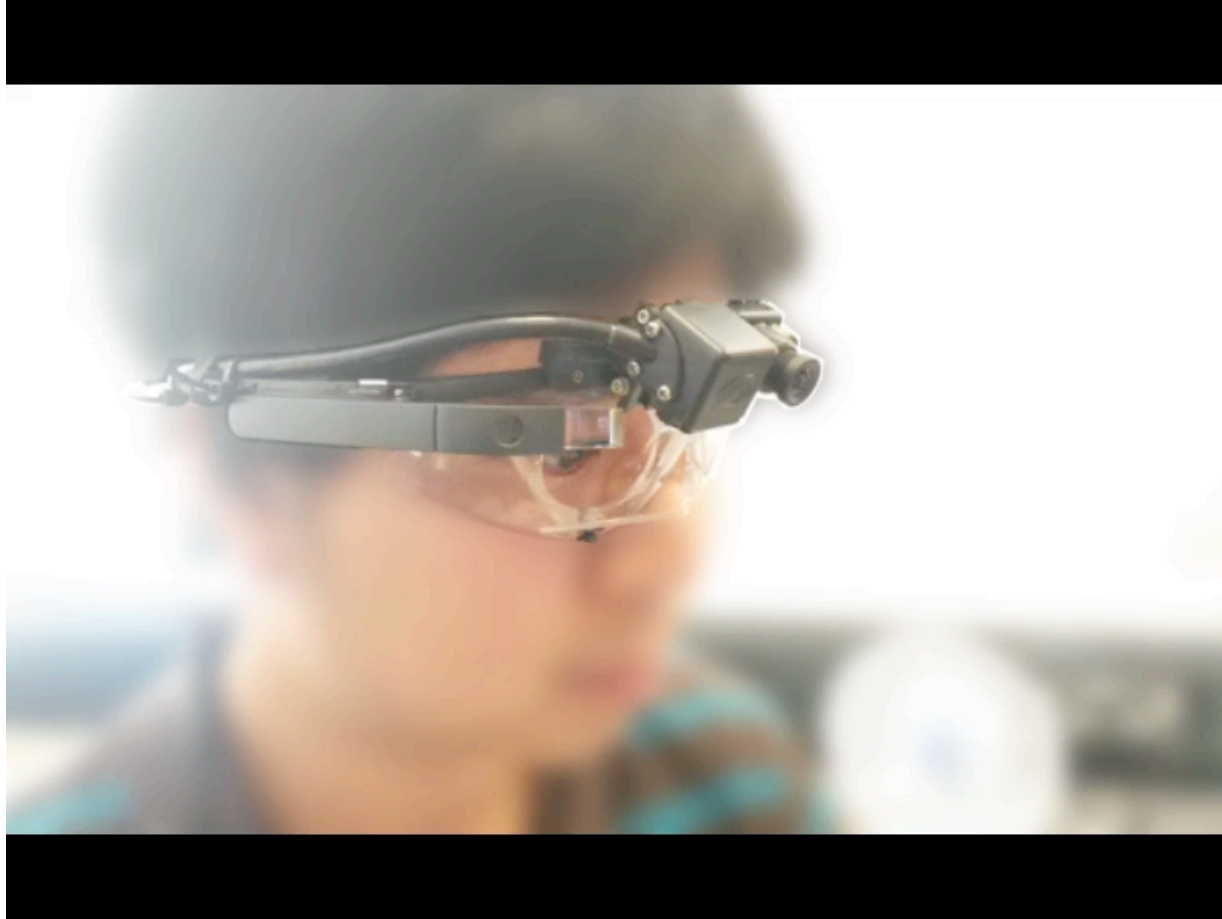
Egocentric Video Summarisation

- Fixation-Driven with Constraints



Egocentric Video Summarisation

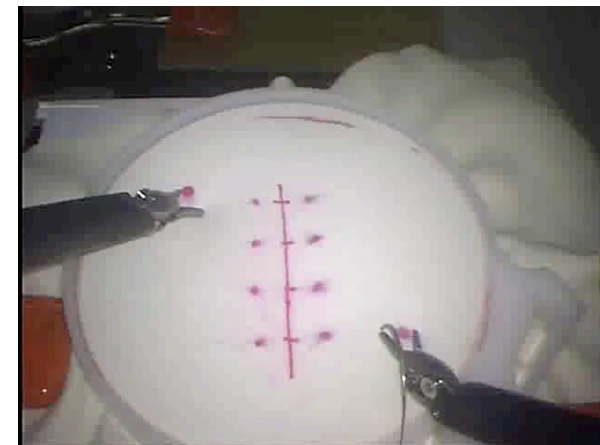
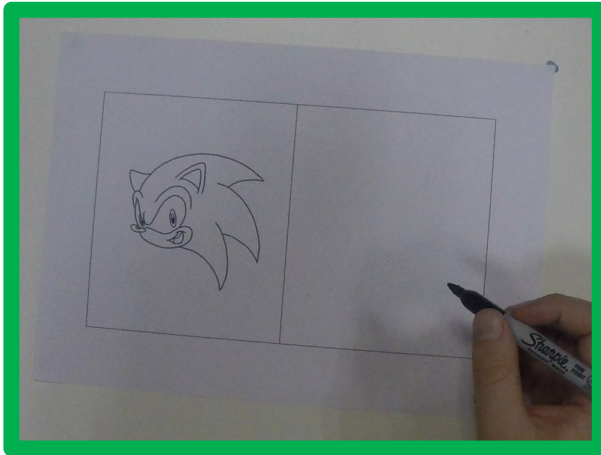
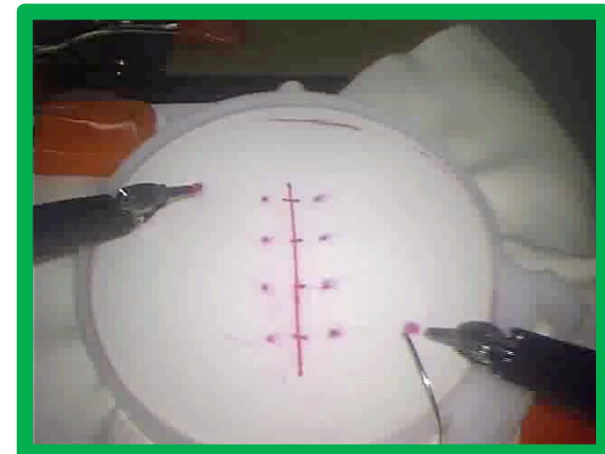
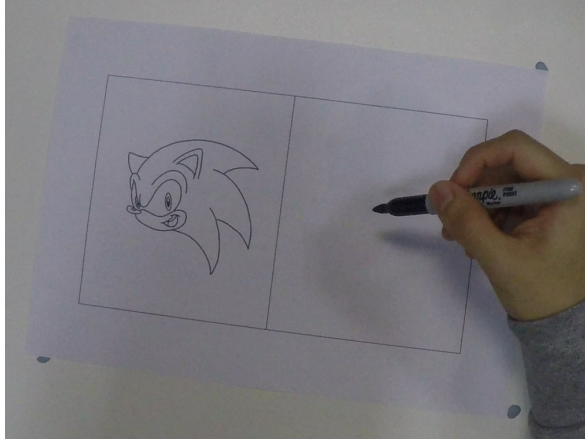
- Fixations from IMUs



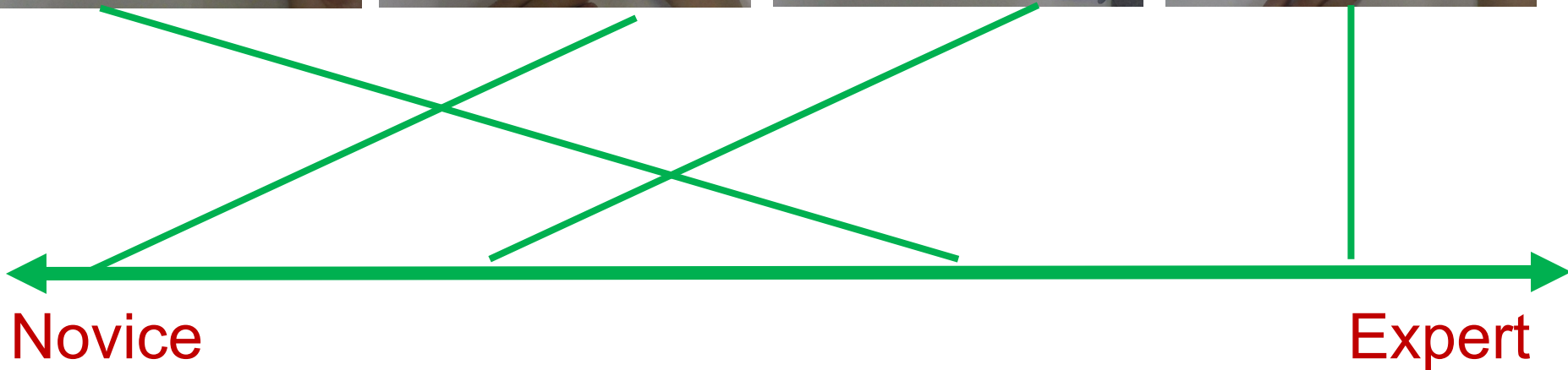
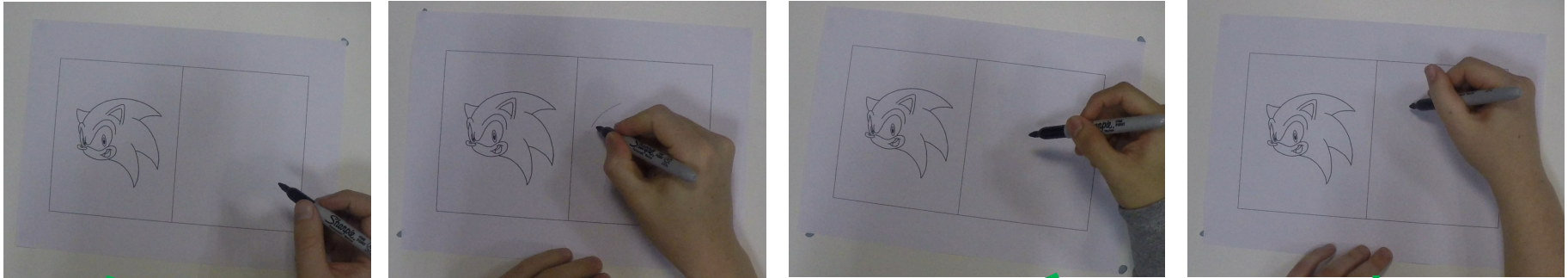
The Unique Applications

2. Skill Determination

Who's Better?



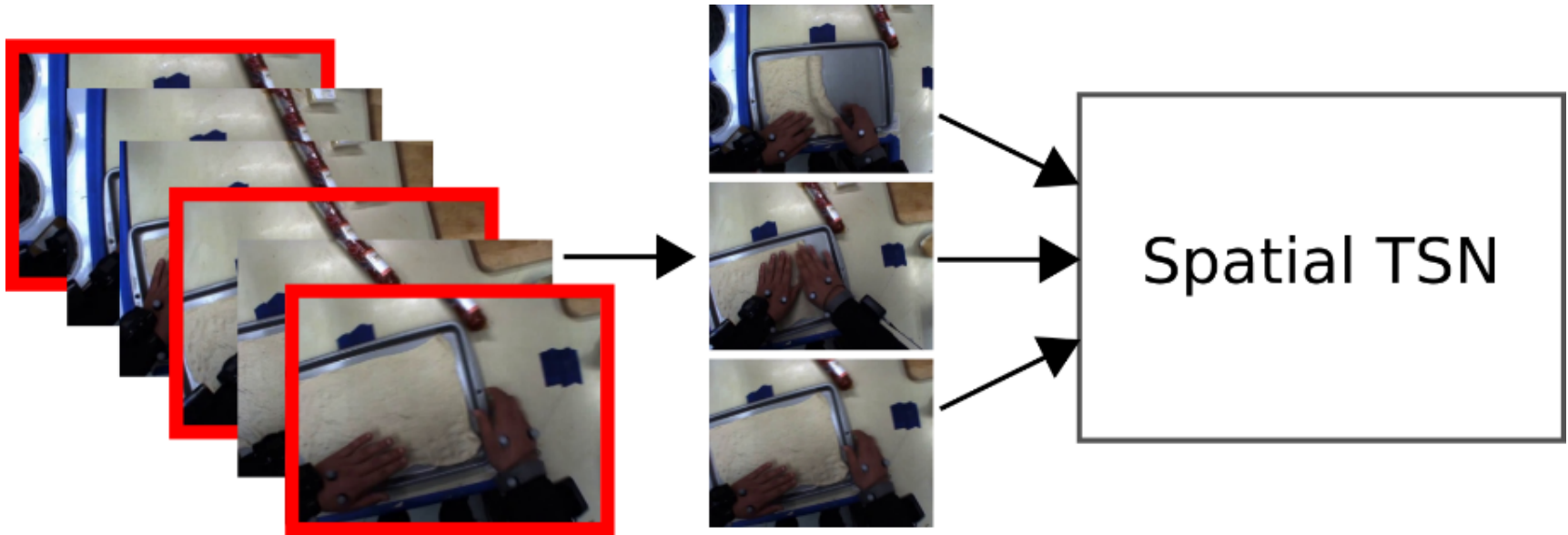
Who's Best?



Who's Better? Who's Best? Skill

with: Hazel Doughty
Walterio Mayol-Cuevas

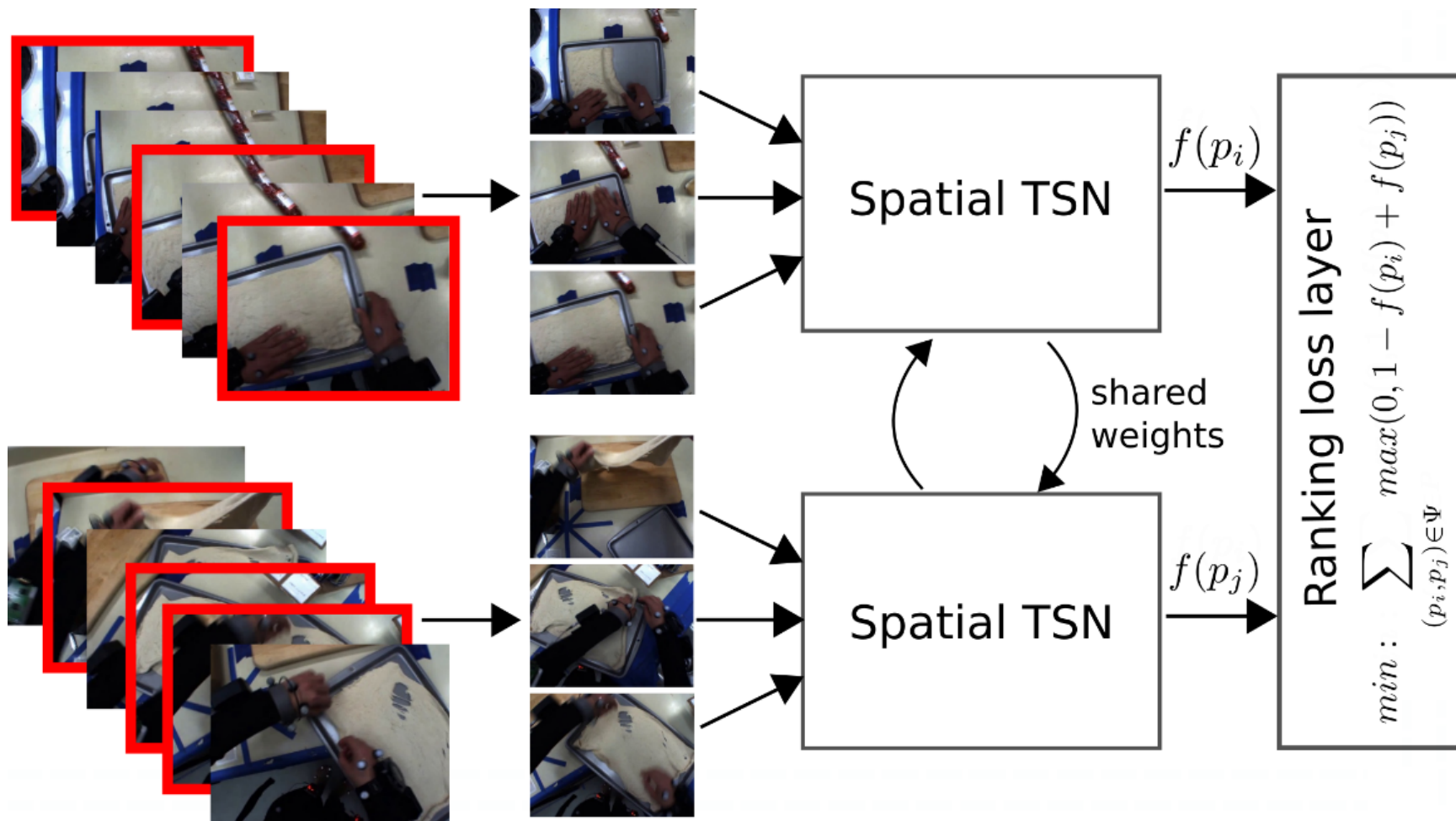
Determination in Video using Deep Ranking



Who's Better? Who's Best? Skill

with: Hazel Doughty
Walterio Mayol-Cuevas

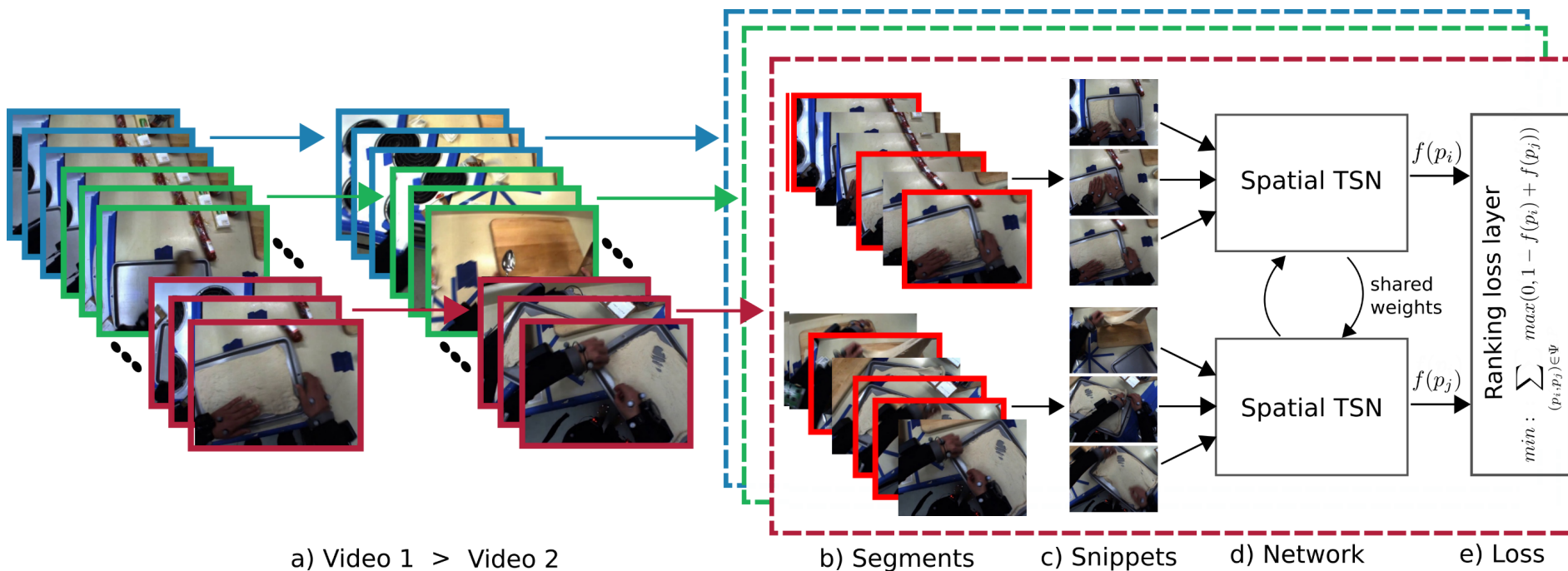
Determination in Video using Deep Ranking



Who's Better? Who's Best? Skill

with: Hazel Doughty
Walterio Mayol-Cuevas

Determination in Video using Deep Ranking



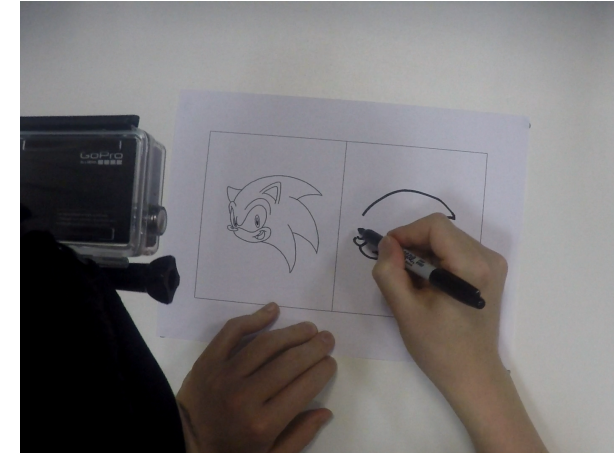
Who's Better? Who's Best? Skill Determination in Video using Deep Ranking

with: Hazel Doughty
Walterio Mayol-Cuevas

Surgery¹



Drawing



Dough-Rolling²



Chopstick-Using



Who's Better? Who's Best? Skill

with: Hazel Doughty
Walterio Mayol-Cuevas

Determination in Video using Deep Ranking

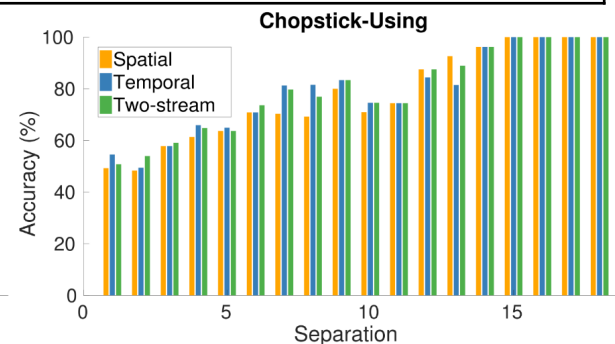
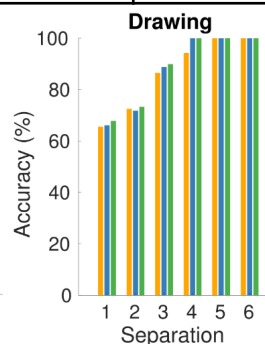
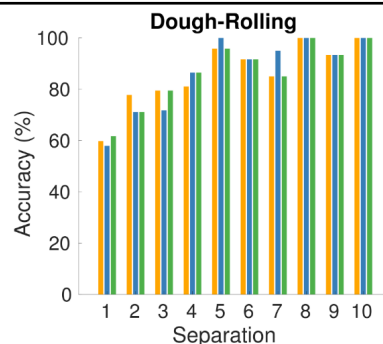
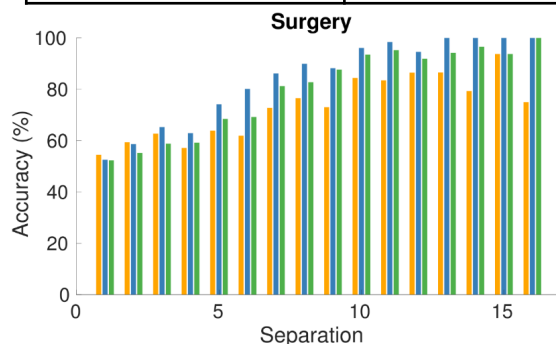
Task	Videos	Max Pairs	#Consistent Pairs	%Consistent Pairs
Surgery (Knot Tying)	36	630	596	95%
Surgery (Needle Passing)	28	378	362	96%
Surgery (Suturing)	39	741	701	95%
Dough-Rolling	33	528	181	34%
Drawing (Sonic)	20	190	118	62%
Drawing (Hand)	20	190	129	68%
Chopstick-Using	40	780	536	69%

Who's Better? Who's Best? Skill

with: Hazel Doughty
Walterio Mayol-Cuevas

Determination in Video using Deep Ranking

Task	Siamese TSN			Siamese TSN with data augmentation		
	Spatial	Temporal	Two-stream	Spatial	Temporal	Two-stream
Surgery	66.5%	74.4%	74.4%	66.5%	75.3%	75.3%
Dough-Rolling	73.9%	76.7%	75.4%	77.0%	76.1%	78.2%
Drawing	75.6%	76.5%	77.4%	76.7%	79.0%	82.1%
Chopstick-Using	67.7%	67.4%	68.1%	66.8%	69.8%	70.0%



Determination in Video using Deep Ranking

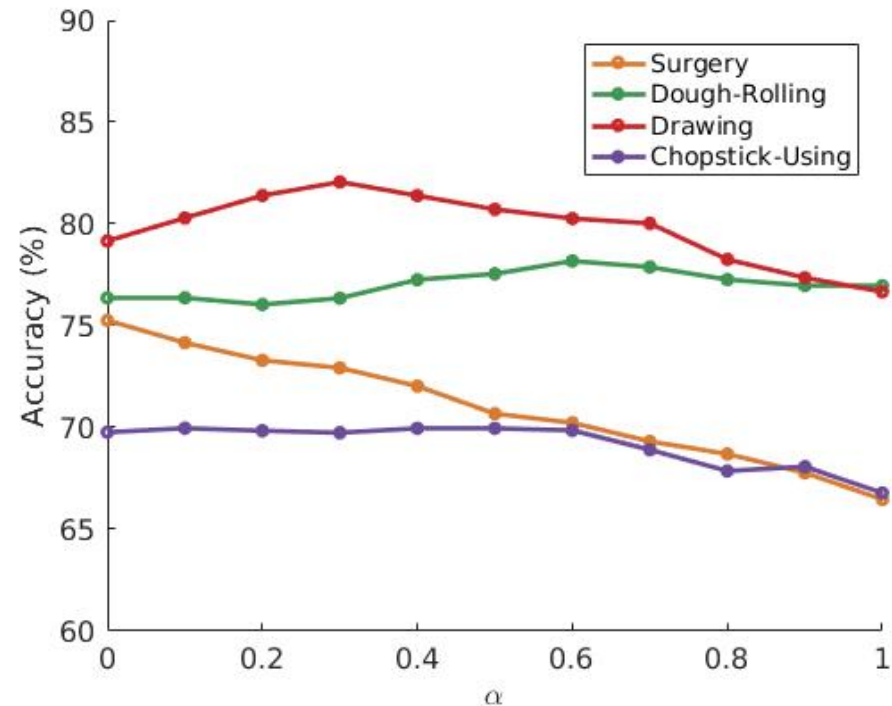
$$\frac{1}{\sigma} \sum_{j=1}^{\sigma} \alpha f_s(p_{ij}) + (1 - \alpha) f_t(p_{ij})$$

Who's Better? Who's Best? Skill

with: Hazel Doughty
Walterio Mayol-Cuevas

Determination in Video using Deep Ranking

$$\frac{1}{\sigma} \sum_{j=1}^{\sigma} \alpha f_s(p_{ij}) + (1 - \alpha) f_t(p_{ij})$$

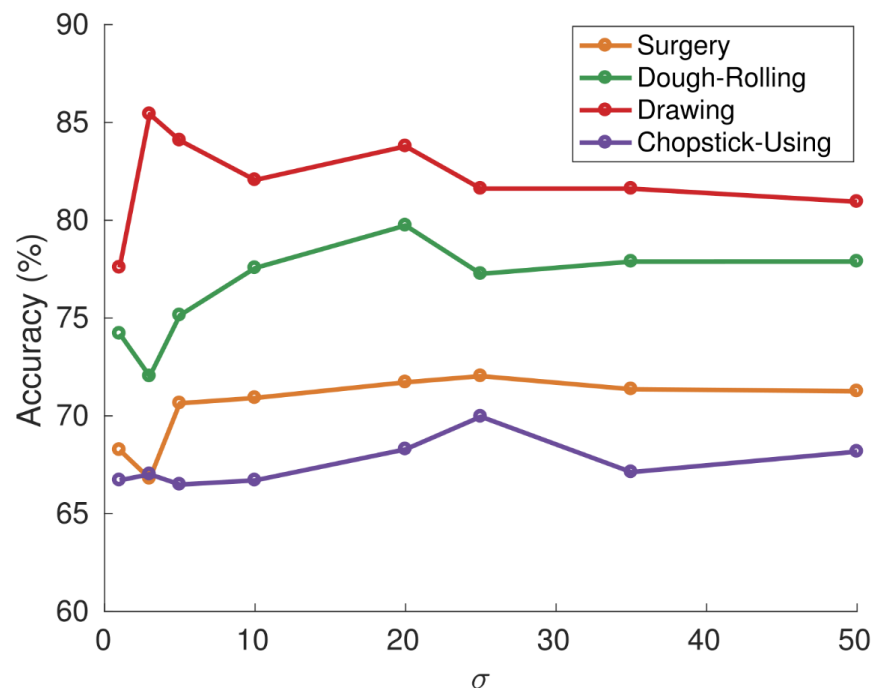


Who's Better? Who's Best? Skill

with: Hazel Doughty
Walterio Mayol-Cuevas

Determination in Video using Deep Ranking

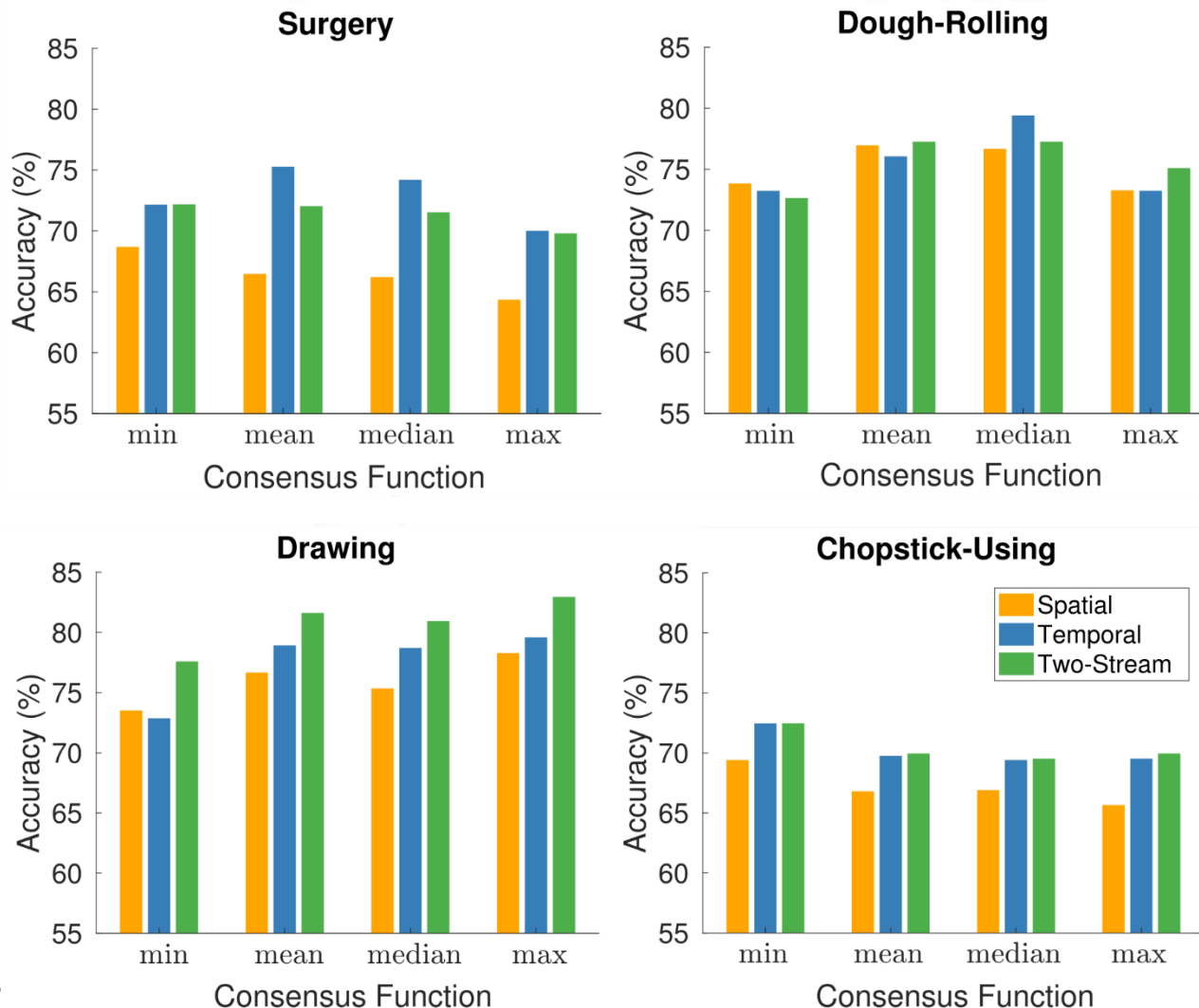
$$\frac{1}{\sigma} \sum_{j=1}^{\sigma} \alpha f_s(p_{ij}) + (1 - \alpha) f_t(p_{ij})$$



Who's Better? Who's Best? Skill

with: Hazel Doughty
Walterio Mayol-Cuevas

Determination in Video using Deep Ranking



Who's Better? Who's Best? Skill

with: Hazel Doughty
Walterio Mayol-Cuevas

Determination in Video using Deep Ranking

Example Rankings



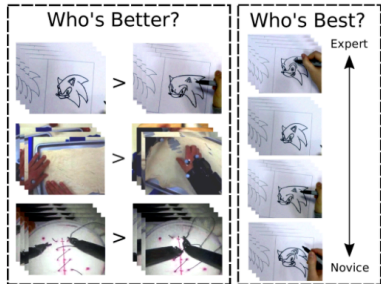
Lowest

Highest

Newly recorded Sonic-Drawing task

More info...

Project Who's Better, Who's Best: Skill Determination in Video



[Video](#)

Who's Better, Who's Best: Skill Determination in Video using Deep Ranking. H Doughty, D Damen, W Mayol-Cuevas. Arxiv (2017). [PDF](#)

The Unique Applications

3. Real-time Solutions

Wearable (Systems)!

- On-the-cloud processing
- On-the-mobile processing
- Onboard processing!

Connecting-to-the-cloud

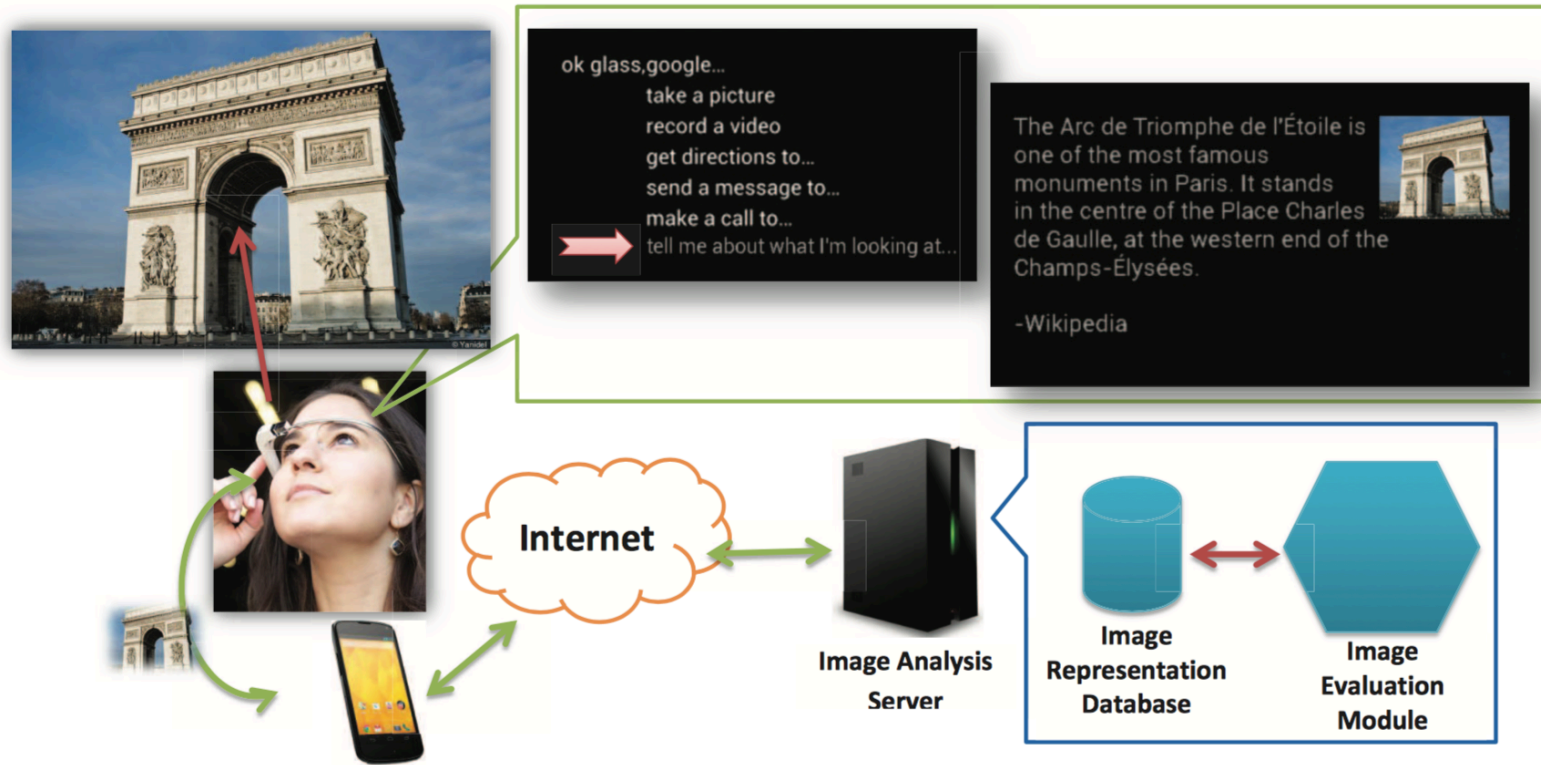


Figure 1. System overview. The user asks the device to inform her about her current view of Arc de Triomphe, and the system responds with the most relevant description in its database.

You Do, I Learn – Google Glass Prototype

GlaciAR
Final Demo

Teesid Leelasawassuk, Dima
Damen and Walterio Mayol
University of Bristol

October 2014

Interactive Conclusions

- Fill in the blanks:
 - Egocentric vision is -----
 - Pick up an action (e.g. open door). Draw a sketch of how it looks like from FPV and TPV
 - The biggest challenge (in your opinion) in egocentric vision is -----
 - The most interesting problem (to you) in egocentric vision is -----

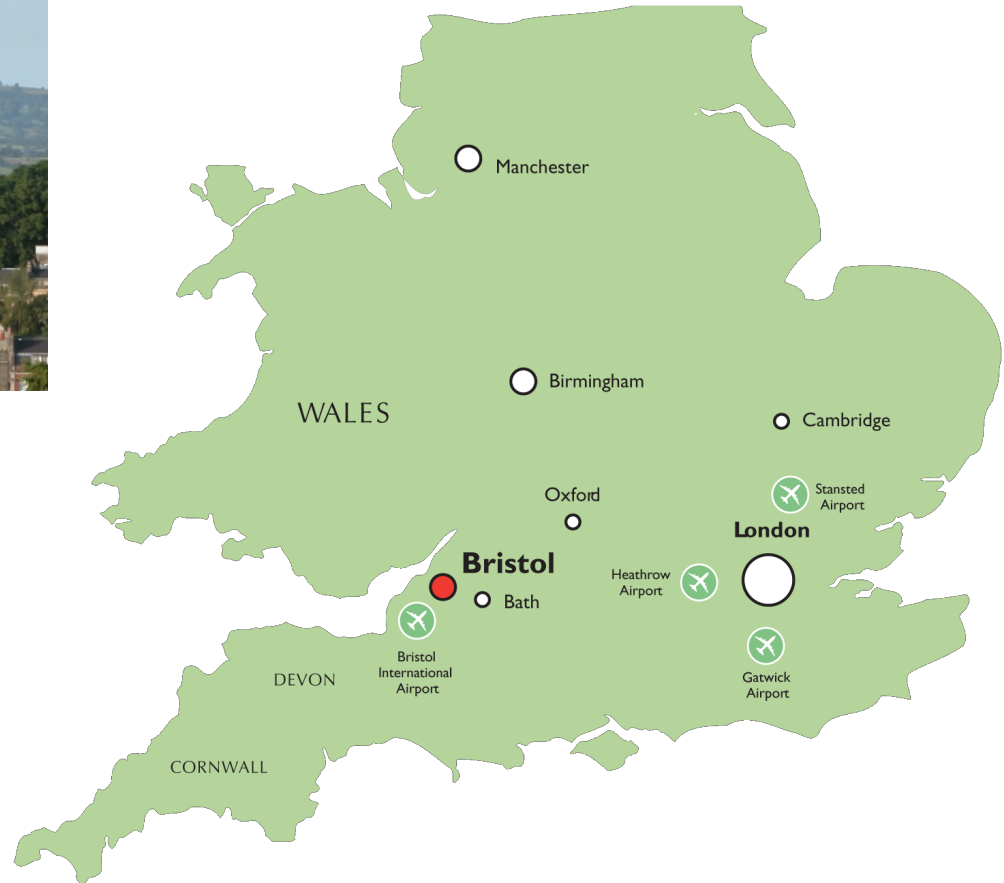
Interested in More?

- Egocentric Perception, Interaction and Computing (EPIC) Workshop Series
 - ECCV 2016 (Amsterdam)
 - ICCV 2017 (Venice – this October)
 - Paper deadline expired
 - Abstract submission still open till Sep

Interested in More?

- Subscribe to the newly introduced mailing list: epic-community@bristol.ac.uk
- Instructions to subscribe:
 - send an email to: sympa@sympa.bristol.ac.uk
 - with the subject: **subscribe epic-community**
 - and blank message content

Bristol and University of Bristol



Bristol and University of Bristol



Thank you...

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

<http://www.cs.bris.ac.uk/~damen>



@dimadamen



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