



Università
di Catania



First Person (Egocentric) Vision for Human-Centric Assistance: History, Building Blocks, and Applications

Antonino Furnari, Francesco Ragusa

Image Processing Laboratory - <http://iplab.dmi.unict.it/>

Department of Mathematics and Computer Science - University of Catania

Next Vision s.r.l., Italy

furnari@dmf.unict.it - <http://www.antoninofurnari.it/>

francesco.ragusa@unict.it - <https://iplab.dmi.unict.it/ragusa/>

<http://iplab.dmi.unict.it/fpv> - <https://www.nextvisionlab.it/>

Before we begin...

The slides of this tutorial are available online at:

<http://www.antoninofurnari.it/talks/iciap2022>



Agenda

- 1) Part I: Definitions, motivations, history and research trends [14.00 - 15.30] – Antonino Furnari
 - a) What is first person vision? What is it for?
 - b) What makes it different from third person vision?
 - c) History of First Person Vision: visions, ideas, research, devices;
 - d) Where do we go from here? Research trends, datasets and challenges.

Coffe Break [15.30 – 16.00]

- 1) **Part II: Building Blocks for First Person Vision Systems [16.00 – 18.00] – Francesco Ragusa**
 - a) **Data Acquisition & Datasets;**
 - b) **Fundamental Task in First Person Vision:**
 - i) **Localization;**
 - ii) **Object Detection and Recognition;**
 - iii) **Egocentric Human-Object Interaction;**
 - iv) **Action/Activities;**
 - v) **Anticipation.**
 - c) **Example Applications;**
 - d) **Conclusion.**

Part 2

Building Blocks for First Person Vision Systems

Data Acquisition – Video Quality

- Try to get a high quality camera to get high quality images!
- Egocentric video is subject to motion blur and exposure issues.

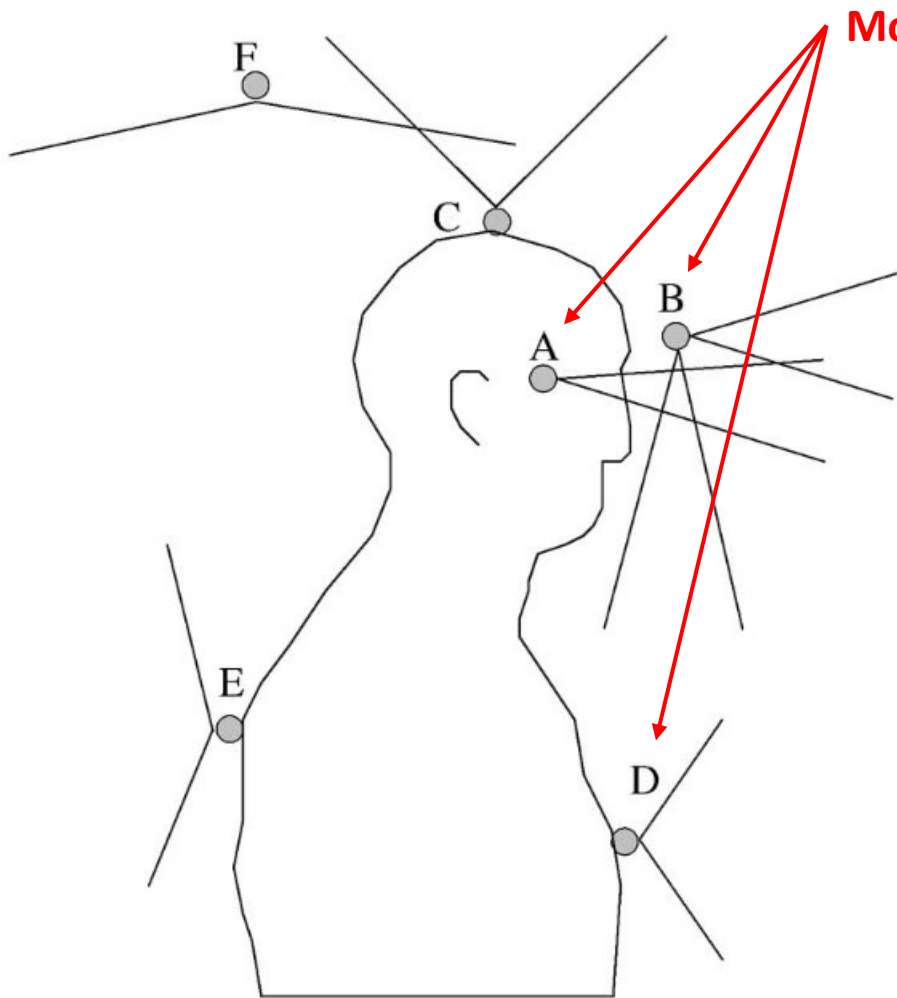
High Quality Video Obtained with a GoPro



Average Quality Video



Data Acquisition – Camera Wearing Modalities



Most Common Wearing Modalities **A,B: head mounted, D: chest mounted**

A



B (frontward)



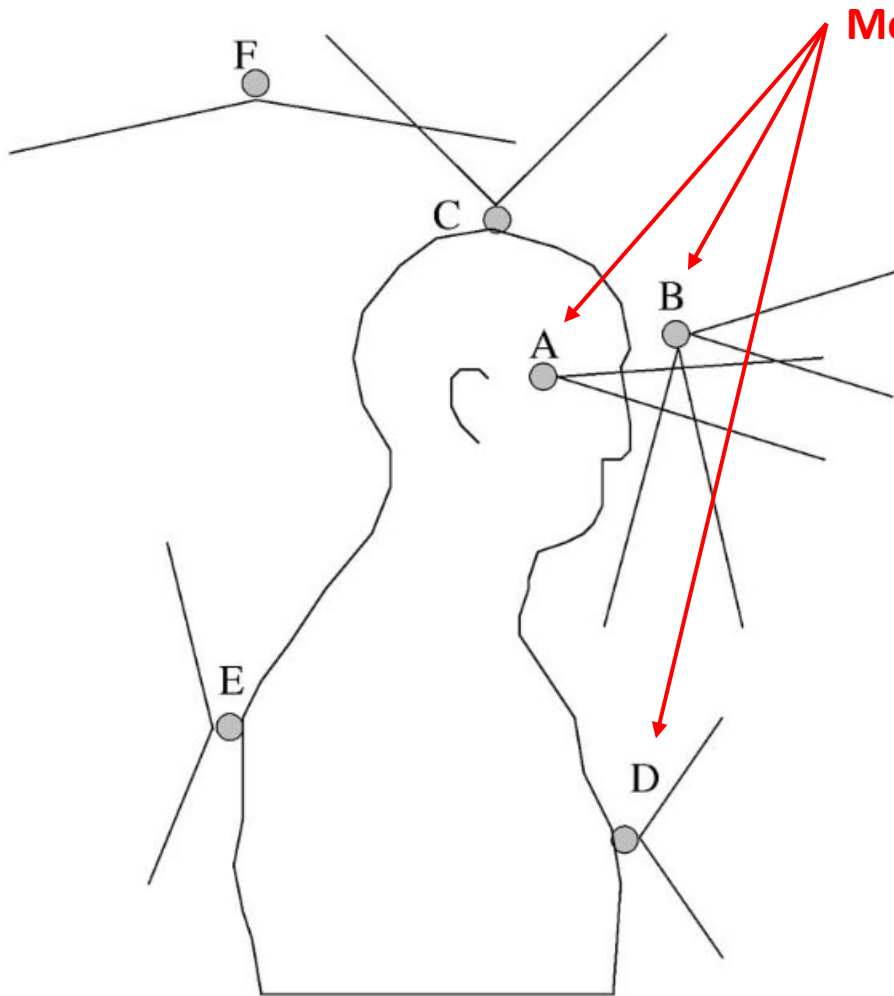
B (downward)



D



Data Acquisition – Camera Wearing Modalities (2)



Most Common Wearing Modalities

- A-B are best to capture objects:
 - A, B (frontward) to capture objects in front of the subjects (e.g., paintings in a museum);
 - B (downward) to capture objects manipulated with hands (e.g., kitchen);
- Chest-mounted cameras (D) are less obtrusive and give stable video, but they may miss details on what the user is looking at;

Data Acquisition – Field of View (FOV)

A wide FOV allows to capture more scene but introduces distortion.

Narrow Angle



Wide Angle



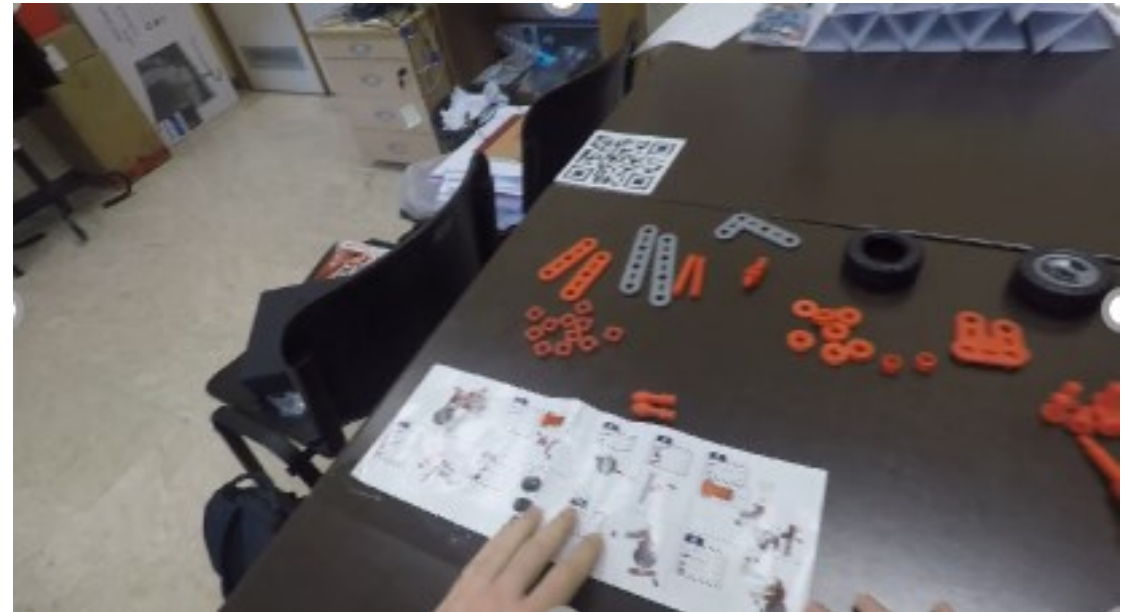
Data Acquisition – Field of View (FOV)

A wide FOV allows to capture more scene but introduces distortion.

Narrow Angle



Wide Angle



Data Acquisition – Other Modalities – Depth

- If you can acquire depth, do it!
- Depth can improve scene understanding by highlighting the position of objects and hands;



Data Acquisition – Other Modalities – Gaze

Gaze can give information on what the user is paying attention to.

However, gaze trackers generally require a calibration process (and some expertise).



F. Ragusa, A. Furnari, S. Livatino, G. M. Farinella. The MECCANO Dataset: Understanding Human-Object Interactions from Egocentric Videos in an Industrial-like Domain. WACV 2021 (ORAL) (<https://arxiv.org/abs/2010.05654>).

Datasets

- If you are trying to solve a specific FPV problem, chances are that someone already collected/labeled data that is suitable for you.
- Search on the internet first!
- In particular, there are quite a few dataset focusing on action/activity recognition;
- In the following, a (non-exhaustive) list of datasets.

Datasets (non-exhaustive)

Dataset	URL	Settings	Annotations	Goal
EGO4D	https://ego4d-data.org/	931 participants performing different activities in different domains.	Different temporal and spatial annotations related to 5 benchmarks	Episodic Memory, Hand-Object Interaction, Audio-Visual Diarization, Social Interactions, Forecasting
EPIC-KITCHENS-100	https://epic-kitchens.github.io/2020-100	Subjects performing unscripted actions in their native kitchens.	Temporal segments	Action recognition, detection, anticipation, retrieval.
MECCANO	https://iplab.dmi.unict.it/MECCANO/	20 subjects assembling a toy motorbike.	Temporal segments, active objects, human-object interactions	Action recognition, Active object detection, Egocentric Human-Object Interaction Detection
ASSEMBLY101	https://assembly-101.github.io/	53 subjects assembling in a cage settings 101 children's toys.	Temporal segments, 3D hand poses	Action recognition, Action Anticipation, Temporal Segmentation

Datasets (non-exhaustive)

Dataset	URL	Settings	Annotations	Goal
EPIC-KITCHENS 2018	https://epic-kitchens.github.io/2018	32 subjects performing unscripted actions in their native environments	action segments, object annotations	Action recognition, Action Anticipation, Object Detection
Charade-Ego	https://allenai.org/plato/charades/	paired first-third person videos	action classes	Action recognition
EGTEA Gaze+	http://ai.stanford.edu/~alireza/GTEA/	32 subjects, 86 sessions, 28 hours	action segments, gaze, hand masks	Understanding daily activities, action recognition
ADL	https://www.csee.umbc.edu/~hpirsiav/papers/ADLdataset/	20 subjects performing daily activities in their native environments	activity segments, objects	Detecting activities of daily living
CMU kitchen	http://www.cs.cmu.edu/~espriggs/cmu-mmacc/annotations/	multimodal, 18 subjects cooking 5 different recipes: brownies, eggs, pizza, salad, sandwich	action segments	Understanding daily activities
EgoSeg	http://www.vision.huji.ac.il/egoseg/	Long term actions (walking, running, driving, etc.)	long term activity	Temporal Segmentation, Indexing

Datasets (non-exhaustive)

Dataset	URL	Settings	Annotations	Goal
First-Person Social Interactions	http://ai.stanford.edu/~alireza/Disney/	8 subjects at disneyworld	Activities: walking, waiting, gathering, sitting, buying something, eating, etc.	Recognizing social interactions
UEC Dataset	http://www.cs.cmu.edu/~kkitani/datasets/	two choreographed datasets with different egoactions (walk, jump, climb, etc.) + 6 youtube sports videos	activities	Unsupervised activity recognition
JPL	http://michaelryoo.com/jpl-interaction.html	interaction with a robot	activities performed on the robot + pose	Interaction recognition/prediction
Multimodal Egocentric Activity Dataset	http://people.sutd.edu.sg/~1000892/dataset	15 seconds clips of 20 activities	activity (walking, elevator, etc.)	Life-logging
LENA: An egocentric video database of visual lifelog	http://people.sutd.edu.sg/~1000892/dataset	13 activities performed by 10 subjects (Google Glass)	activity (walking, elevator, etc.)	Life-logging

Datasets (non-exhaustive)

Dataset	URL	Settings	Annotations	Goal
FPPA	http://tamaraberg.com/prediction/Prediction.html	Five subjects performing 5 daily actions	activity (drinking water, putting on clothes, etc.)	Temporal prediction
UT Egocentric	http://vision.cs.utexas.edu/projects/egocentric/index.html	3-5 hours long videos capturing a person's day	important regions	Summarization
VINST/ Visual Diaries	http://www.csc.kth.se/cvap/vinst/NovEgoMotion.html	31 videos capturing the visual experience of a subject walkin from metro station to work	location id, novel egomotion	Novelty detection
Bristol Egocentric Object Interaction (BEOID)	https://www.cs.bris.ac.uk/~damen/BEOID/	8 subjects, six locations. Interaction with objects and environment	gaze, objects, mode of interaction (pick, plug, etc.)	Provide assistance on object usage
Object Search Dataset	https://github.com/Mengmi/deepfuturegaze_gan	57 sequences of 55 subjects on search and retrieval tasks	gaze	gaze prediction

Datasets (non-exhaustive)

Dataset	URL	Settings	Annotations	Goal
UNICT-VEDI	http://iplab.dmi.unict.it/VEDI/	different subjects visiting a museum	location, observed objects	localizing visitors of a museum and estimating their attention
UNICT-VEDI-POI	http://iplab.dmi.unict.it/VEDI_POIs/	different subjects visiting a museum	object bounding boxes annotations, observed objects	recognizing points of interest observed by the visitors
Simulated Egocentric Navigations	http://iplab.dmi.unict.it/SimulatedEgocentricNavigations/	simulated navigations of a virtual agent within a large building	3-DOF pose of the agent in each image	egocentric localization
EgoCart	http://iplab.dmi.unict.it/EgocentricShoppingCartLocalization/	egocentric images collected by a shopping cart in a retail store	3-DOF pose of the shopping cart in each image	egocentric localization
Unsupervised Segmentation of Daily Living Activities	http://iplab.dmi.unict.it/dailylivingactivities	egocentric videos of daily activities	activities	unsupervised segmentation with respect to the activities

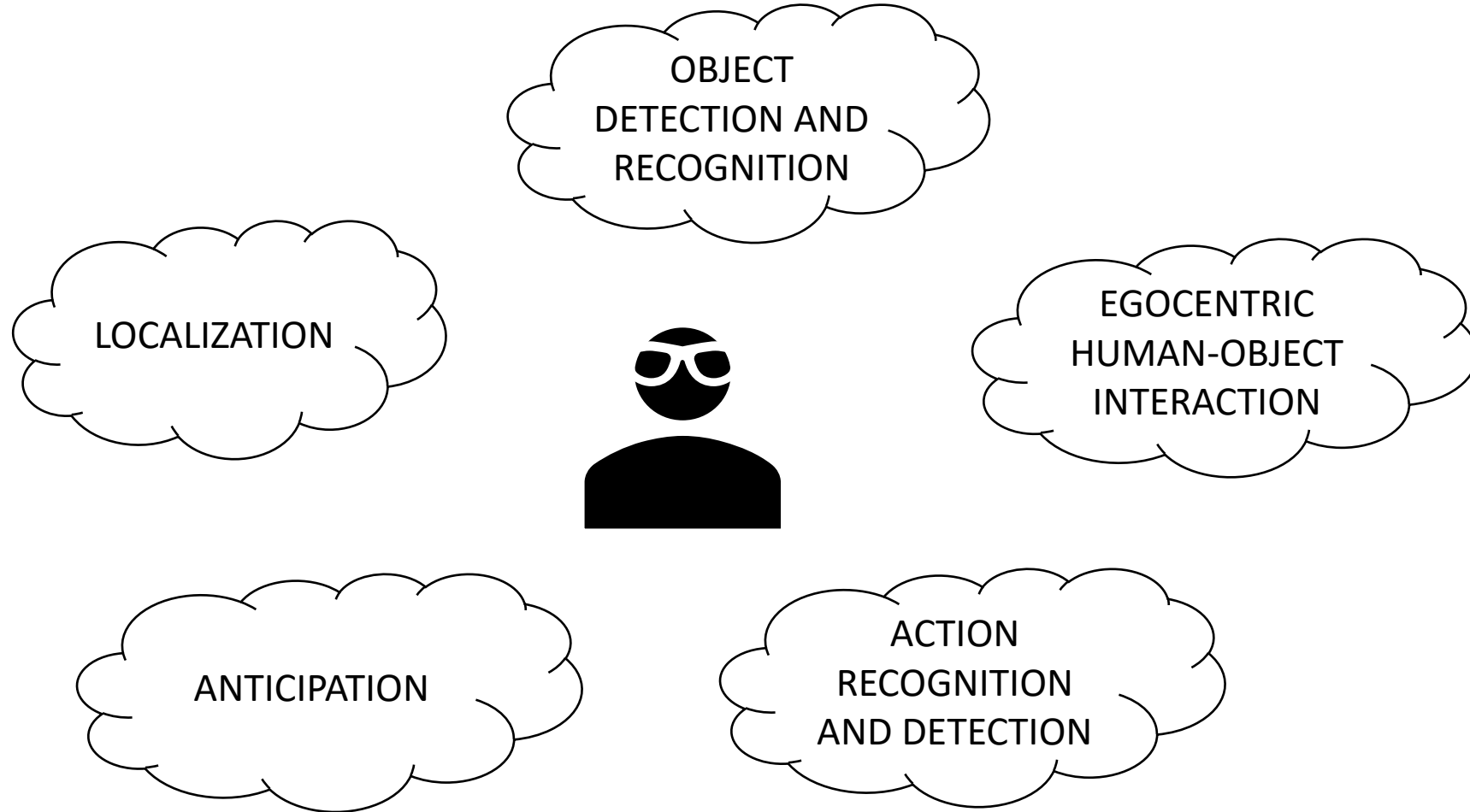
Datasets (non-exhaustive)

Dataset	URL	Settings	Annotations	Goal
Visual Market Basket Analysis	http://iplab.dmi.unict.it/vmba/	egocentric images collected by a shopping cart in a retail store	class-location of each image	egocentric localization
Location Based Segmentation of Egocentric Videos	http://iplab.dmi.unict.it/PersonalLocationSegmentation/	egocentric videos of daily activities	location classes	egocentric localization, video indexing
Recognition of Personal Locations from Egocentric Videos	http://iplab.dmi.unict.it/PersonalLocations/	egocentric videos clips of daily activities	location classes	recognizing personal locations
EgoGesture	http://www.nlpr.ia.ac.cn/iva/yfzhang/datasets/egogesture.html	2k videos from 50 subjects performing 83 gestures	Gesture labels, depth	Gesture recognition
EgoHands	http://vision.soic.indiana.edu/projects/egohands/	48 videos of interactions between two people	Hand segmentation masks	Egocentric hand segmentation
DoMSEV	http://www.verlab.dcc.ufmg.br/sema-ntic-hyperlapse/cvpr2018-dataset/	80 hours/different activities	Scene/Action labels with IMU, GPS mad depth	Summarization

Datasets (non-exhaustive)

Dataset	URL	Settings	Annotations	Goal
EGO-HPE	http://imagelab.ing.unimore.it/imagelab2015/researchactivity.asp?idAttivita=23	Egocentric videos for head pose estimation	Head pose of the subjects	Head-pose estimation
EGO-GROUP	http://imagelab.ing.unimore.it/imagelab2015/researchactivity.asp?idAttivita=23	18 videos of people engaging social relationships	Social relationships	Understanding social relationships
DR(eye)VE	http://aimagelab.ing.unimore.it/dreyeve	74 videos of people driving	Eye fixations	Autonomous and assisted driving
THU-READ	http://ivg.au.tsinghua.edu.cn/dataset/THU_READ.php	8 subjects performing 40 actions with a head-mounted RGBD camera	Action segments	RGBD egocentric action recognition
EGO-CH	https://iplab.dmi.unict.it/EGO-CH/	70 subjects visiting two cultural sites in Sicily, Italy.	Temporal segments, room-based localization, objects	Room-based localization, Object detection, Behavioral analysis

Fundamental Tasks of a First Person Vision System



Localization in First Person Vision

- Knowing the location of the user for a First Person Vision system is important to implement contextual awareness
 - Behave differently depending on the environment
 - Generate reminders when I get to a particular place
 - «remember to do the laundry when you get home»;
 - Turn notifications on or off when you are in given environments:
 - Put in silent mode when I am in a conference room;
 - Help localize/navigate the user
 - E.g., in a retail store or in a museum;
 - Implement augmented reality
 - Show location-specific information when I get to a place (e.g., a room in a museum)

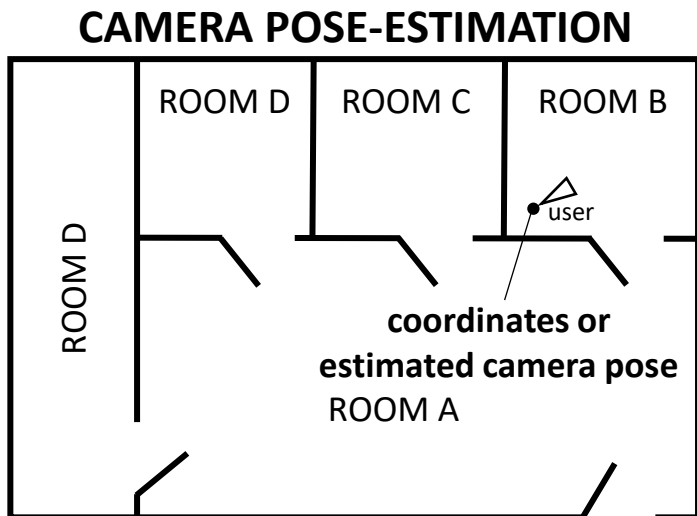
Localization – Levels of Granularity

SCENE RECOGNITION

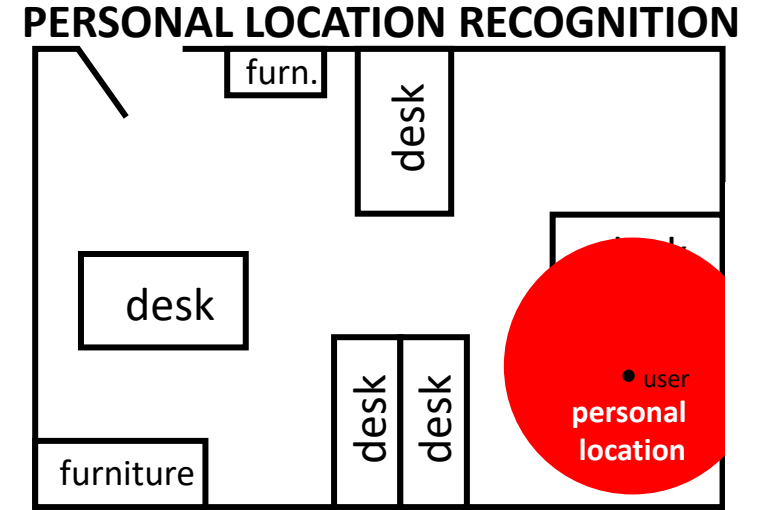
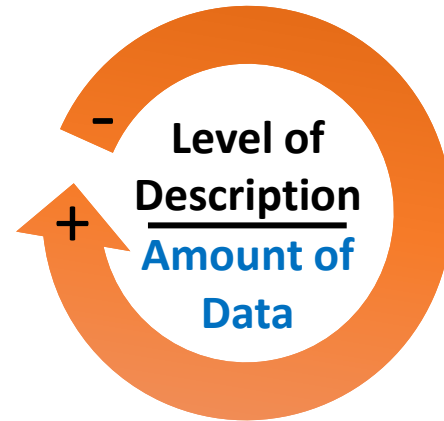


INSIDE CITY

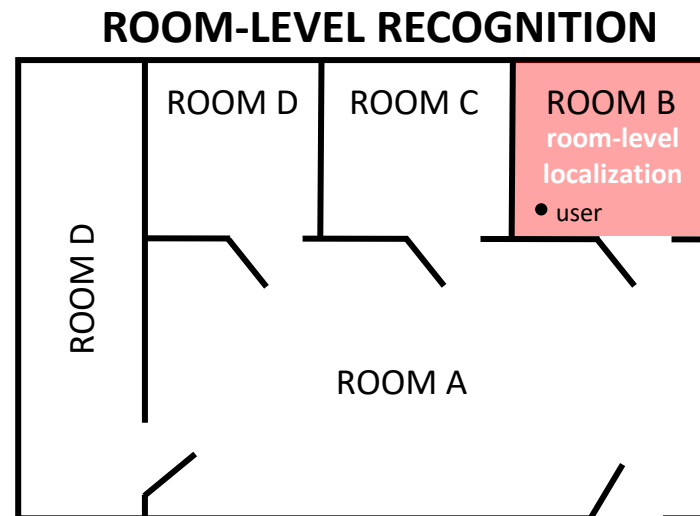
off-the-shelf detectors



3D reconstruction of the building



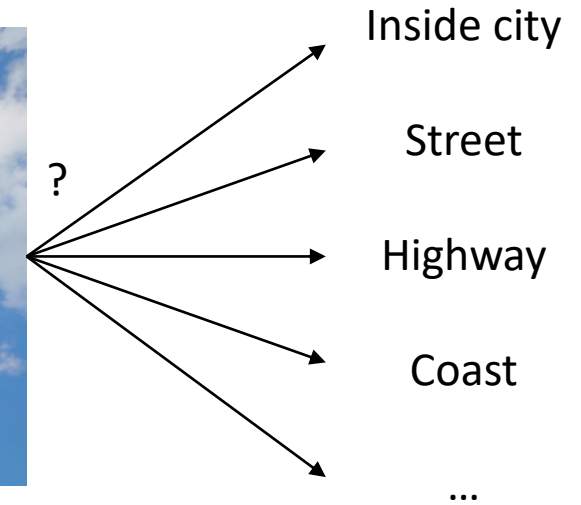
few training data



moderate amount of training data

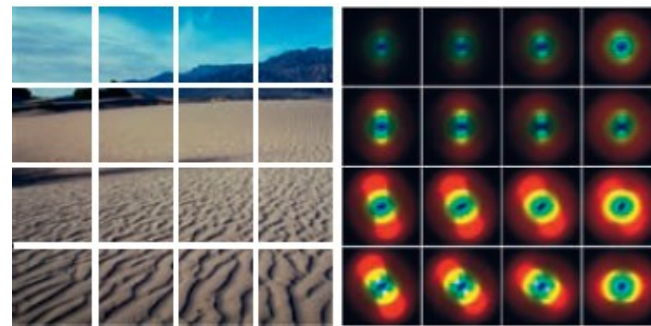
Scene Recognition

- The most basic form of localization;
- Tells what kind of scene the user is in;
- Useful to distinguish between (even for unseen places) :
 - indoor/outdoor
 - natural/artificial
 - conf. room
 - Office
- Can use off-the-shelf detections.



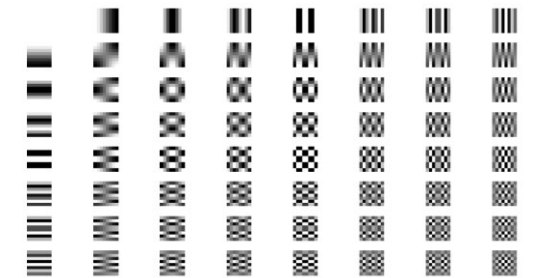
COMPUTATIONALLY INEXPENSIVE ALGORITHMS

GIST Descriptor



Oliva, Aude, and Antonio Torralba. "Modeling the shape of the scene: A holistic representation of the spatial envelope." *International journal of computer vision* 42.3 (2001): 145-175.

DCT-GIST (runs on the IGP pipeline)



G. M. Farinella, D. Ravi, V. Tomaselli, M. Guarnera, S. Battiato, "Representing scenes for real-time context classification on mobile devices", *Pattern Recognition*, Elsevier, ISSN 0031-3203, Vol. 48, N. 4, pp. 1082-1096, doi: 10.1016/j.patcog.2014.05.014, 2015

DATA & CODE HERE -> <http://places2.csail.mit.edu/>

Scene Recognition – Places



GT: cafeteria

top-1: cafeteria (0.179)

top-2: restaurant (0.167)

top-3: dining hall (0.091)

top-4: coffee shop (0.086)

top-5: restaurant patio (0.080)

- Places is a large (10M images – 400+ classes) dataset for scene recognition;
- CNN models trained to recognize 365 scene classes available for download;
- Can be used off-the-shelf!

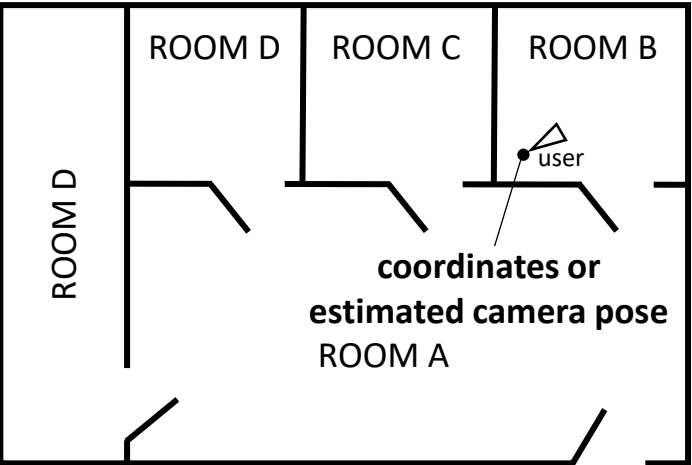
Localization – Levels of Granularity

SCENE RECOGNITION

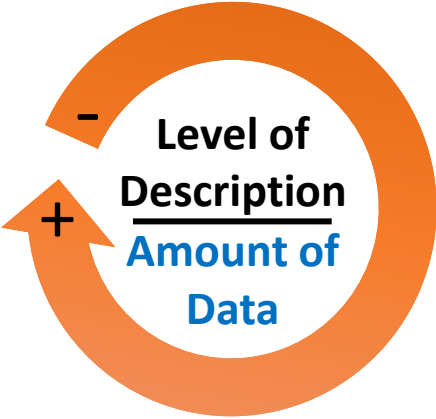


off-the-shelf detectors

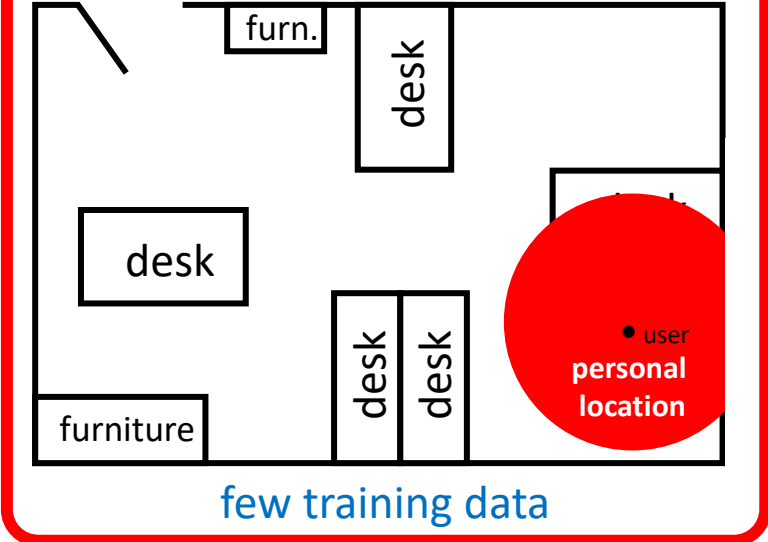
CAMERA POSE-ESTIMATION



3D reconstruction of the building

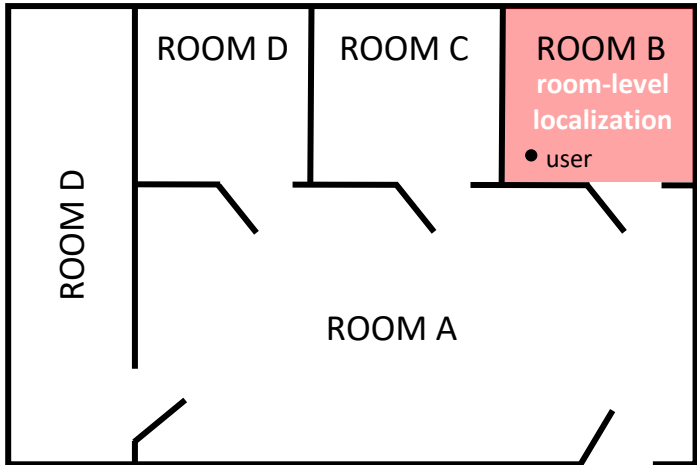


PERSONAL LOCATION RECOGNITION



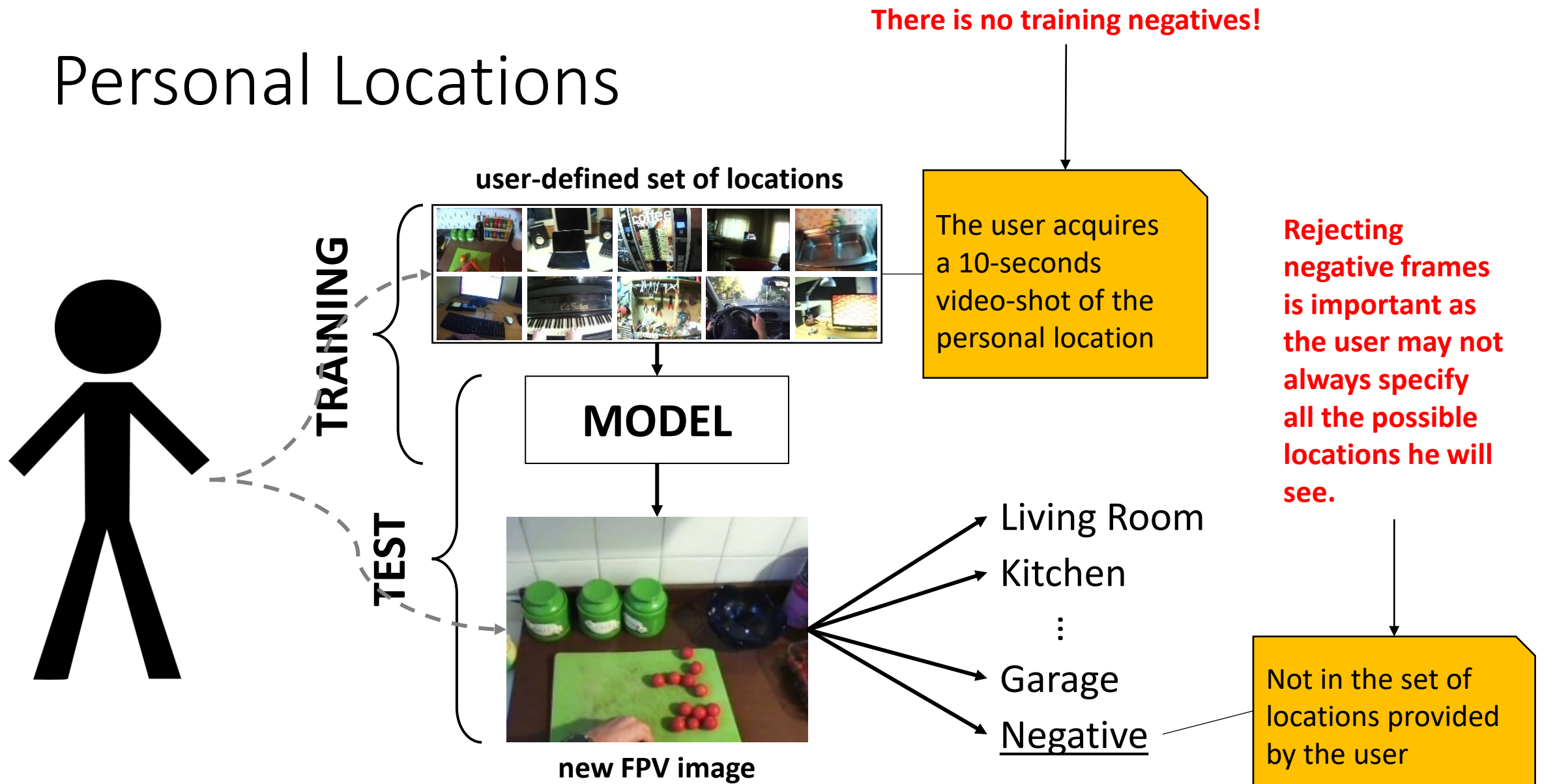
few training data

ROOM-LEVEL RECOGNITION

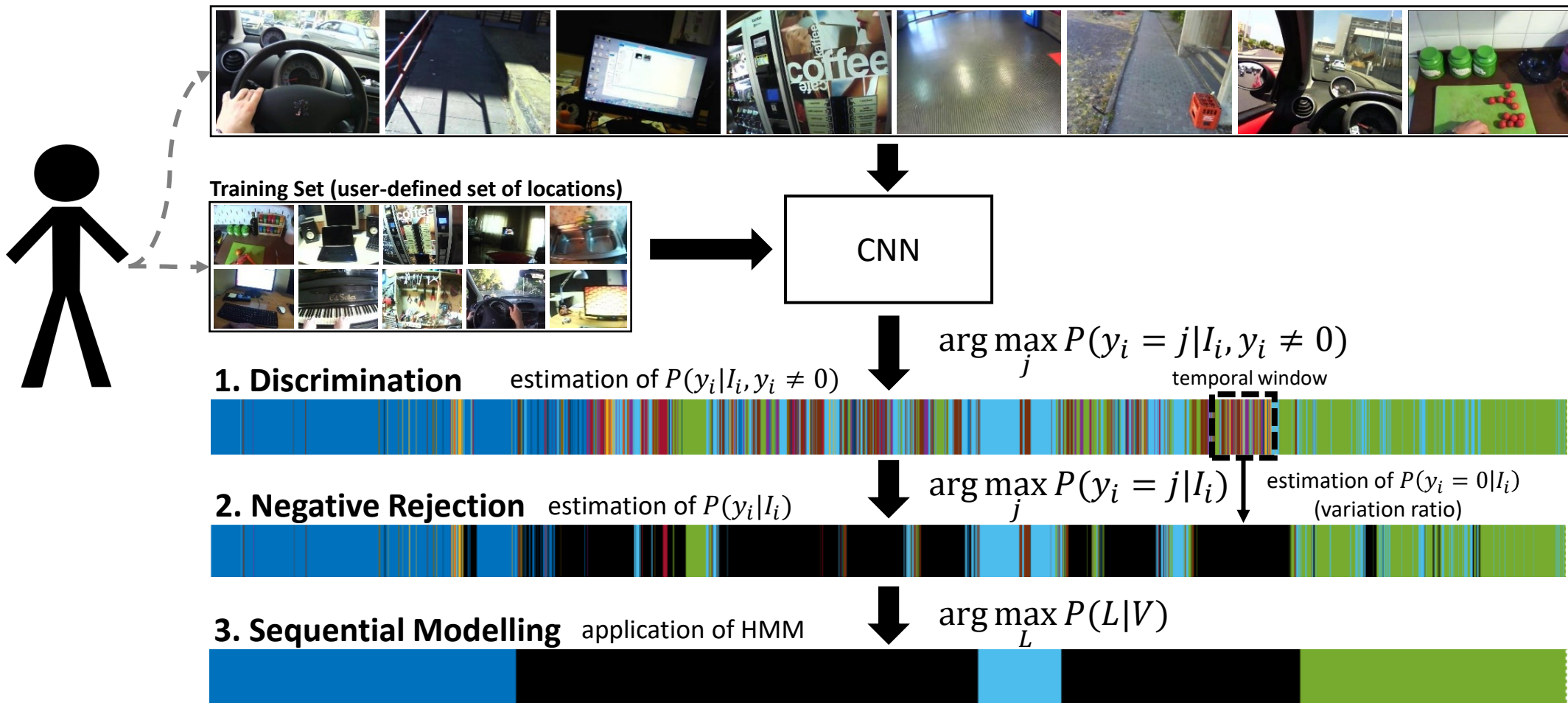


moderate amount of training data

Personal Locations



Personal Locations – Full Model



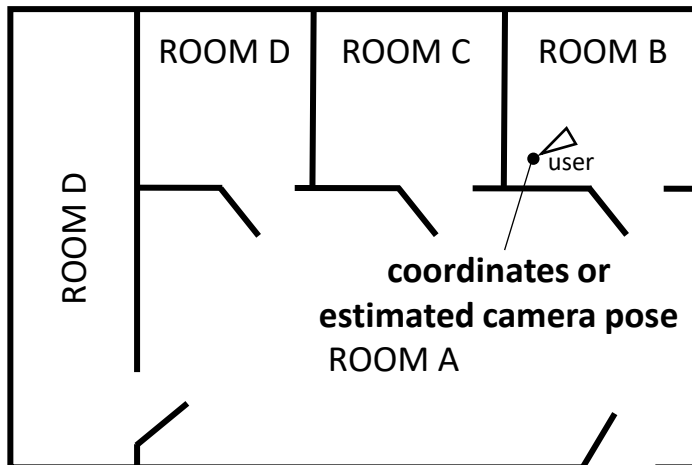
Localization – Levels of Granularity

SCENE RECOGNITION



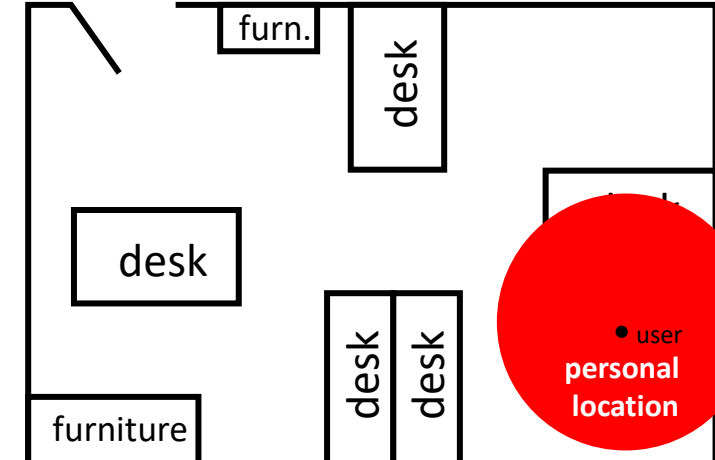
off-the-shelf detectors

CAMERA POSE-ESTIMATION



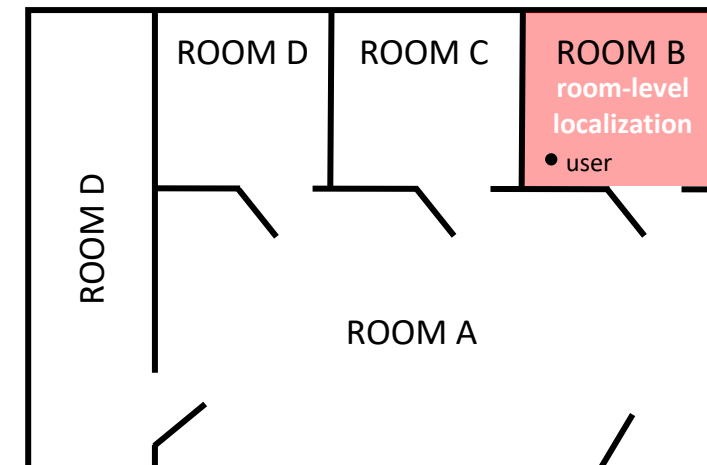
3D reconstruction of the building

PERSONAL LOCATION RECOGNITION

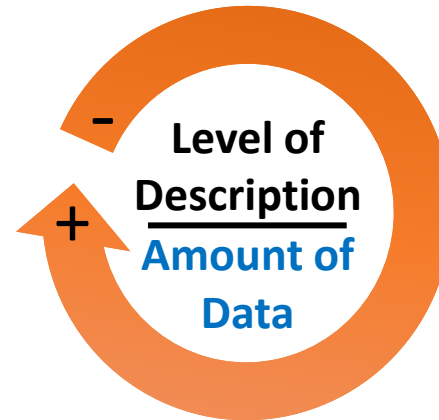


few training data

ROOM-LEVEL RECOGNITION



moderate amount of training data



Room-Level Localization

Localizing the user in a larger environment (e.g., a museum).

Extending Personal Location Recognition to Room-Level Localization:

- Collect a longer training video for each room including different points of view;
- Same algorithm as before



VEDI – Vision Exploitation for Data Interpretation, PON MISE Horizon 2020
 F. Ragusa, A. Furnari, S. Battiato, G. Signorello, G. M. Farinella

Time Spent at Location
 LOC EST GT

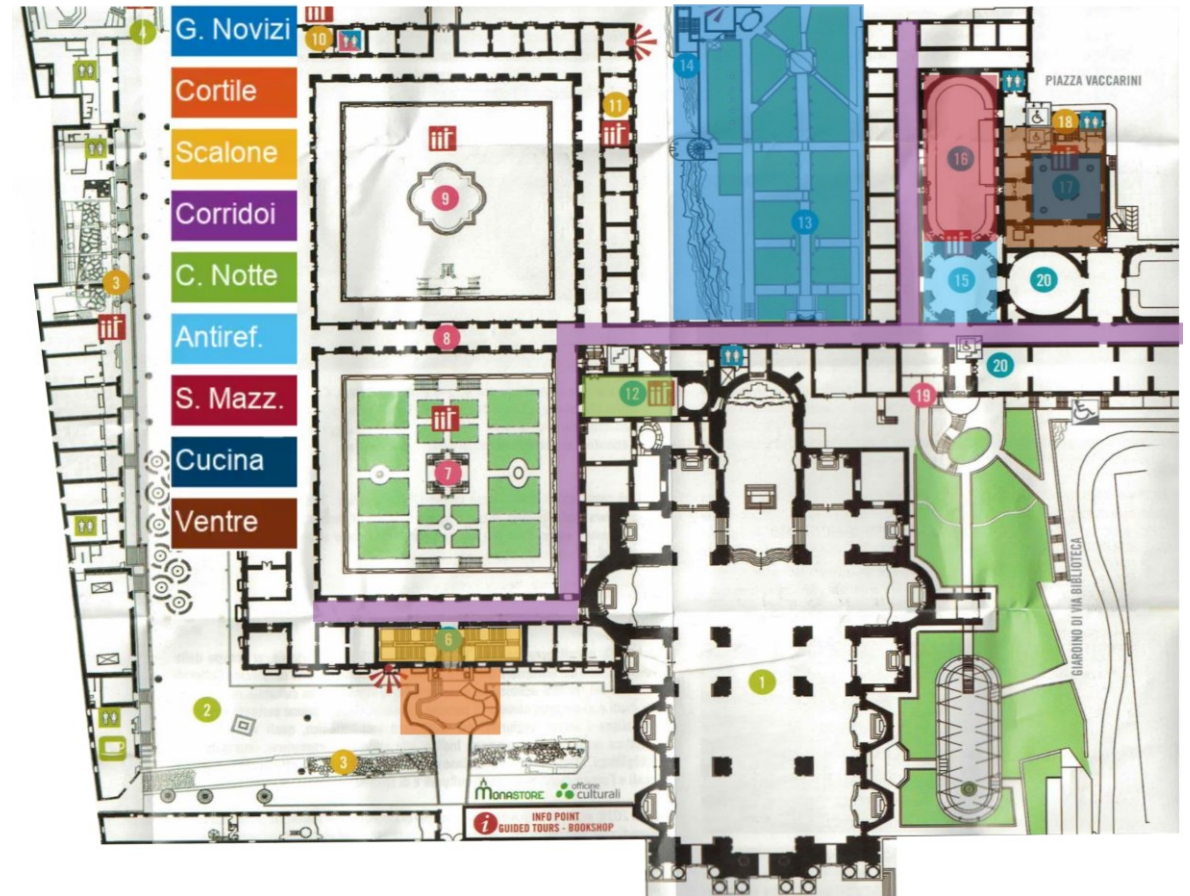
G. Novizi	00:00	00:00
Cortile	00:03	00:03
Scalone	00:00	00:00
Corridoi	00:00	00:00
C. Notte	00:00	00:00
Antiref.	00:00	00:00
S. Mazz.	00:00	00:00
Cucina	00:00	00:00
Ventre	00:00	00:00
Negative	00:00	00:00



Detected Shots for Storyboard Summary



Estimated Probabilities	Predicted Class	GT Class
Giardino dei Novizi		
Cortile	●	●
Scalone Monumentale		
Corridoi		
Coro di Notte		
Antirefettorio		
Aula Santo Mazzarino		
Cucina		
Ventre		
Negative		



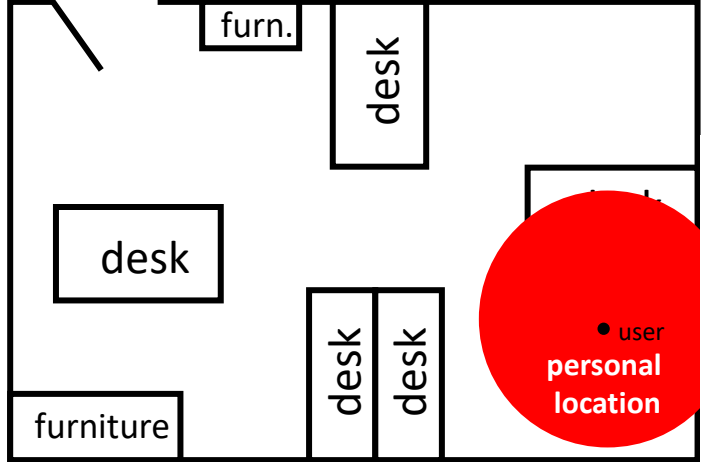
Localization – Levels of Granularity

SCENE RECOGNITION



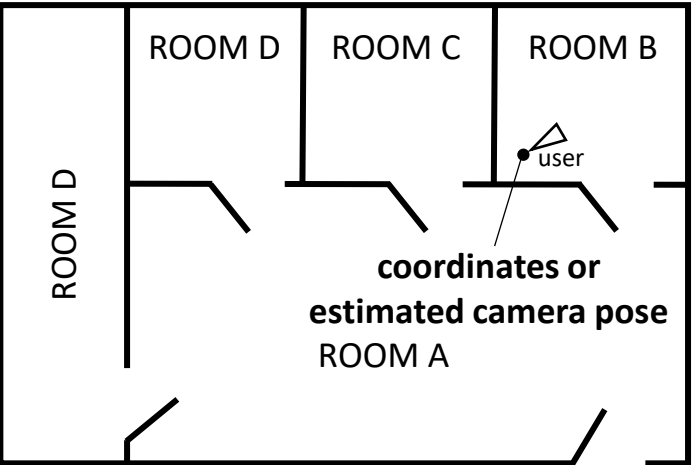
off-the-shelf detectors

PERSONAL LOCATION RECOGNITION

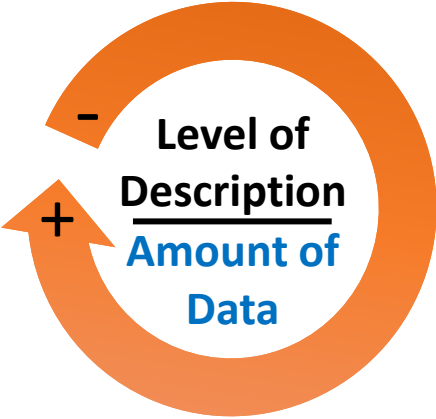


few training data

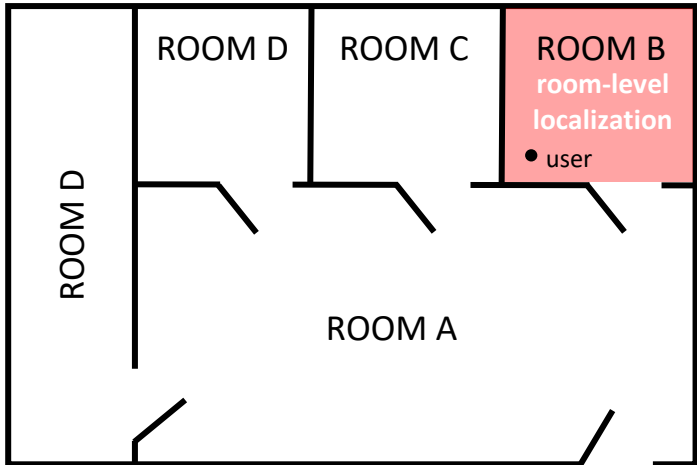
CAMERA POSE-ESTIMATION



3D reconstruction of the building



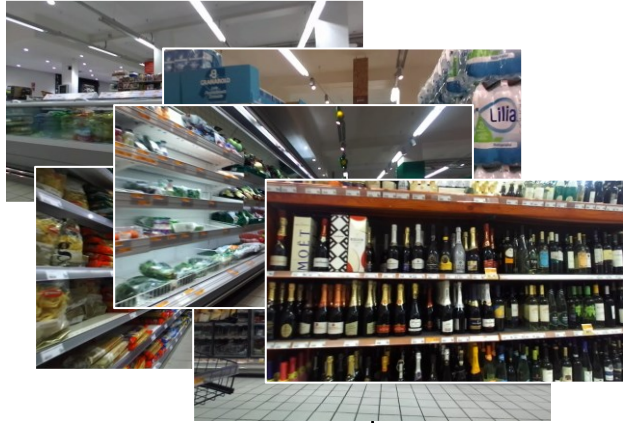
ROOM-LEVEL RECOGNITION



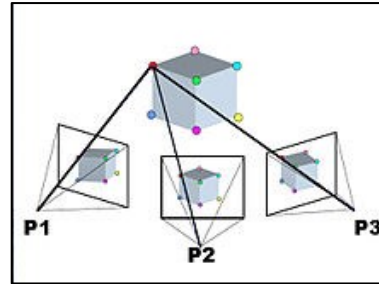
moderate amount of training data

Camera Pose Estimation – Dataset Creation

Images



Structure from Motion (SfM)



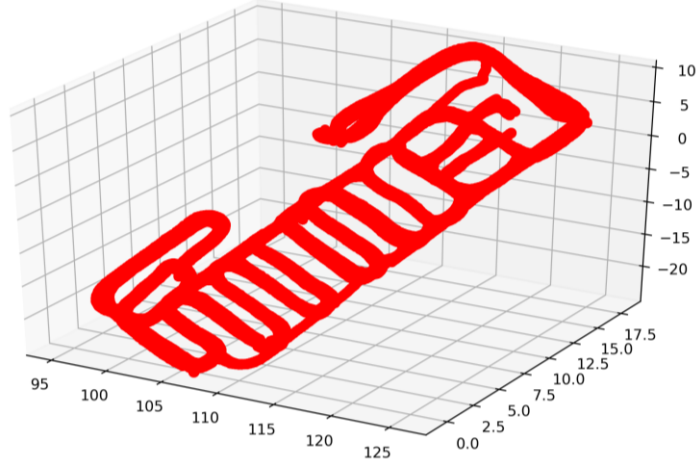
(P,Q)

Attach estimated 6DOF pose to each image

3D Model

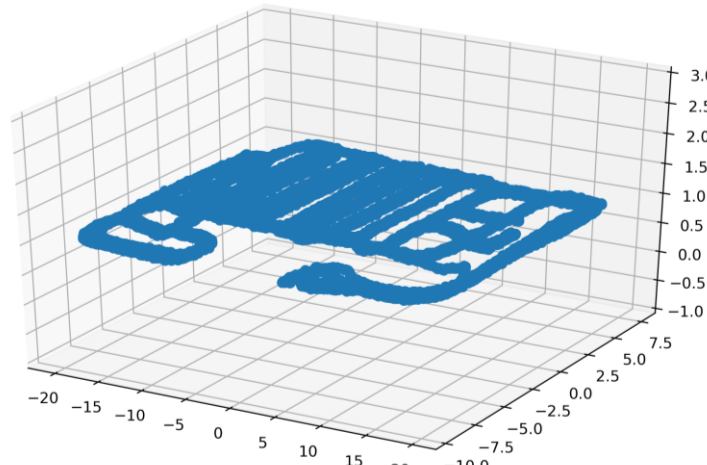


Arbitrary Coordinate System (pose/scale)

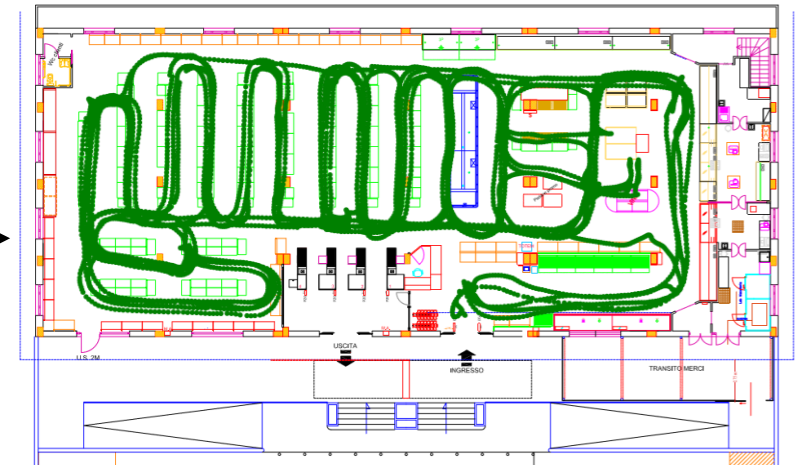


camera poses

PCA



rotated poses



scaled/aligned poses

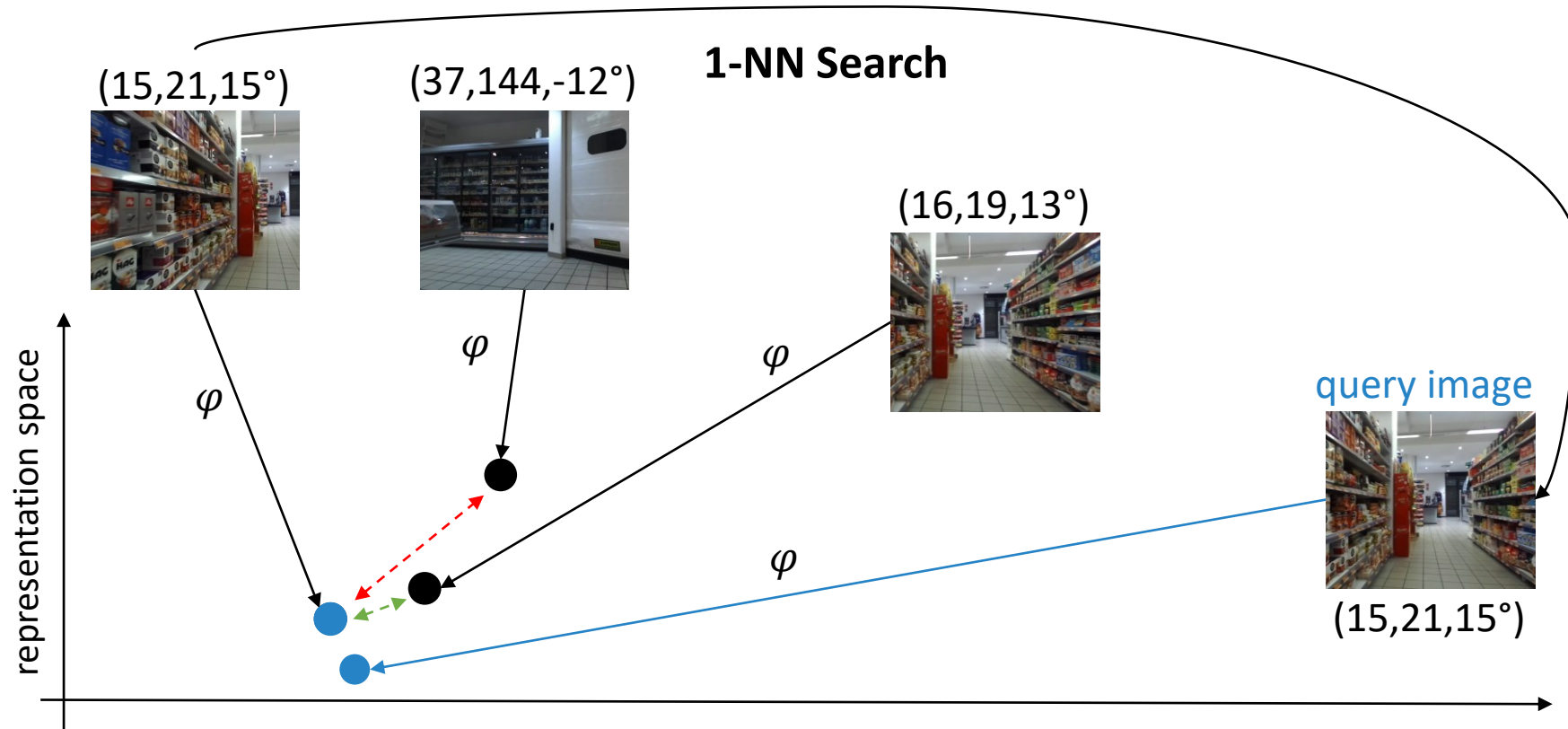
Structure from Motion (SfM) Softwares

Many options available:

- COLMAP (free)
 - <https://colmap.github.io/>
- Visual SFM (free)
 - <http://ccwu.me/vsfm/>
- 3D Zephyr (paid)
 - <https://www.3dflow.net/it/3df-zephyr-pro-3d-models-from-photos/>

Camera Pose Estimation – Retrieval Approach

Use deep metric learning to learn a representation function φ which maps close to each other images of nearby locations

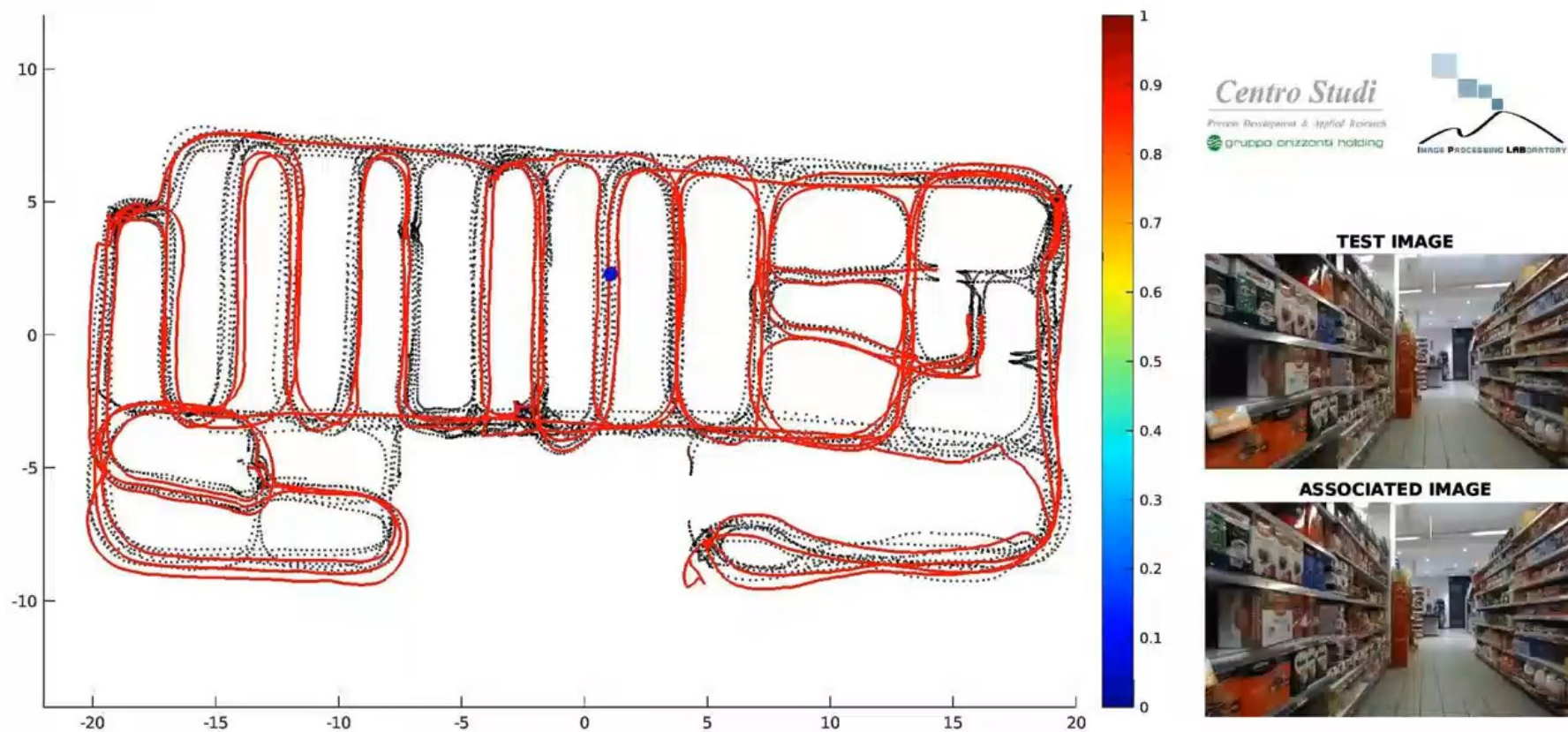


Camera Pose Estimation – Demo

EGOCENTRIC SHOPPING CART LOCALIZATION

Emiliano Spera, Antonino Furnari, Sebastiano Battiato, Giovanni Maria Farinella

<http://iplab.dmi.unict.it/EgocentricShoppingCartLocalization/>



Other approaches to visual localization



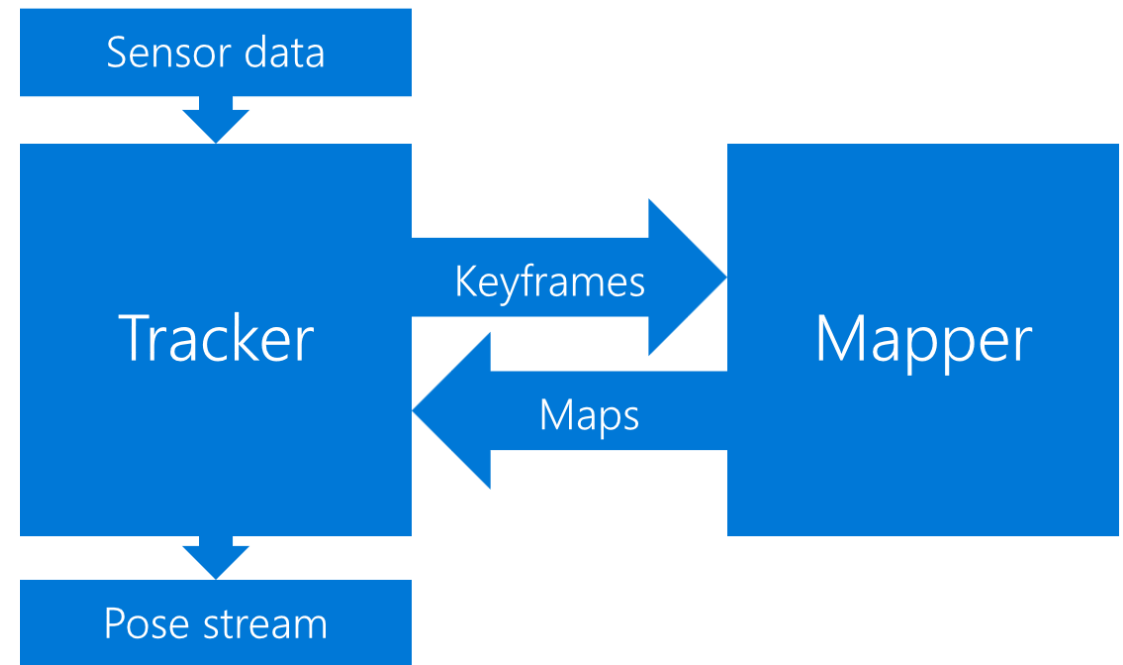
Literature is rich. See here -> <https://sites.google.com/view/lsvpr2019/home>

Camera Pose Estimation - HoloLens

Microsoft HoloLens implements a localization system which is used for augmented reality.

- The mapper continuously refines a map of the environment;
- The tracker sees small local sub-maps;

Microsoft provides an API to access HoloLens's location information, which can be used by other apps.

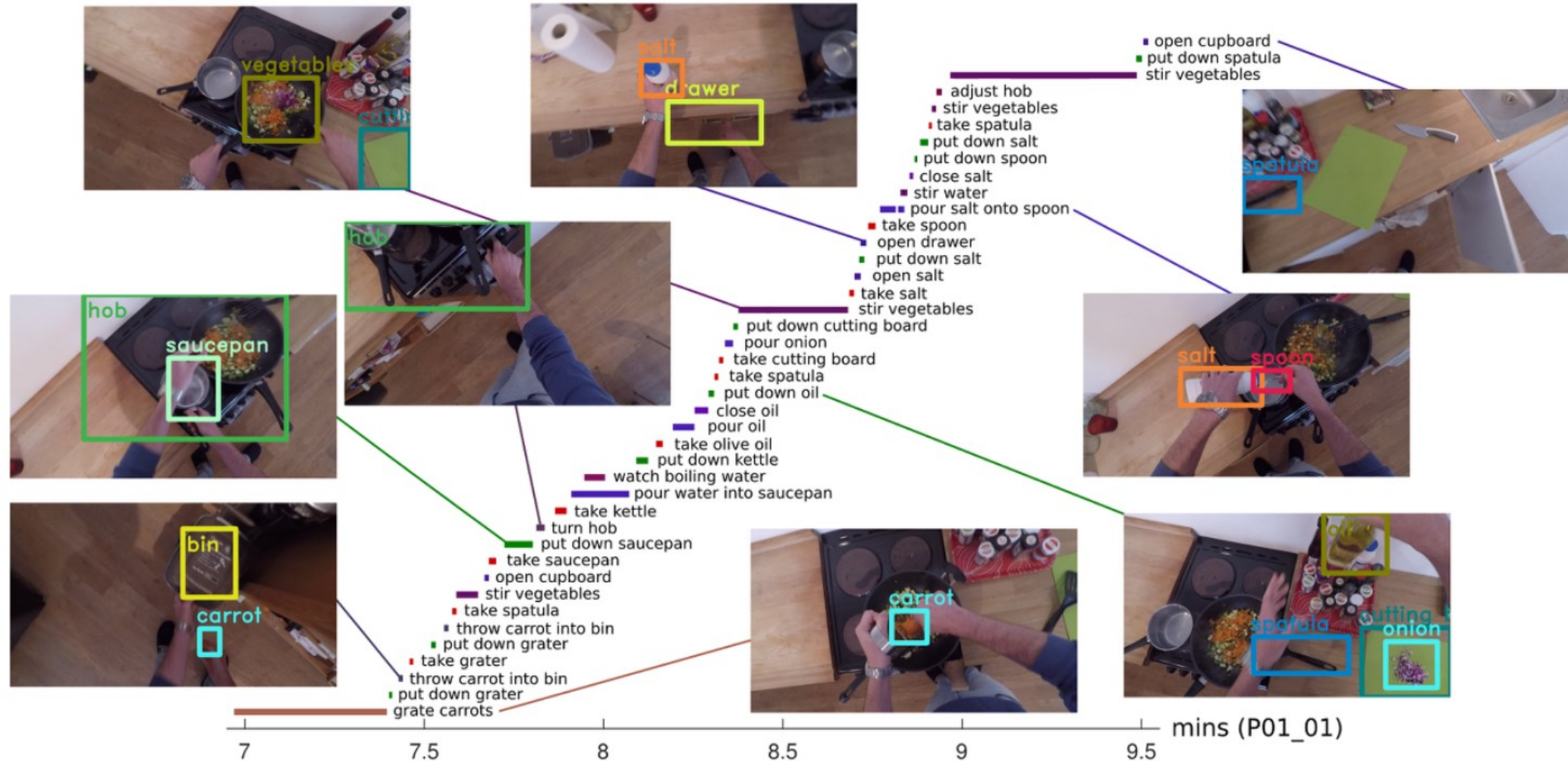


Objects and Actions are tight!

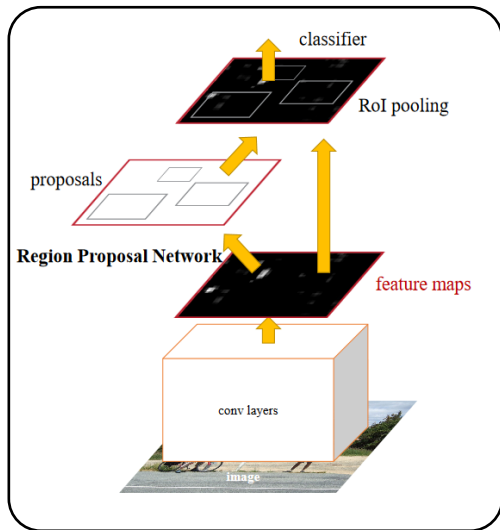
Useful to know what is in the scene

Useful to know what actions can be performed

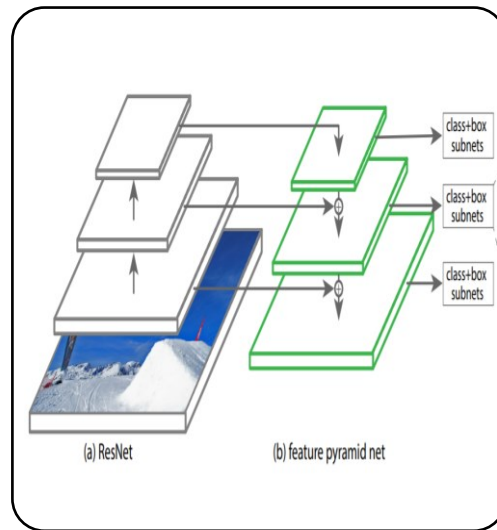
Object Detection



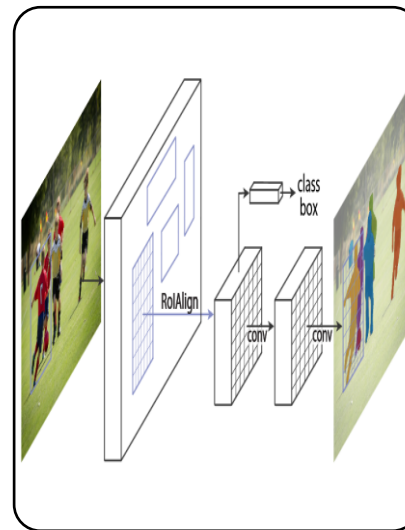
Off-the-shelf object detectors



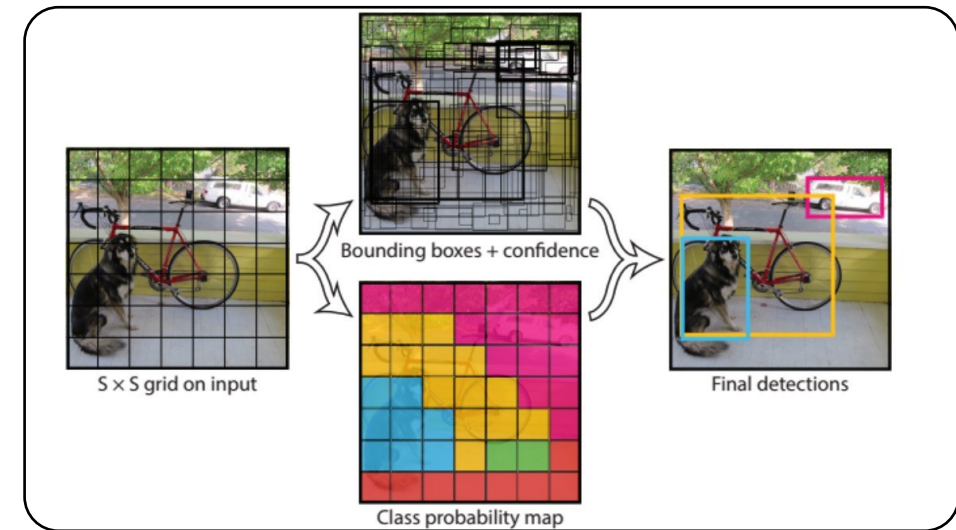
Faster-RCNN
(bounding boxes)



RetinaNet
(bounding boxes - faster)



Mask-RCNN
(boxes + segments)



YOLO
(much faster, but less accurate)

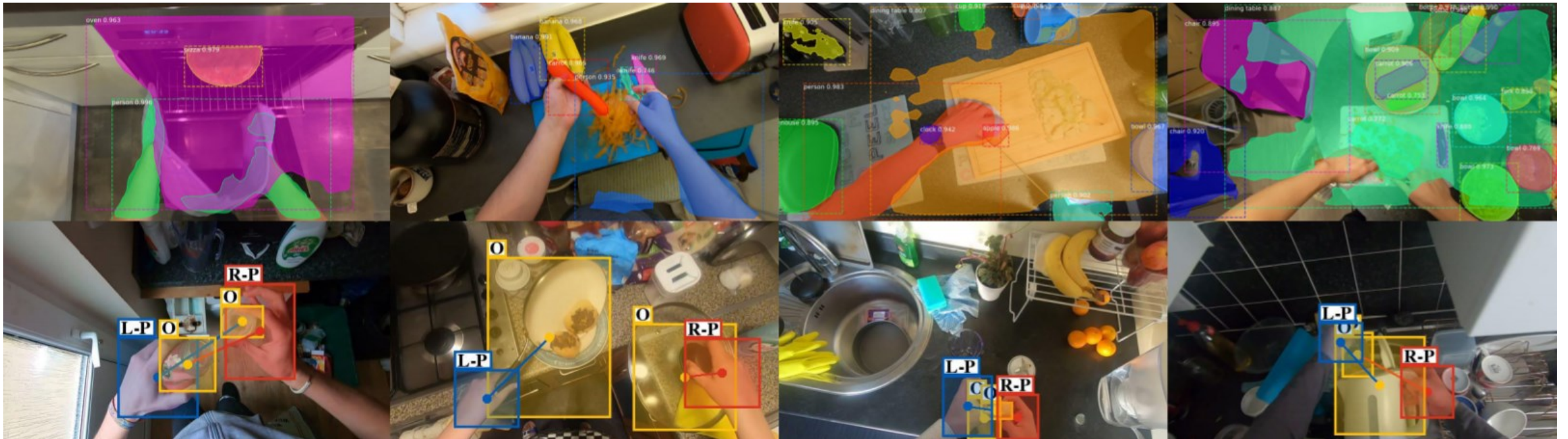
<https://github.com/facebookresearch/detectron2>

<https://pjreddie.com/darknet/yolo/>

Ren, S., He, K., Girshick, R., & Sun, J. (2015). Faster r-cnn: Towards real-time object detection with region proposal networks. In *NIPS*.
Joseph Redmon, Ali Farhadi, YOLO9000: Better, Faster, Stronger, The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2017
He, K., Gkioxari, G., Dollár, P., & Girshick, R. (2017, October). Mask r-cnn. In *Computer Vision (ICCV), 2017* (pp. 2980-2988). IEEE.

Off-the-shelf detectors on EPIC-KITCHENS

Depending on the scenario, off-the-shelf detectors can be a starting point, but they are not always accurate.



Damen, Doughty, Farinella, Furnari, Kazakos, Moltisanti, Munro, Price, Wray (2020). Rescaling Egocentric Vision. *arXiv preprint arXiv:2006.13256* (2020).

Train/Finetune your own object detector



<https://www.csee.umbc.edu/~hpirsiav/papers/ADLdataset/>



<http://epic-kitchens.github.io/>



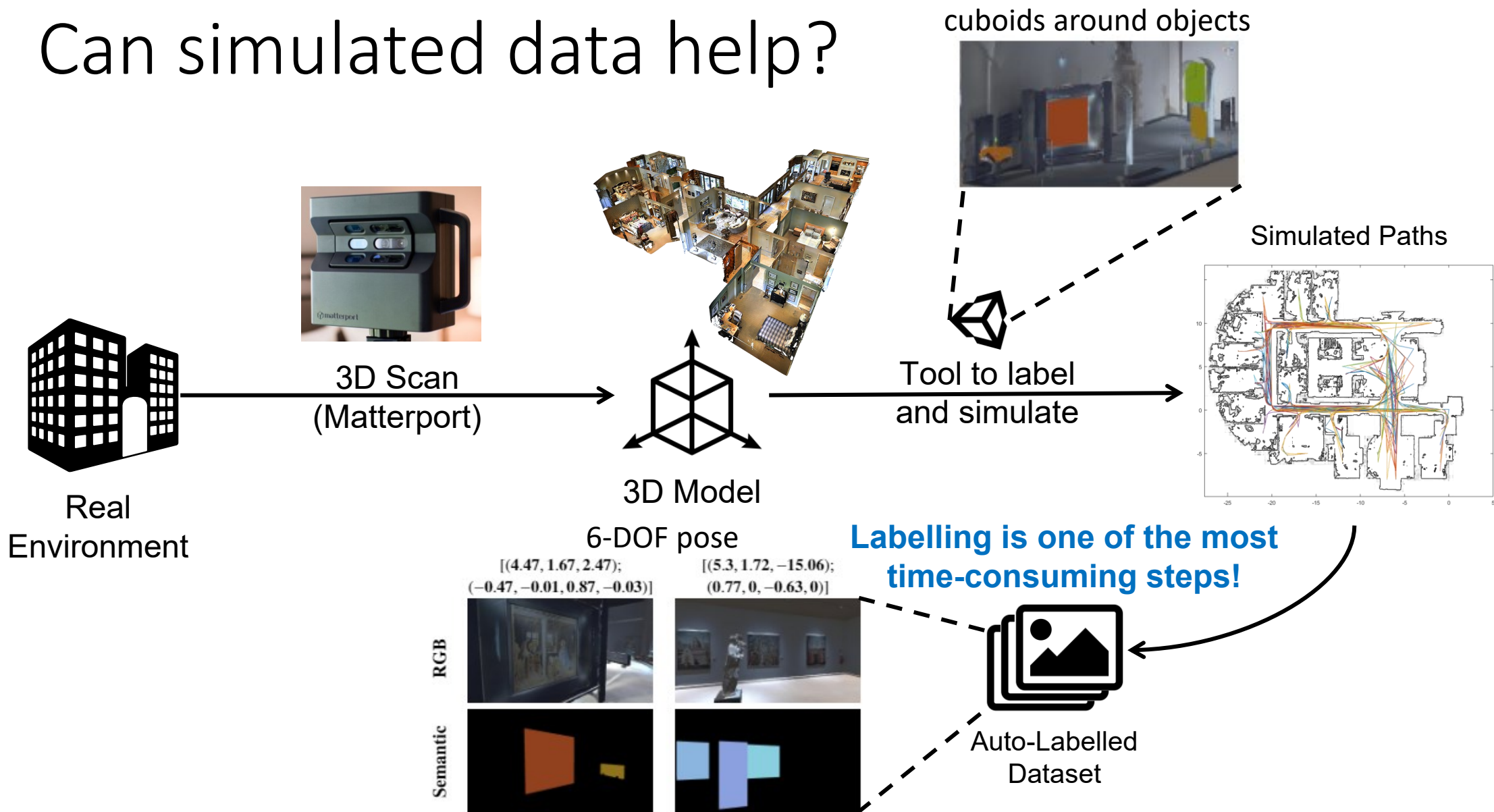
<https://iplab.dmi.unict.it/EGO-CH/>



<https://iplab.dmi.unict.it/MECCANO/>

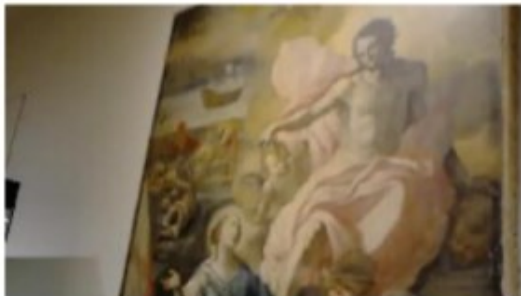
- In some scenario, it could be necessary to fine-tune an object-detector with application-specific data.
- On the left: main egocentric datasets providing bounding box annotations.
- Recently, EGO4D has been released and it has been annotated with bounding boxes.

Can simulated data help?

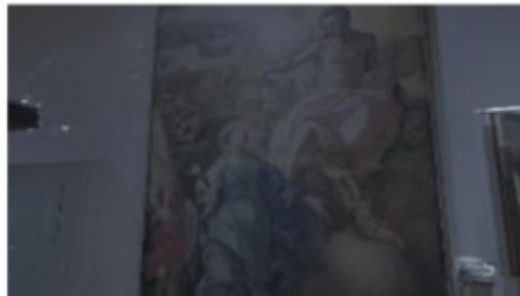
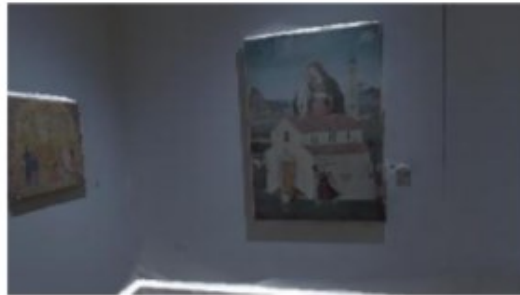


Domain Gap Between Real and Synthetic Images

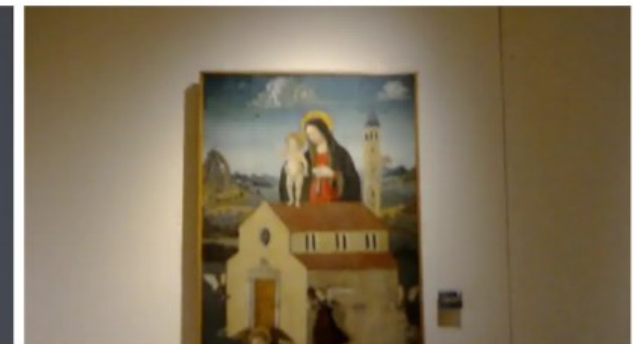
Real Image



Ref. Synthetic

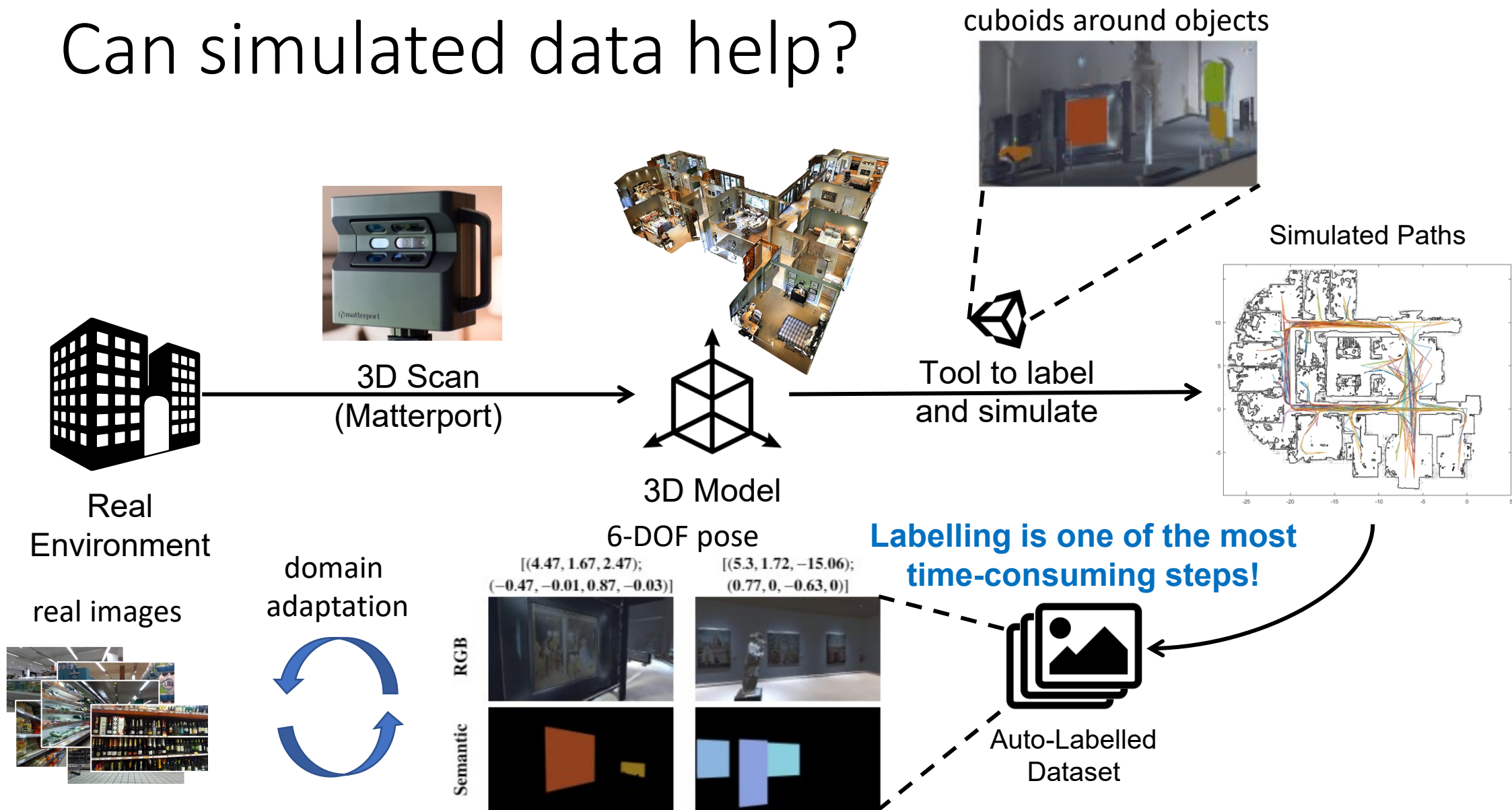


- Synthetic Images are Easier to collect.
- However, there is significant domain GAP between real and synthetic images;
- As a result, methods trained only on synthetic images will not work directly with real ones.



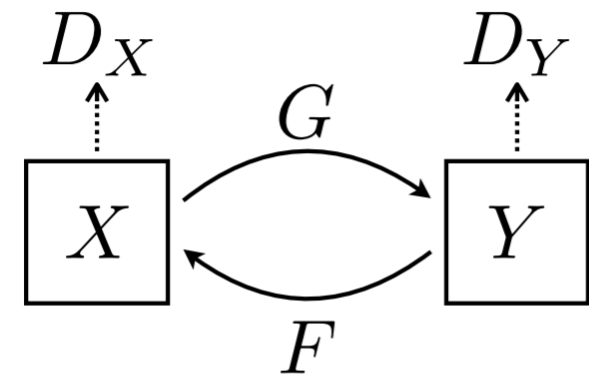
An object detection method trained on synthetic images (left), does not perform well on real images (right)

Can simulated data help?



Domain Adaptation Through Image to Image Translation

- One way to reduce the domain gap is by making the synthetic and real images look similar.
- This can be done using image to image translation.
- CycleGAN uses two discriminators to learn two mappings (from synthetic to real and vice versa) which make transformed images undistinguishable.

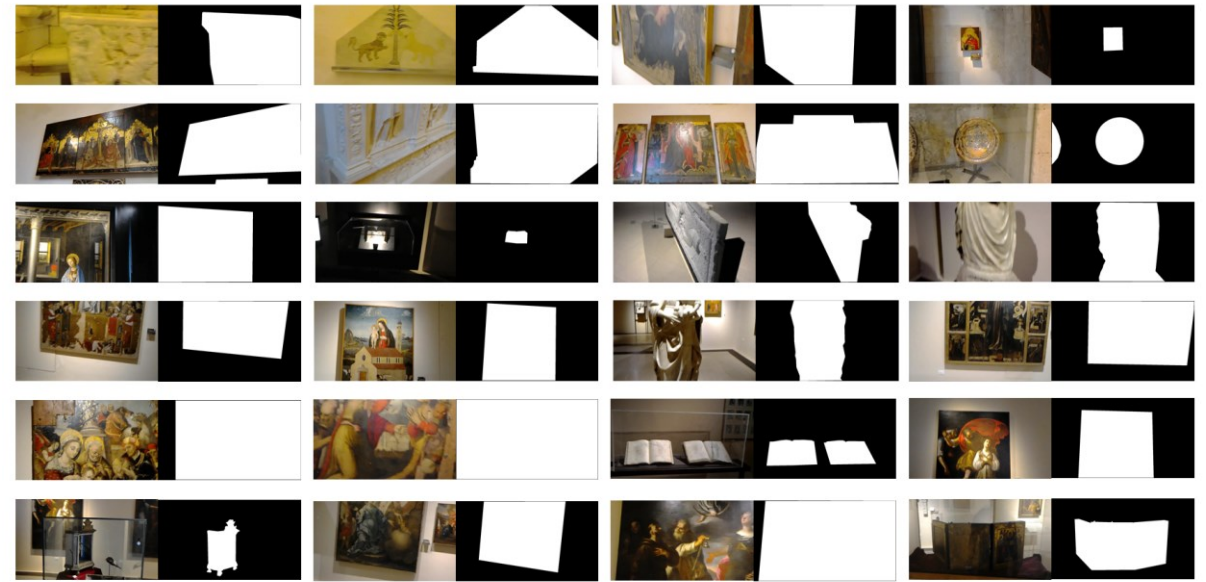


DATA HERE -> <https://iplab.dmi.unict.it/EGO-CH-OBJ-SEG/>

Domain Adaptation for Semantic Object Segmentation Dataset



Synthetic Images



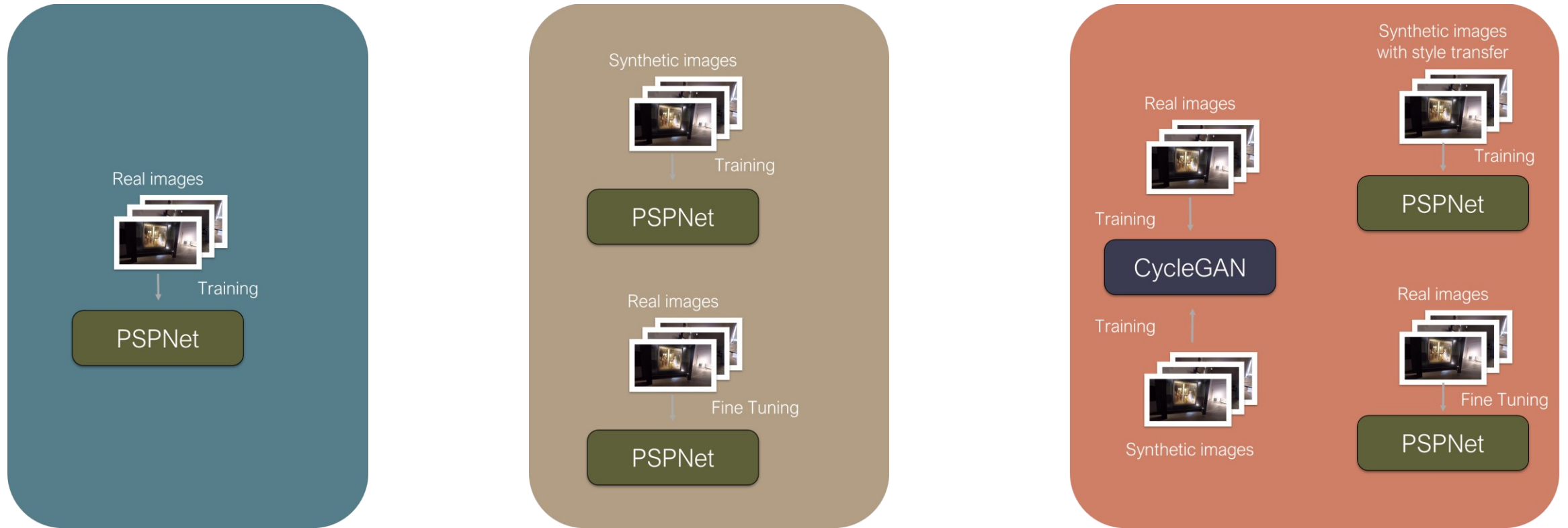
Real Images

24 objects, ~25k synthetic images, ~5k real labeled images, semantic segmentations masks

Francesco Ragusa, Daniele DiMauro, Alfio Palermo, Antonino Furnari, Giovanni Maria Farinella (2020). Semantic Object Segmentation in Cultural Sites using Real and Synthetic Data. International Conference on Pattern Recognition (ICPR).

DATA HERE -> <https://iplab.dmi.unict.it/EGO-CH-OBJ-SEG/>

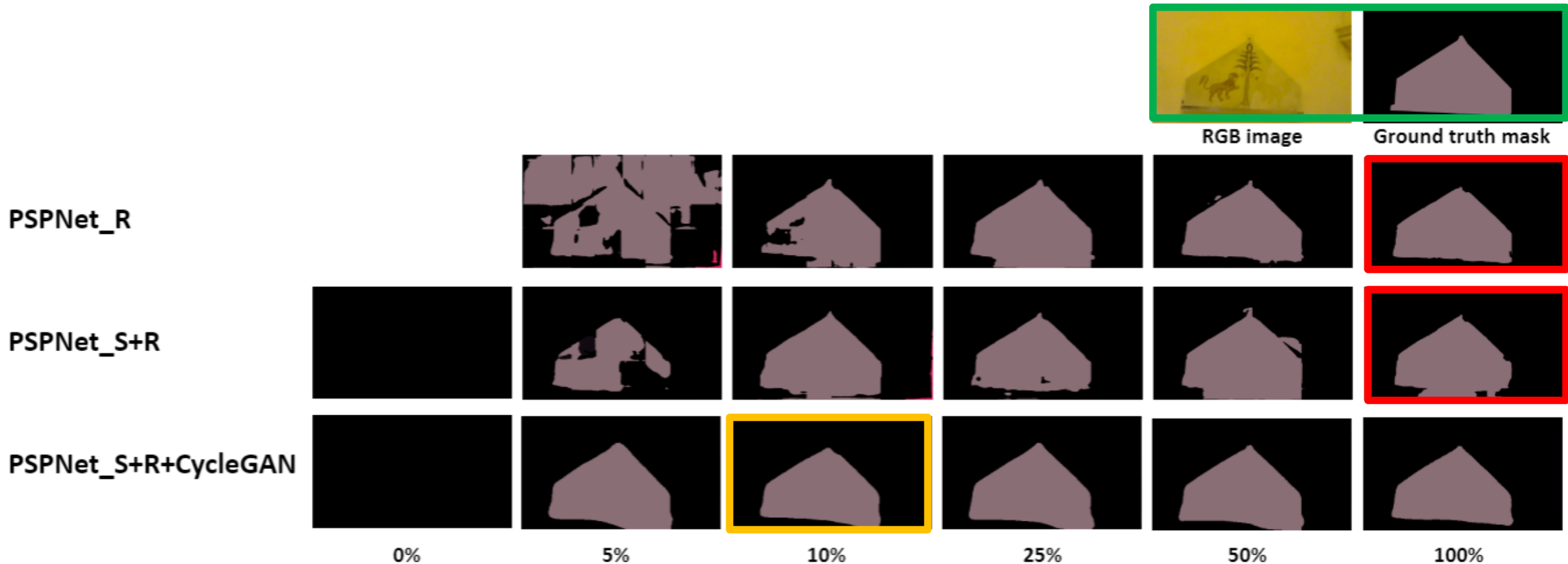
Domain Adaptation for Semantic Object Segmentation Dataset



Francesco Ragusa, Daniele DiMauro, Alfio Palermo, Antonino Furnari, Giovanni Maria Farinella (2020). Semantic Object Segmentation in Cultural Sites using Real and Synthetic Data. International Conference on Pattern Recognition (ICPR).

DATA HERE -> <https://iplab.dmi.unict.it/EGO-CH-OBJ-SEG/>

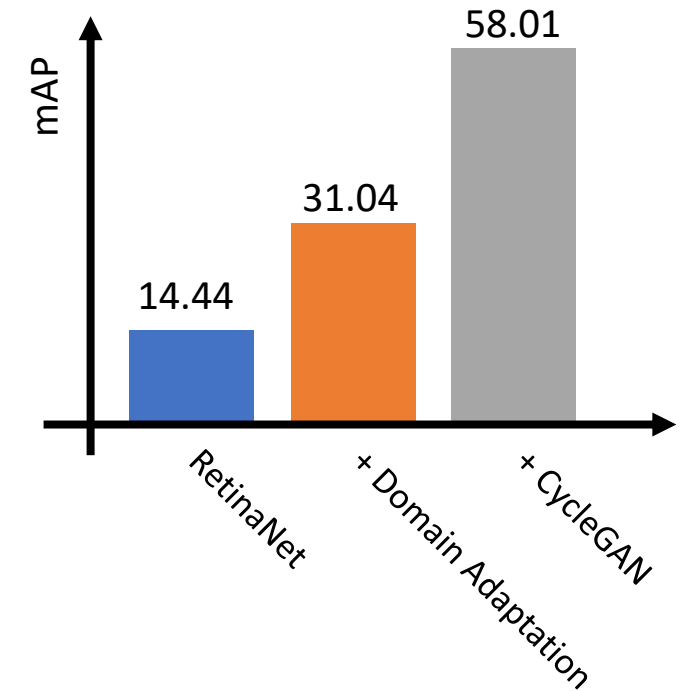
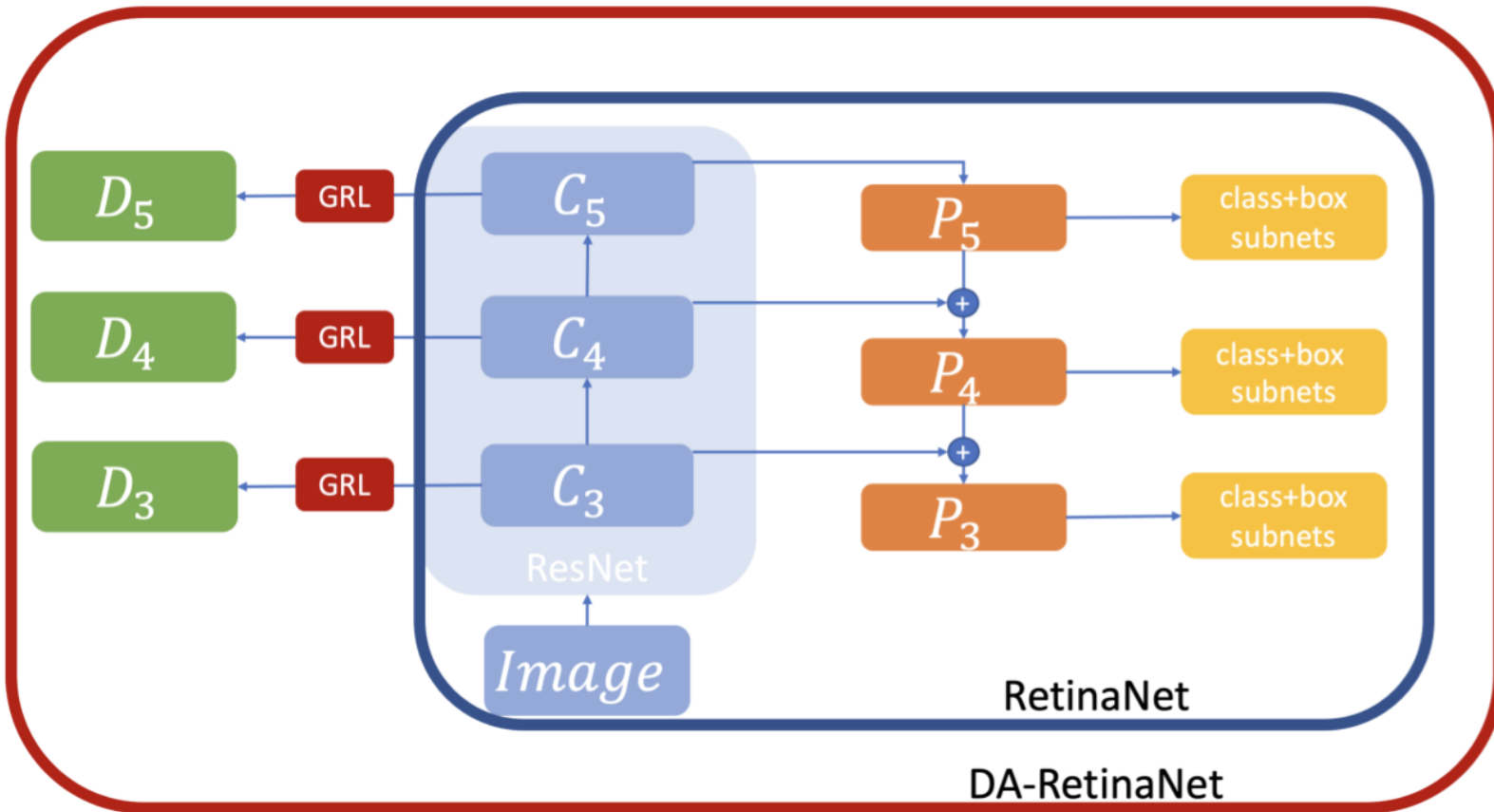
Domain Adaptation for Semantic Object Segmentation Dataset



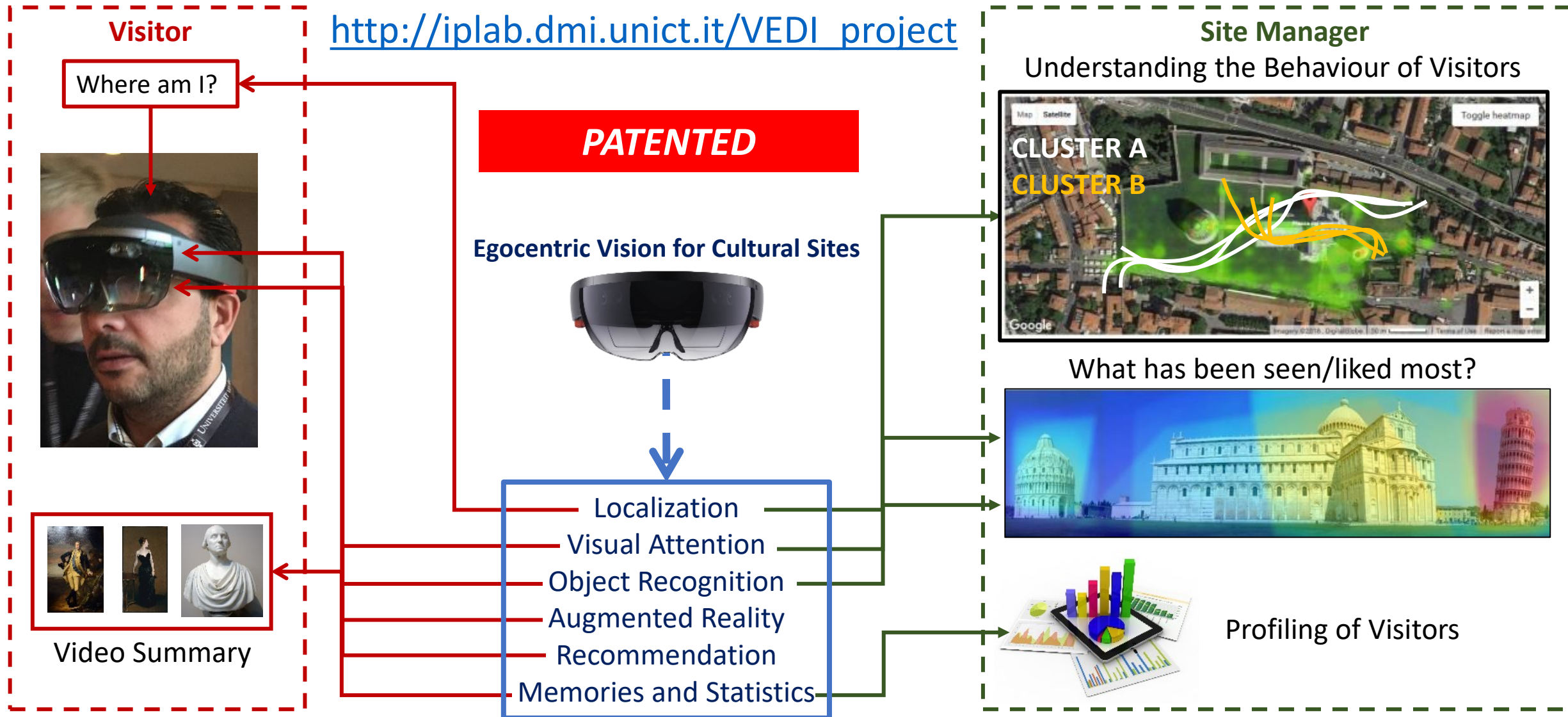
Francesco Ragusa, Daniele DiMauro, Alfio Palermo, Antonino Furnari, Giovanni Maria Farinella (2020). Semantic Object Segmentation in Cultural Sites using Real and Synthetic Data. International Conference on Pattern Recognition (ICPR).

Unsupervised Domain Adaptation for Object Detection

We applied domain adaptation through adversarial learning to the RetinaNet object detector and combined it with image-to-image translation.



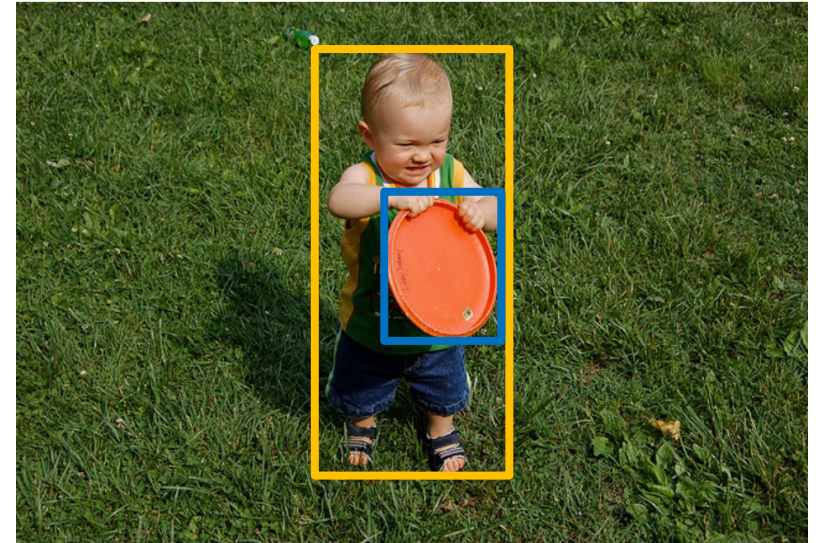
Vision Exploitation for Data Interpretation (VEDI)



Human-Object Interaction



<human, talks, cellphone>



<human, holds, frisbee>

Egocentric Human-Object Interaction

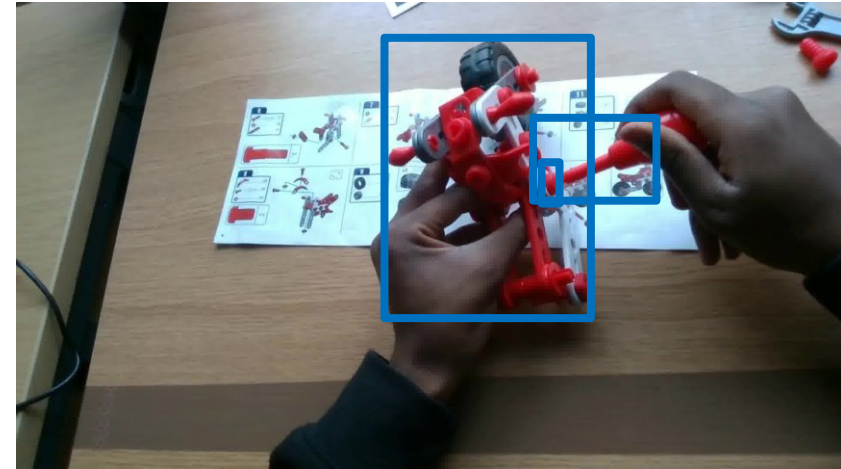
$$O = \{o_1, o_2, \dots, o_n\}$$

$$V = \{v_1, v_2, \dots, v_m\}$$

$$e = (v_h, \{o_1, o_2, \dots, o_i\})$$

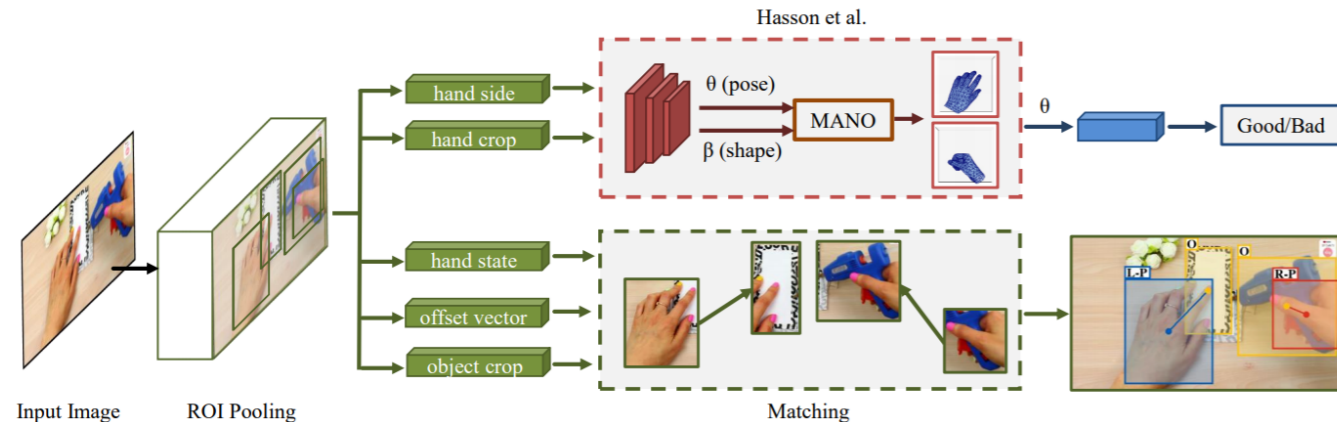
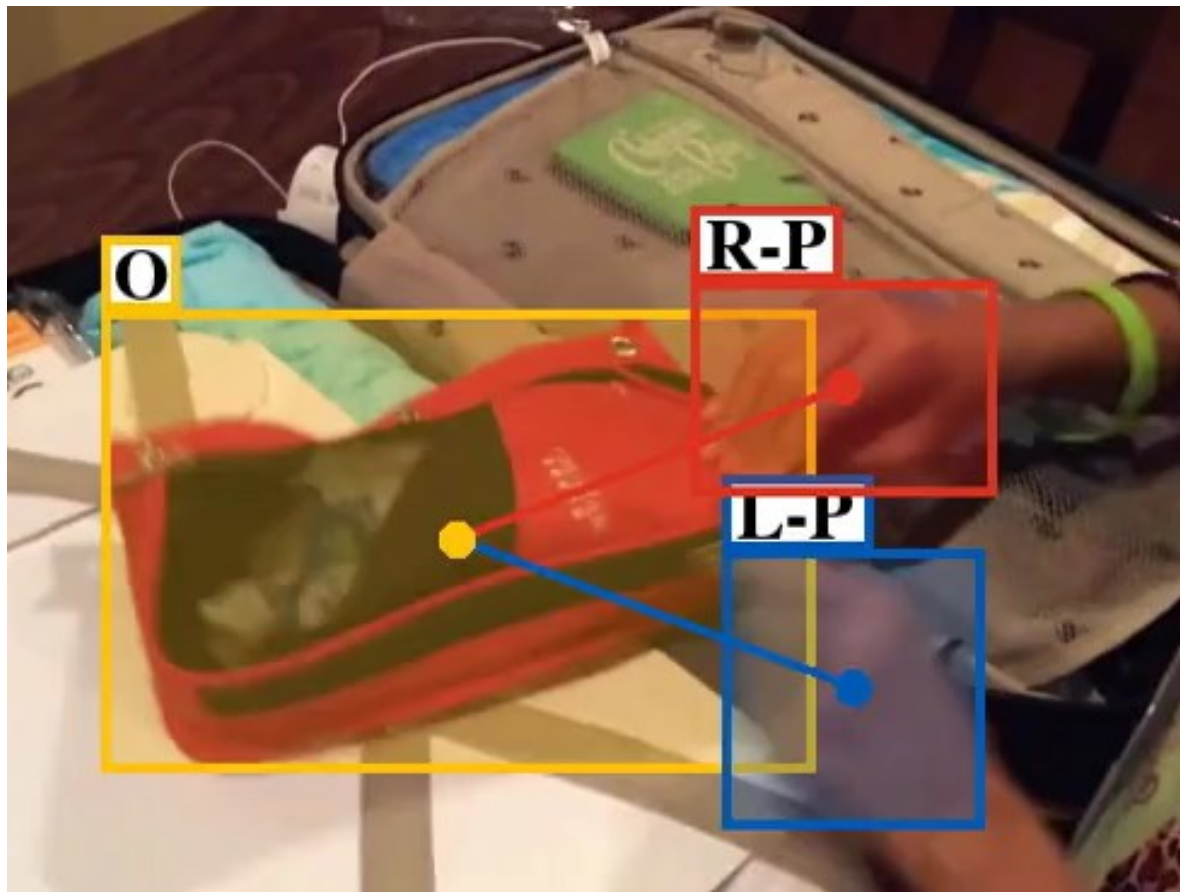


<take, screwdriver>



<screw, {screwdriver, screw, partial_model}>

Hands in Contact – Hands + Objects

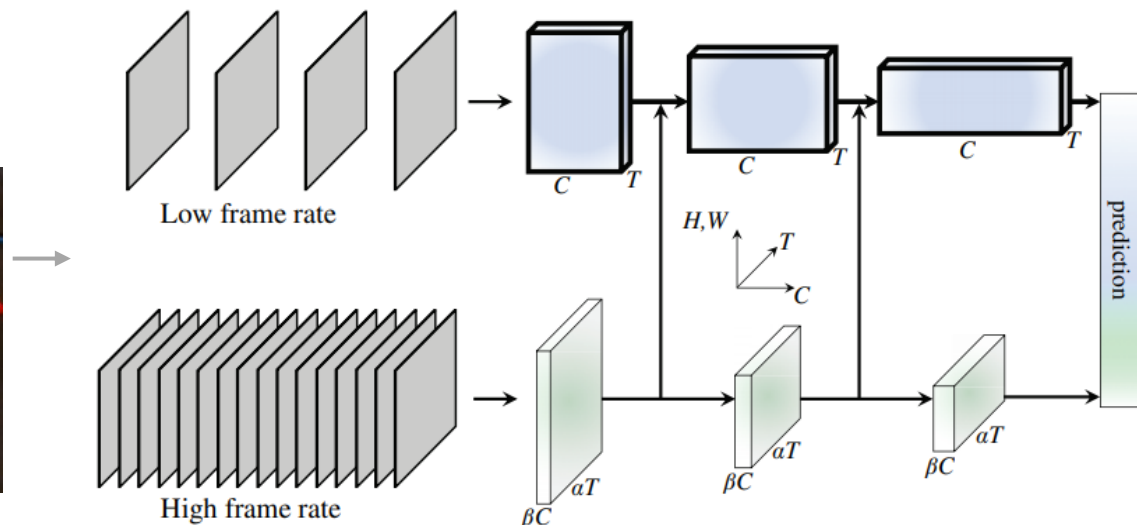


An «augmented» detector which recognizes:

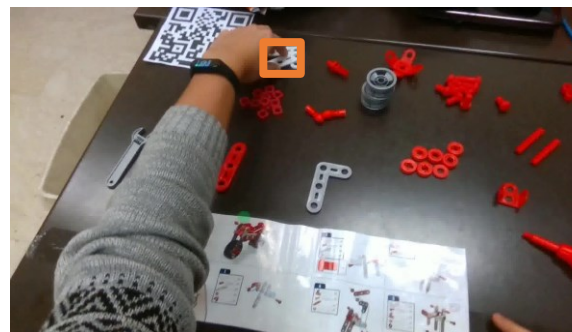
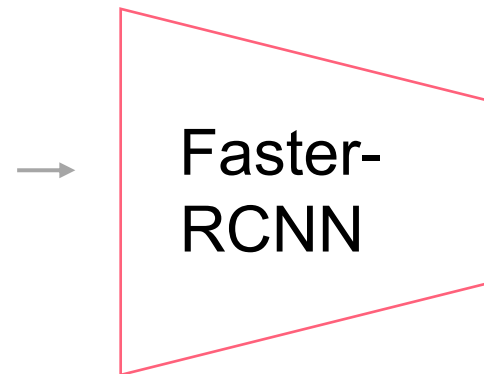
- The left hand;
- The right hand;
- The interacted object.

Egocentric Human-Object Interaction

Video Based

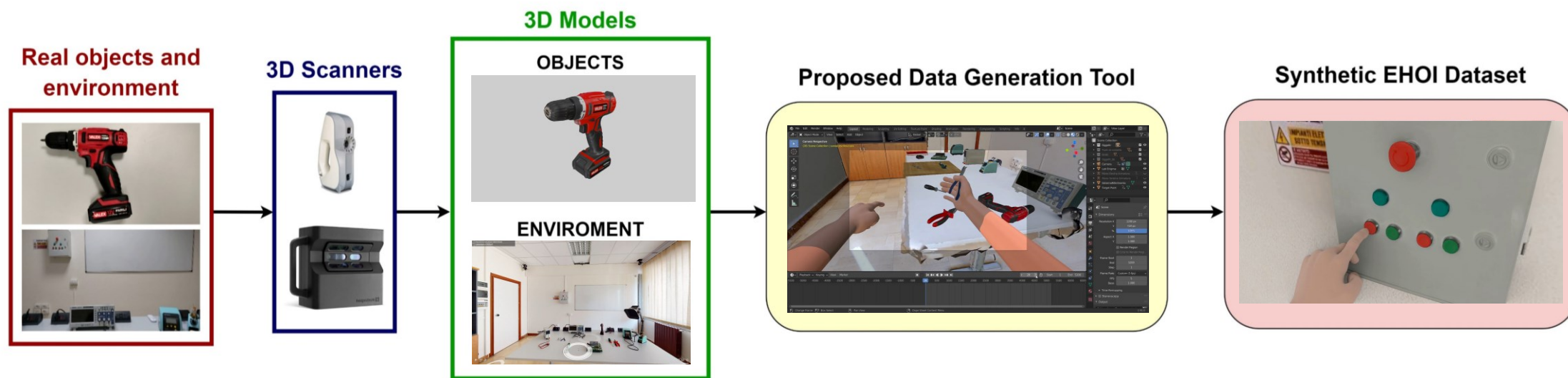


<take>

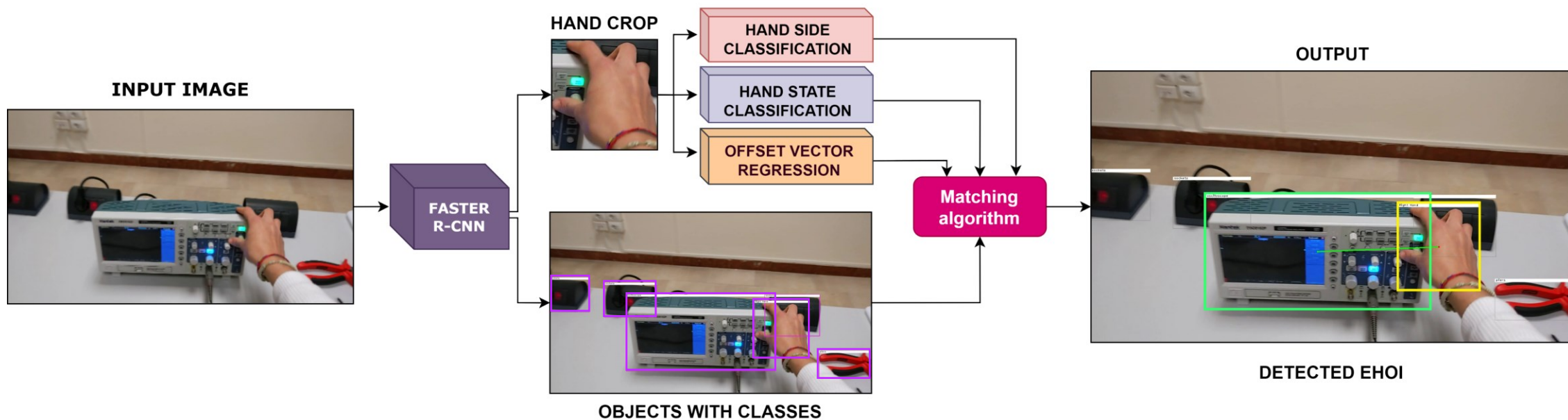


<white_bar>

Can simulated data help?



Can simulated data help?

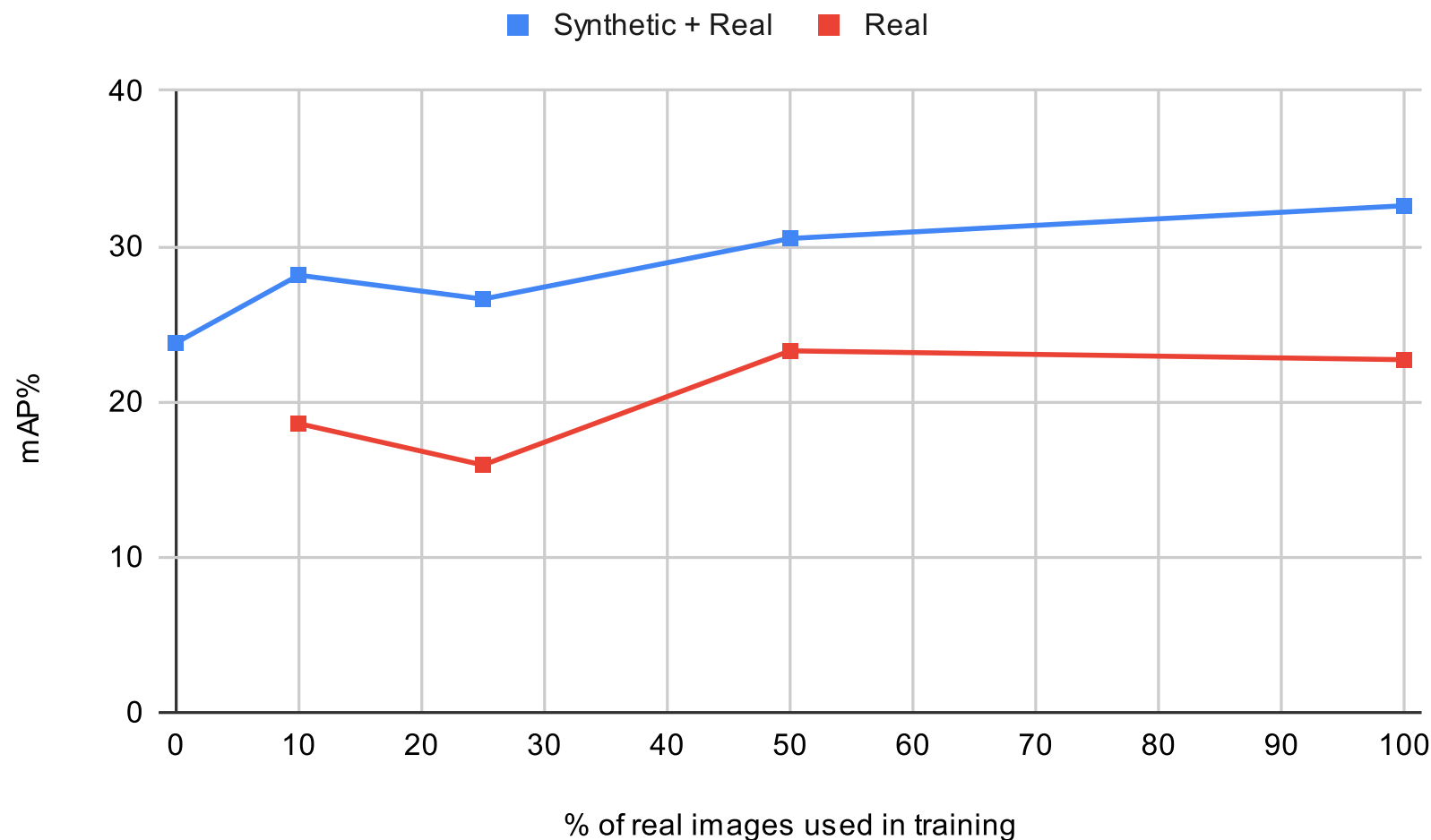


DATA HERE -> https://iplab.dmi.unict.it/EHOI_SYNTH/

IN THIS CONFERENCE

Can simulated data help?

ID 221: Poster
Tomorrow 15.00 – 16.00



Understanding Actions

- Recognizing and detecting the actions performed by user allows to understand what happens in the the video;
- This can be useful to:
 - Segment the video into coherent temporal units for:
 - Summarization;
 - Video understanding;
 - Understand the user's goals to assist them;

Relation between Action and Interaction

TAKE SCREWDRIVER



Relation between Action and Interaction

TAKE SCREWDRIVER



Start Action

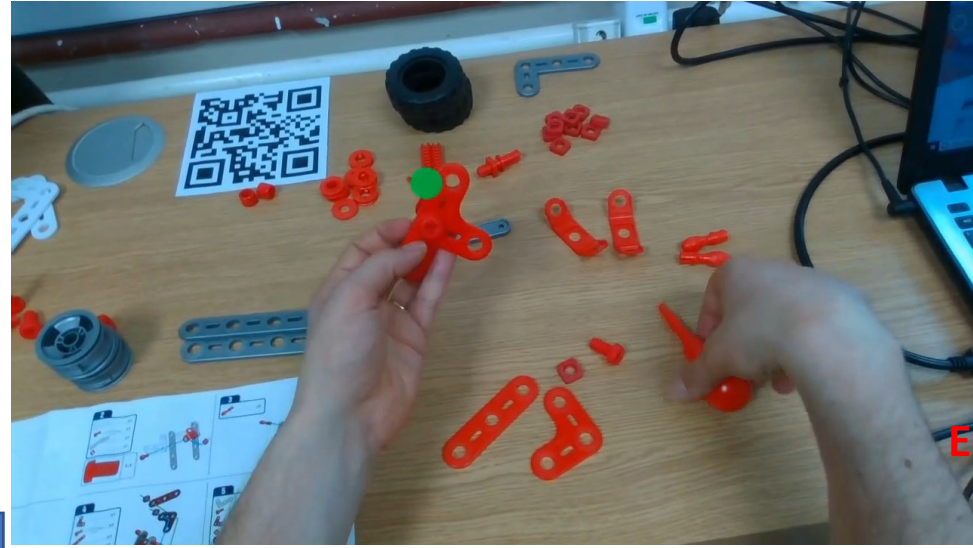
Start Interaction (H-O)



Frame of Contact

Relation between Action and Interaction

TAKE SCREWDRIVER

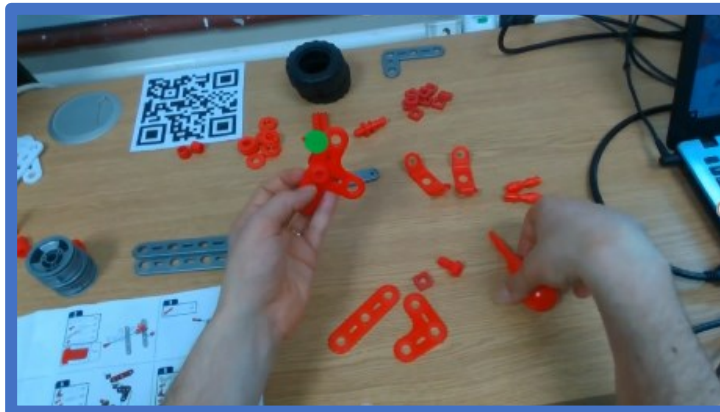


End Interaction

Start Action

Start Interaction (H-O)

End Action












Frame of Contact



Frame of Decontact

Relation between Action and Interaction

Relation	Verbs	MECCANO verbs
	pat, hit, kick	//
	pick up	take, fit, align, plug, pull
	close, open, turn on, press, push	browse
	walk, jump, run	//
	wring out, wash, cut, mix	pull
	throw, leave, place	put
	move	browse
	twist, rip	screw, unscrew, tighten, loosen
	stretch, knead, write, watch	check



Model

VERB

NOUN



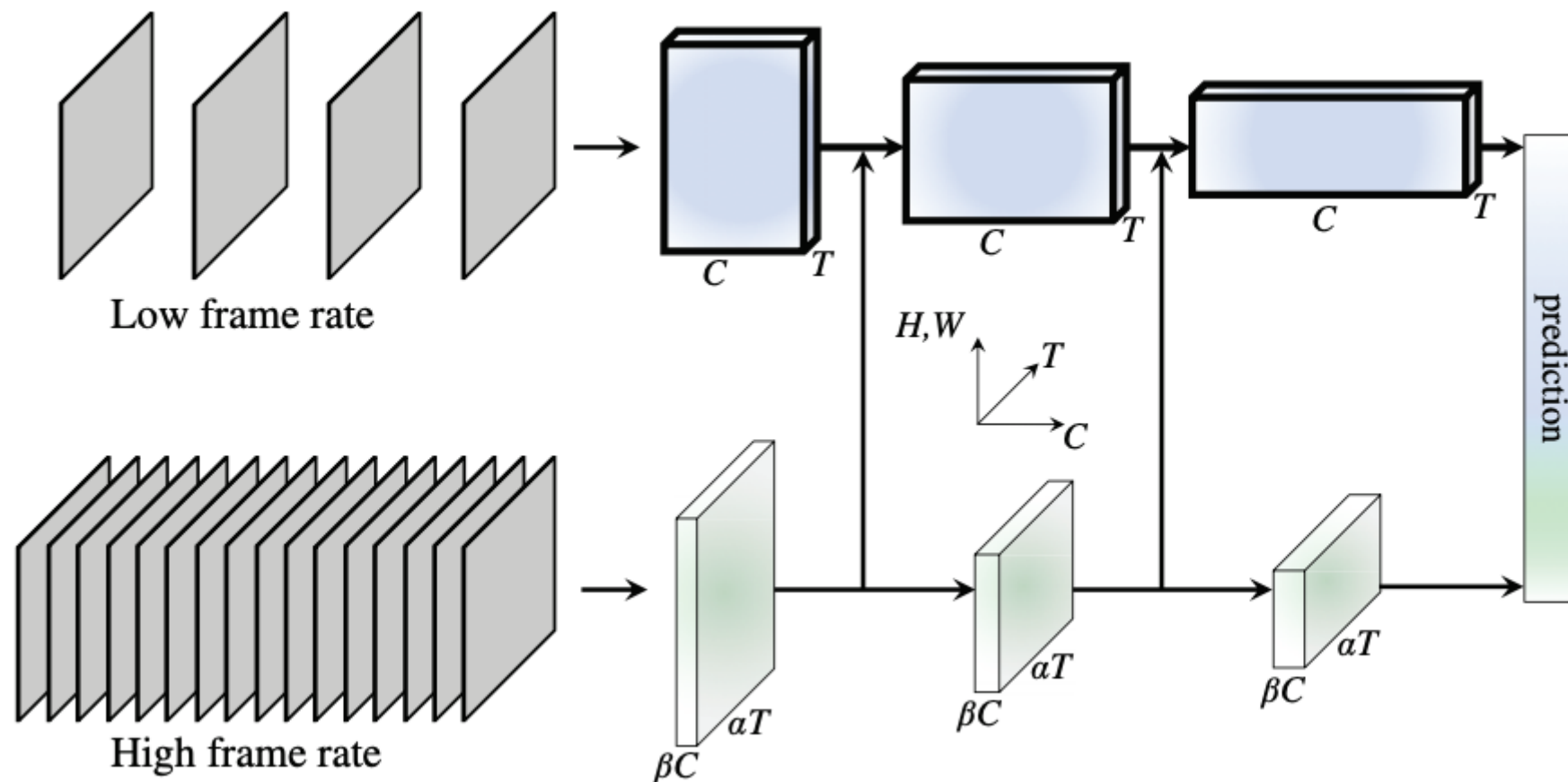
Open - Box
 $v = 3$ $n = 23$



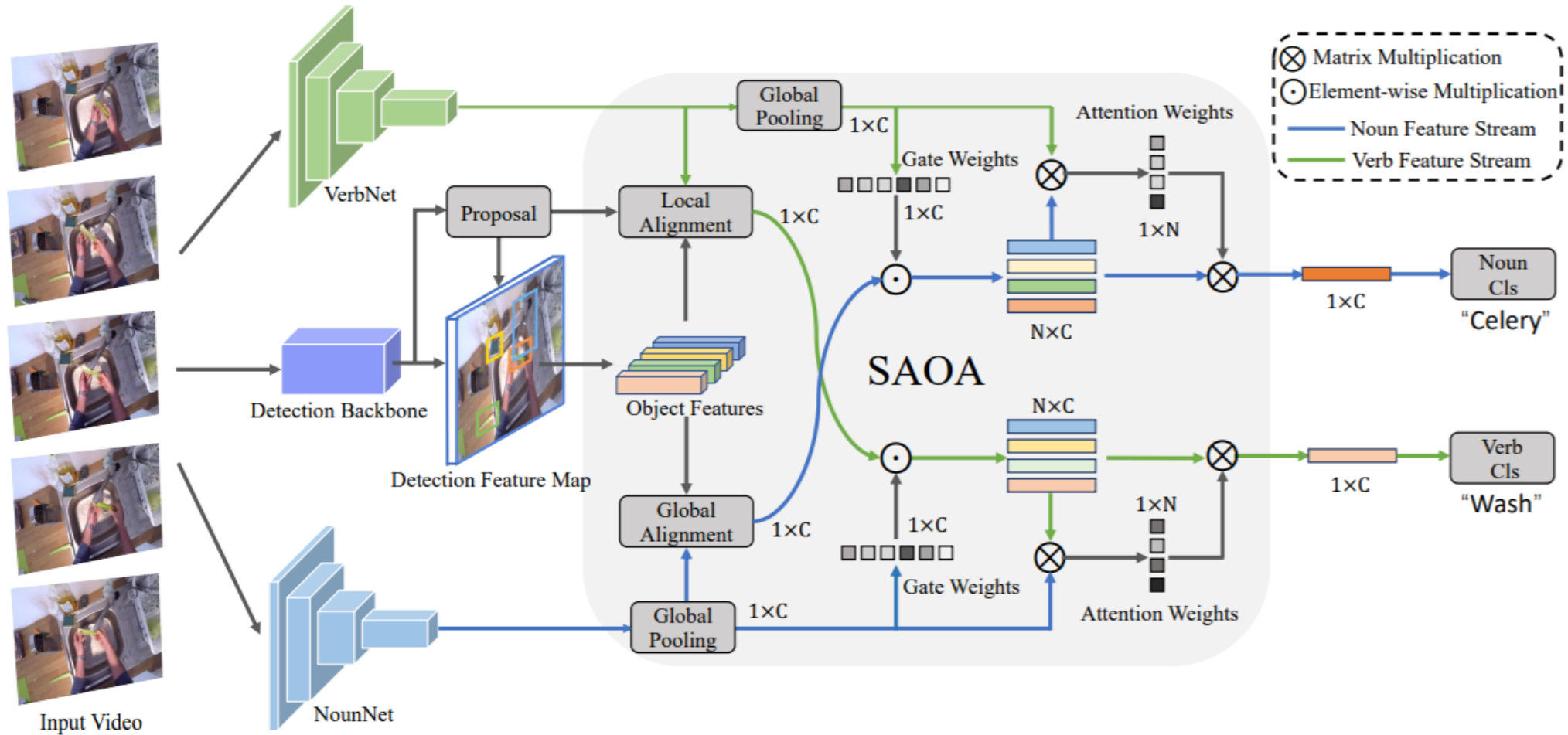
"observe a trimmed segment denoted by start and end time and classify the action present in the clip"

As defined in EPIC-KITCHENS-2020

SlowFast Networks for Video Recognition



Object-centric Egocentric Action Recognition



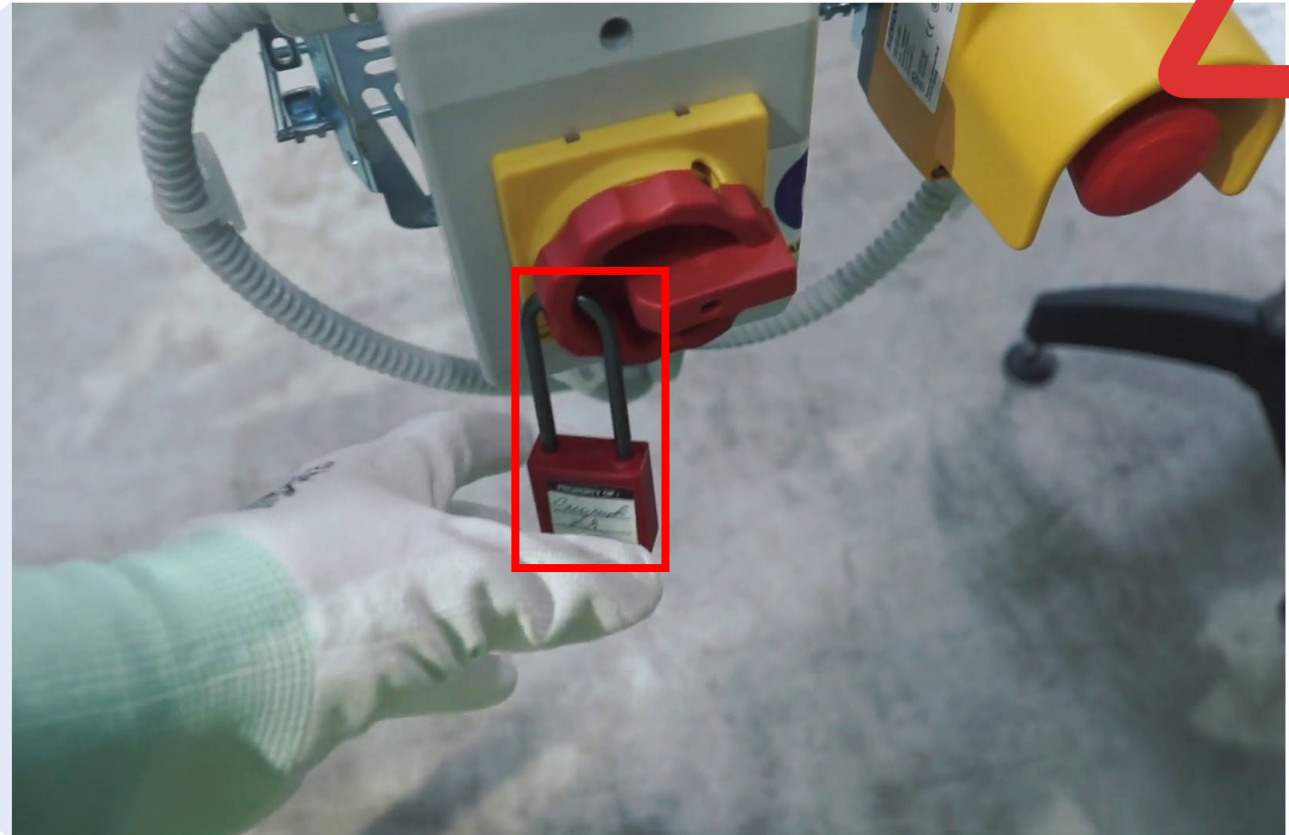
Wang, X., Zhu, L., Wu, Y., & Yang, Y. (2020). Symbiotic attention for egocentric action recognition with object-centric alignment. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.

Personal assistants and Future Predictions

Intelligent assistants should be able to understand what are the user's goals and what is going to happen in the future.

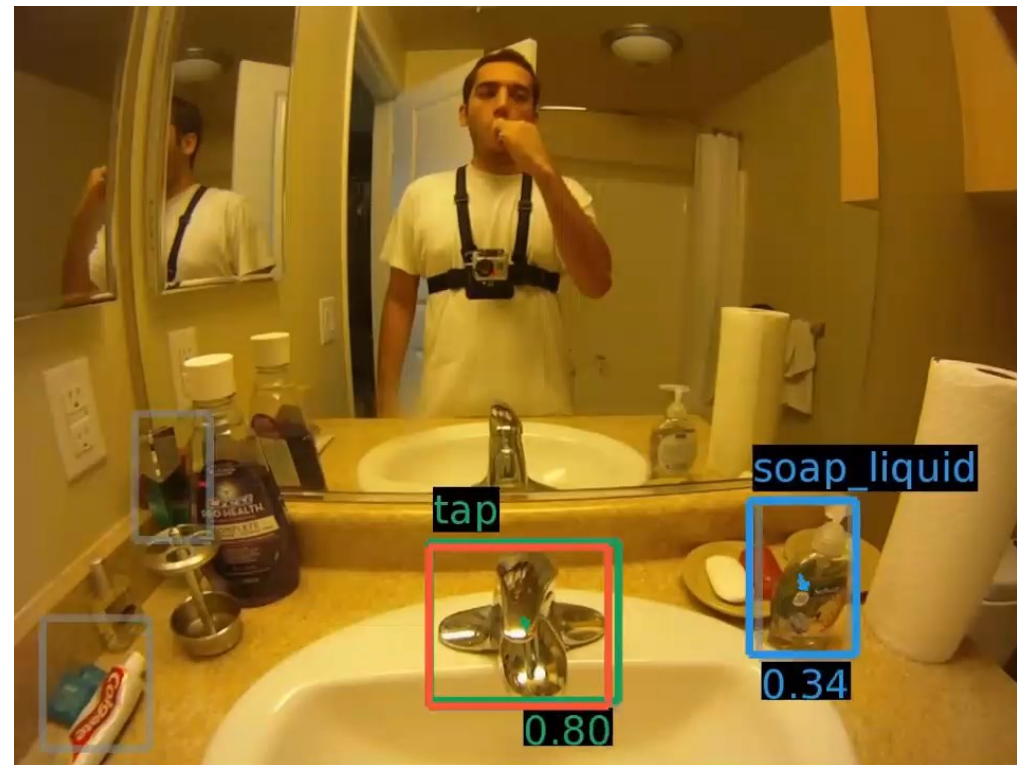
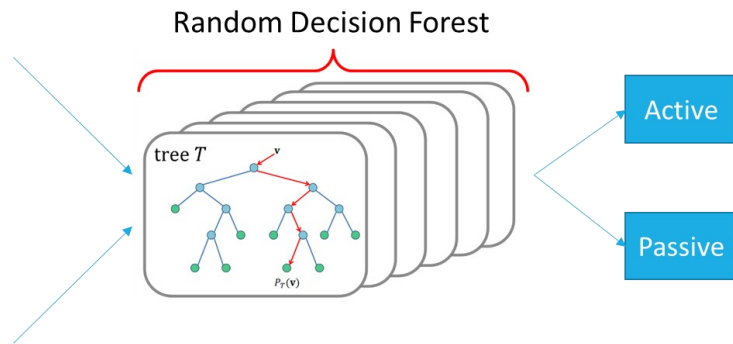
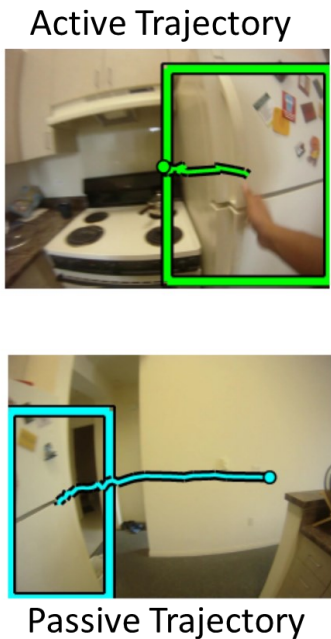
Next-active-object: **LOCKER**

Next action: **OPEN LOCKER**



Anticipation – Next-Active-Objects

Use egocentric object trajectories to distinguish passive from next-active-objects (i.e., those which will be used soon by the user).



THE UNIVERSITY OF TEXAS AT AUSTIN
IMAGE PROCESSING LABORATORY
Next Active Object Prediction from Egocentric Videos
<http://iplab.dmi.unict.it/NextActiveObjectPrediction/>

SUCCESS EXAMPLES

object class
positive predictions (score>0.5)

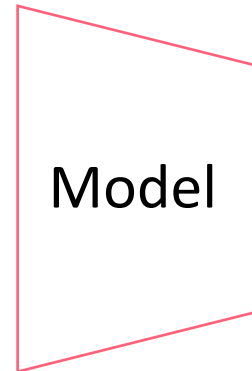
object class
negative predictions (score<=0.5)

discarded objects

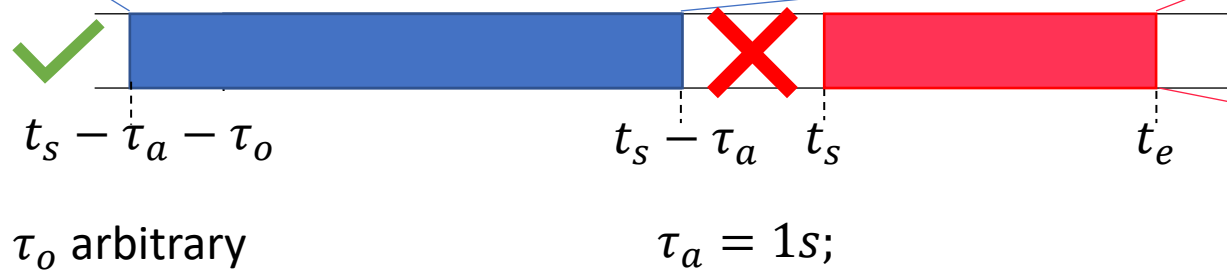
gt next active object

Action Anticipation Task - Definition

(observed video)



Take - Plate



(unobserved)

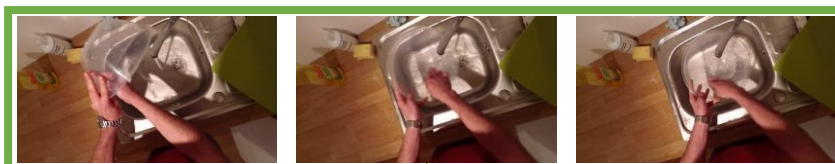
Action Anticipation

MULTILABEL PROBLEM
SINGLE LABEL DATASET

=

MULTILABEL CLASSIFICATION
WITH MISSING LABELS

observed segment



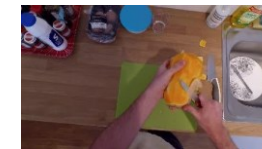
same video

HOW SHOULD WE EVALUATE?
WHAT IS A GOOD ARCHITECTURE?
WHAT IS A GOOD LEARNING OBJECTIVE?

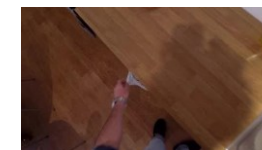
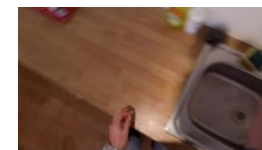
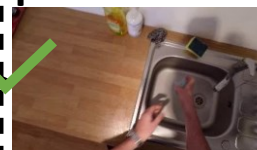
1 sec

Future Actions

peel squash



put down rag



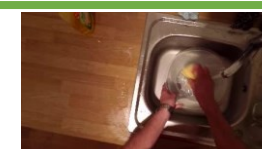
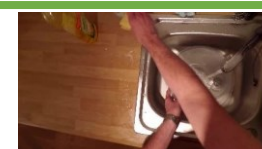
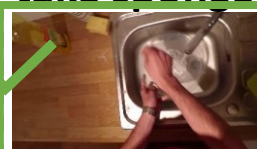
put down sponge



wash bowl



take sponge



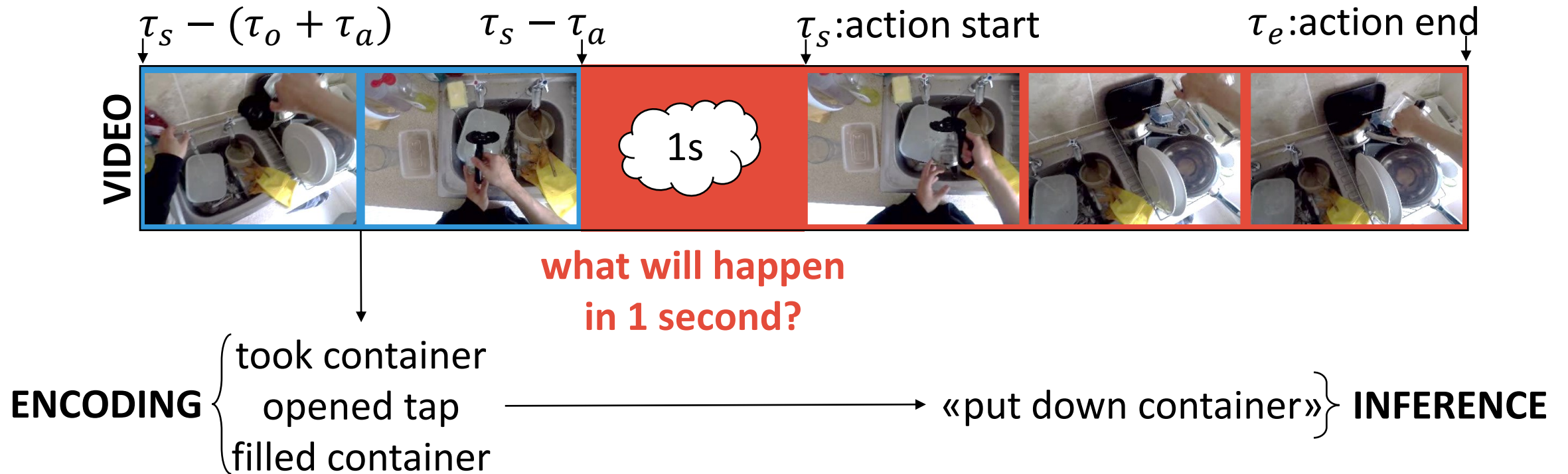
Mean Top-5 Recall

- Since multiple futures are possible, TOP-1 accuracy is not a good measure;
- TOP-5 Accuracy still suffers from a major problem due to the unbalanced nature of the dataset:
 - A high Top-5 accuracy is possible by ranking higher the most frequent classes;
- To overcome this issue, we introduced Mean Top-5 Recall:
 - An action anticipation is correct if the ground truth label is among the Top-5 predicted actions:



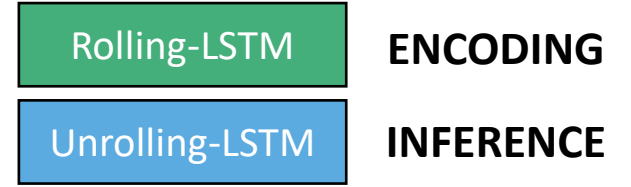
- Results are computed **per-class** then averaged to obtain a single measure.

Action Anticipation: Encoding vs Inference

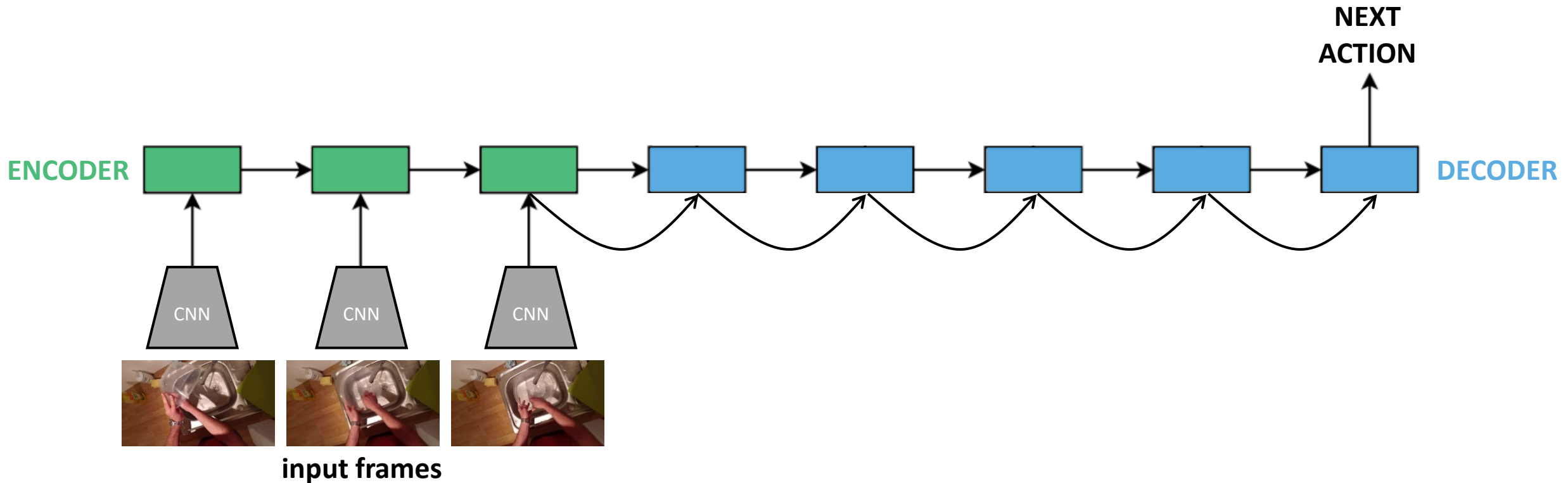


The two tasks are usually performed jointly
 We propose to disentangle them by using two separate LSTMs

Sequence to Sequence Models

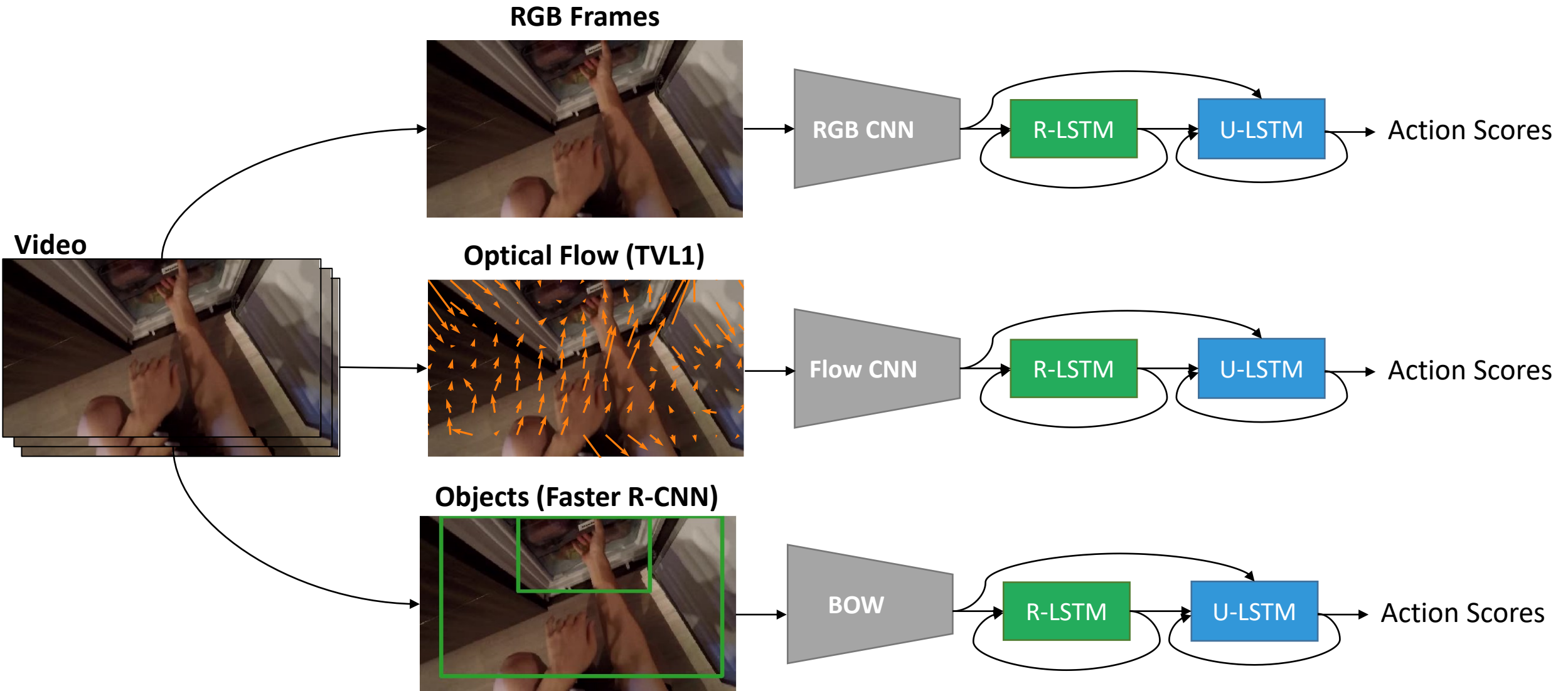


We take inspiration from sequence to sequence models.



Three Modalities

How to fuse?

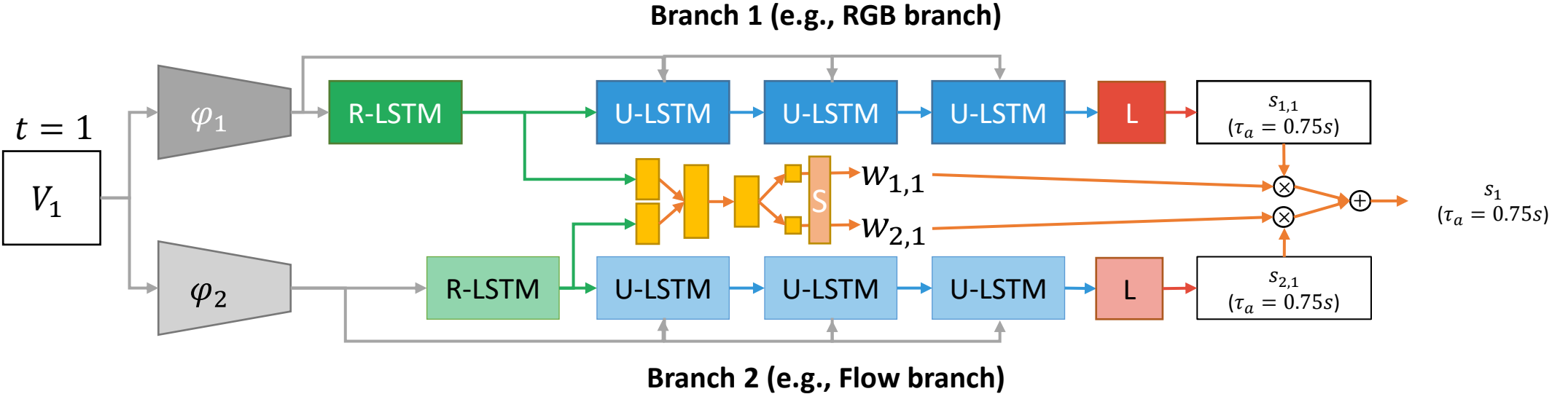


A. Furnari, G. M. Farinella, What Would You Expect? Anticipating Egocentric Actions with Rolling-Unrolling LSTMs and Modality Attention. ICCV 2019 (ORAL).

A. Furnari, G. M. Farinella. Rolling-Unrolling LSTMs for Action Anticipation from First-Person Video. TPAMI 2020. <http://iplab.dmi.unict.it/rulstm>

Modality ATtention (MATT)

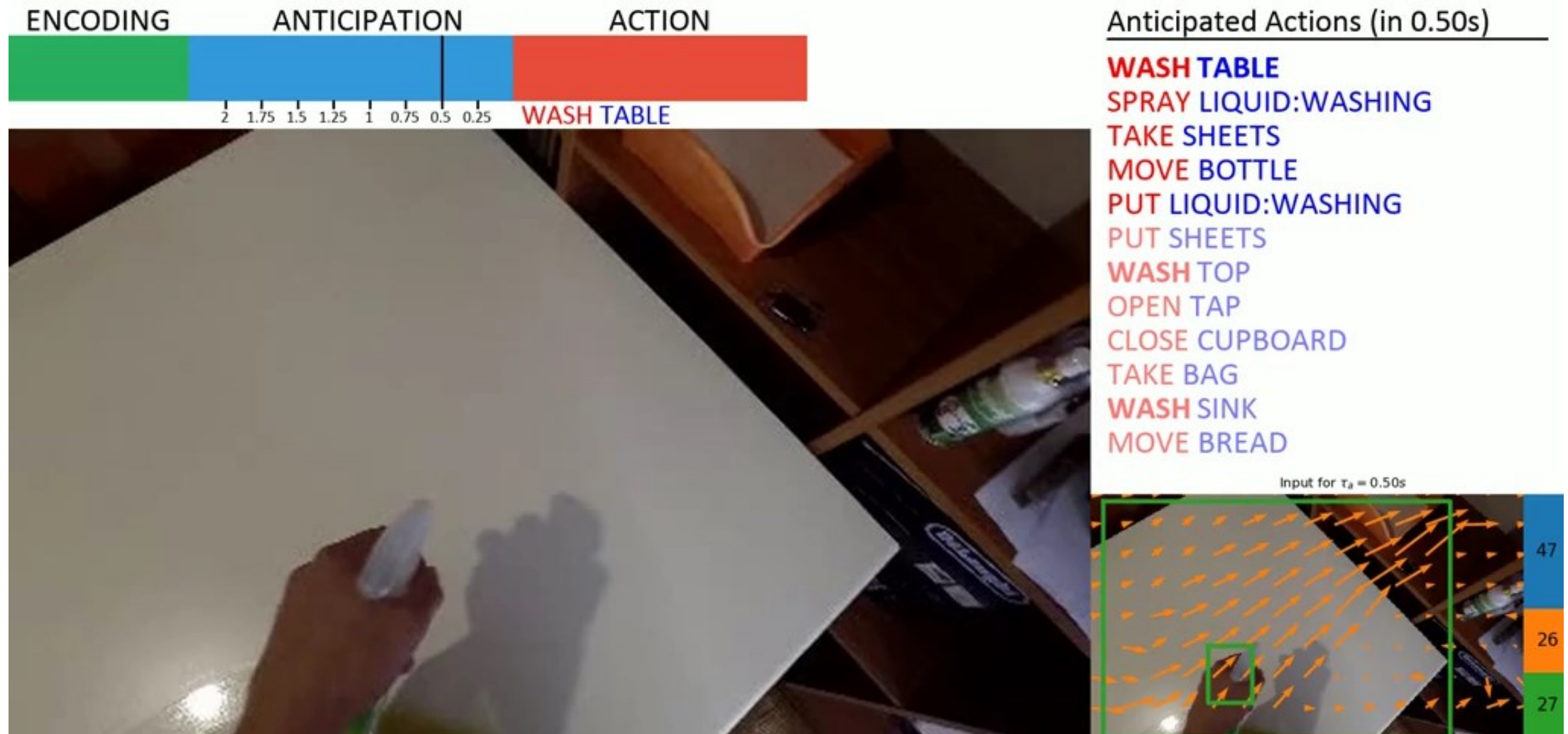
The relative importance of each modality may depend on the observed sample.



V_t Input video snippets L Linear transformation MATT Modality Attention Network (MATT) S SoftMax
 → Message passing

A. Furnari, G. M. Farinella, What Would You Expect? Anticipating Egocentric Actions with Rolling-Unrolling LSTMs and Modality Attention. ICCV 2019 (ORAL).
 A. Furnari, G. M. Farinella. Rolling-Unrolling LSTMs for Action Anticipation from First-Person Video. TPAMI 2020. <http://iplab.dmi.unict.it/rulstm>

Demo Video: Egocentric Action Anticipation



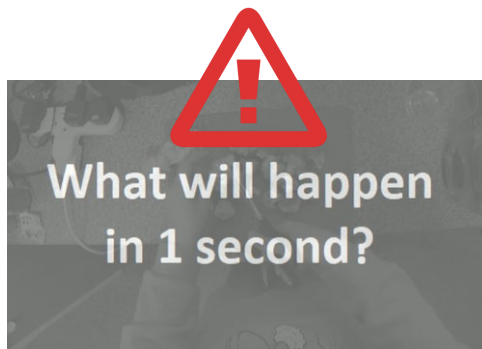
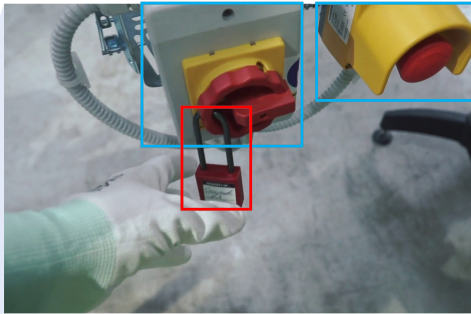
A. Furnari, G. M. Farinella, What Would You Expect? Anticipating Egocentric Actions with Rolling-Unrolling LSTMs and Modality Attention. ICCV 2019 (ORAL).

A. Furnari, G. M. Farinella. Rolling-Unrolling LSTMs for Action Anticipation from First-Person Video. TPAMI 2020. <http://iplab.dmi.unict.it/rulstm>

Can we bring egocentric vision to industry?

Next-active-object: **LOCKER**

Next action: **OPEN LOCKER**



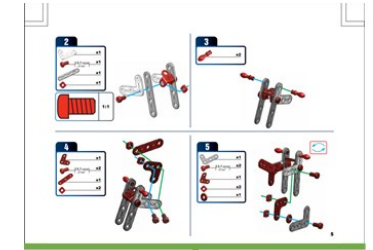
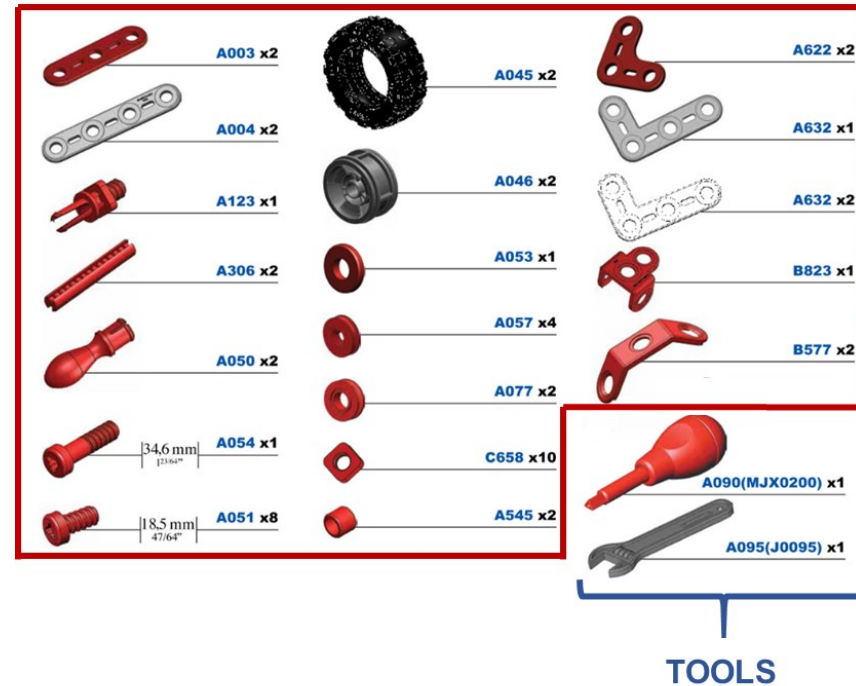
- The factory is a natural place for a wearable assistant;
- Closed-world assumption;
- Current research has considered different scenarios;
- No datasets in industrial-like scenarios;

The MECCANO Dataset

Data HERE -> <https://iplab.dmi.unict.it/MECCANO/>

We asked subjects to record egocentric videos while assembling a toy motorbike.

The assembly required to interact with several parts and two tools.



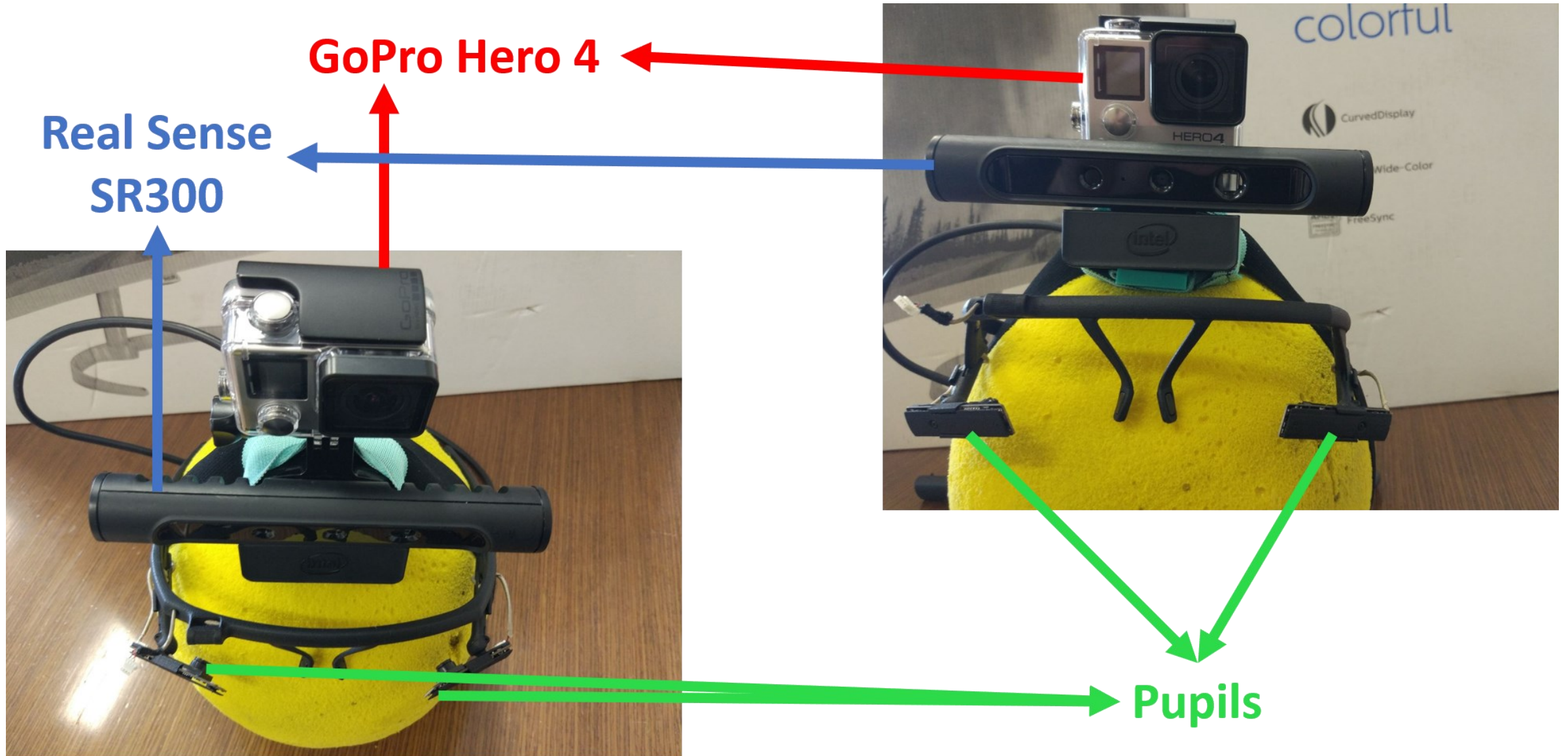
BOOKLET

COMPONENTS



The scenario is industrial-like, with subjects undertaking interactions with tiny objects and tools in a sequential fashion to reach a goal.

Data Collection



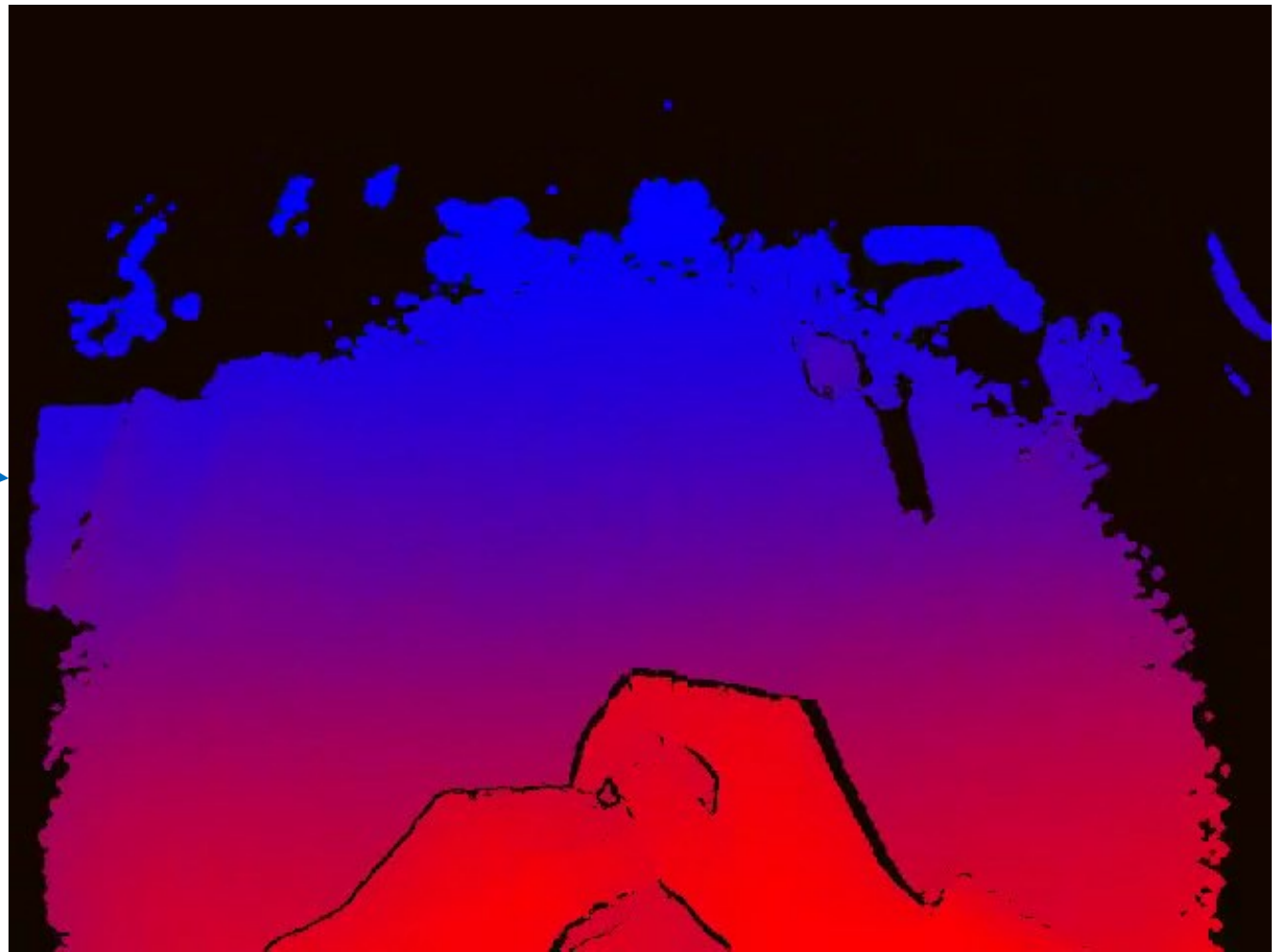
The MECCANO Dataset

RGB



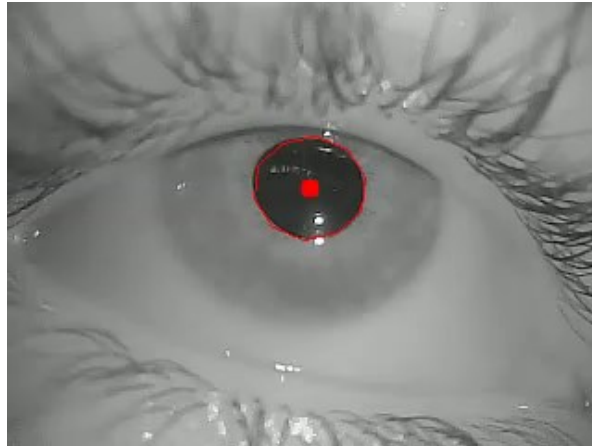
The MECCANO Dataset

Depth



The MECCANO Dataset

Gaze



The MECCANO Dataset: Statistics



20 Subjects



3 Modalities



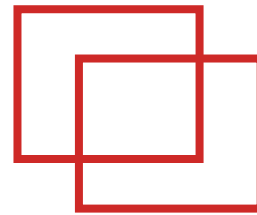
20 min. avg. Video length



5 Tasks



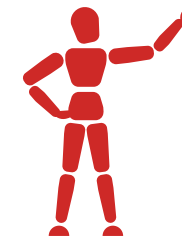
8858 Segments



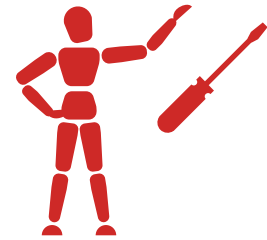
64349 Boxes



20 Objects



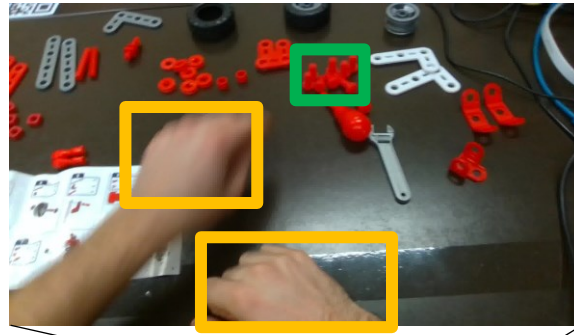
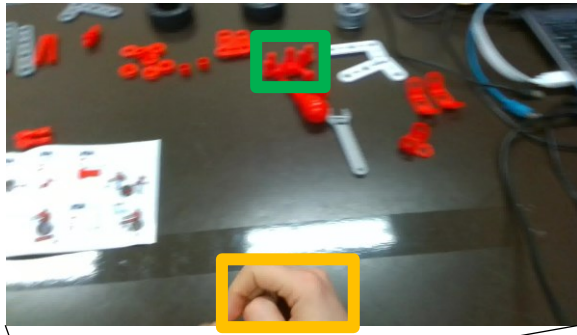
12 Verbs



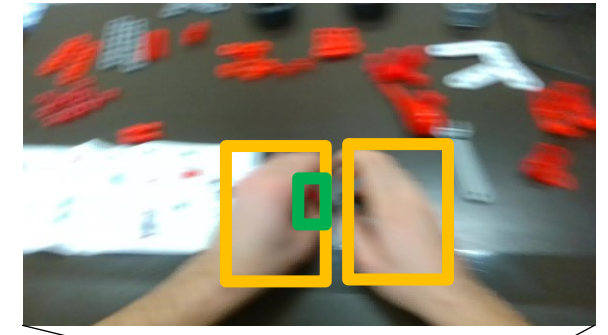
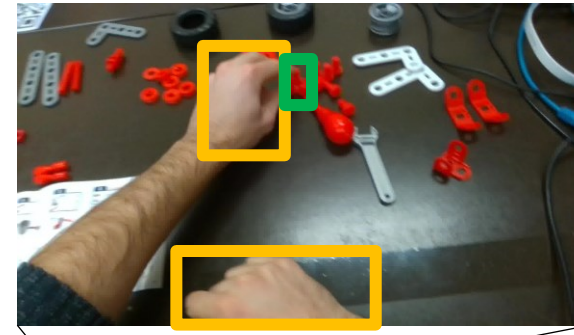
61 Actions

The MECCANO Dataset: Hands and Future Objects

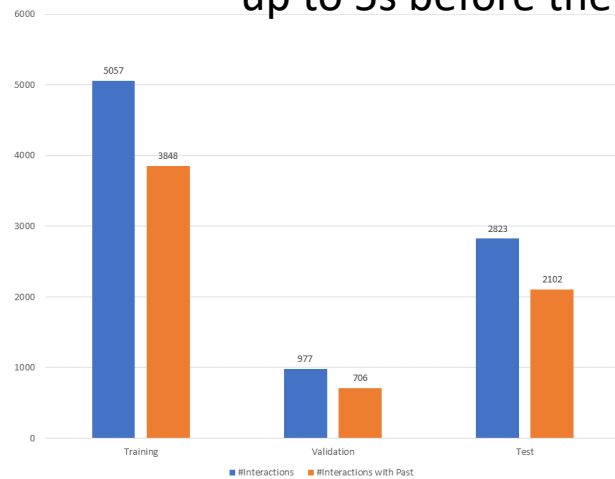
Hands + Next Active Objects



Action Boundaries + Active Objects + Hands

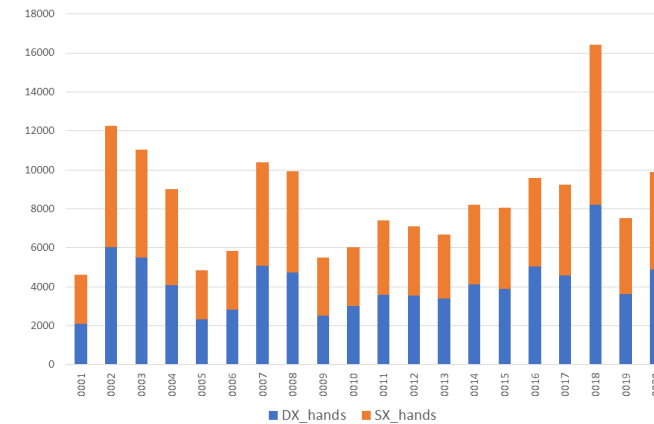


up to 5s before the interaction



Video	Interactions	Interactions with past
0001	319	257
0002	586	452
0003	573	429
0004	485	372
0005	251	200
0006	307	234
0007	493	367
0008	550	384
0009	289	289
0010	304	194
0011	400	310
0012	384	258
0013	313	244
0014	434	297
0015	425	324
0016	576	436
0017	484	339
0018	788	603
0019	400	294
0020	496	373
Total	8857	6656

Human-Object Interaction



The MECCANO Dataset: Tasks

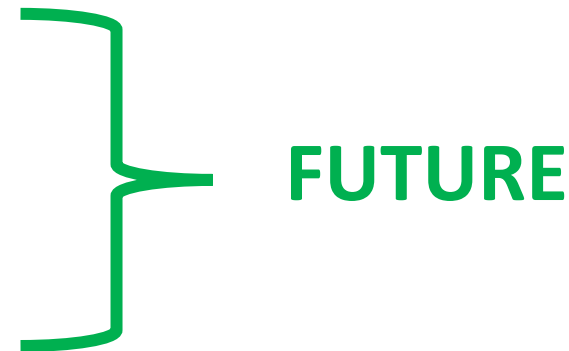
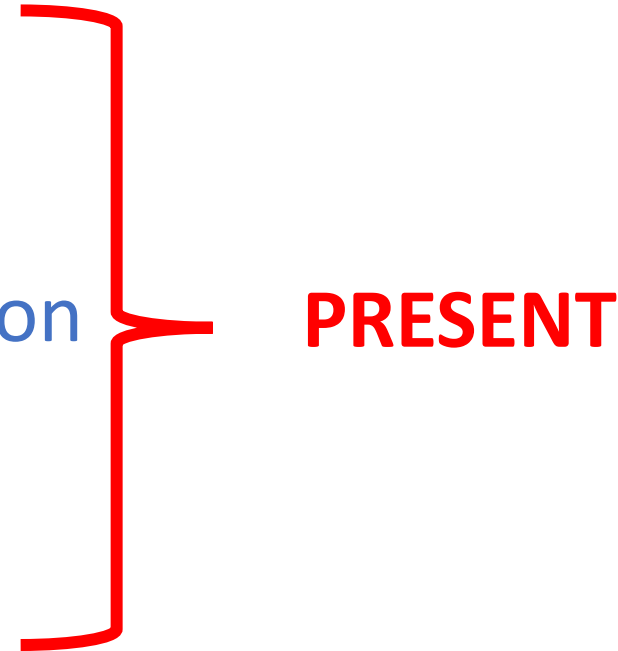
1) Action Recognition

2) Active Object Detection and Recognition

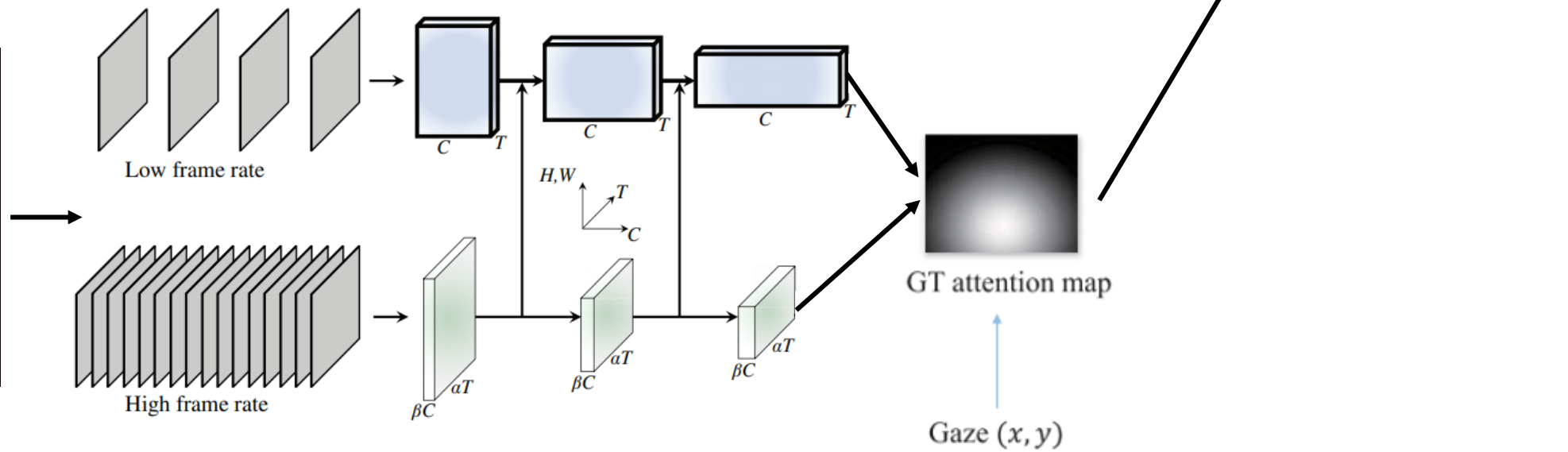
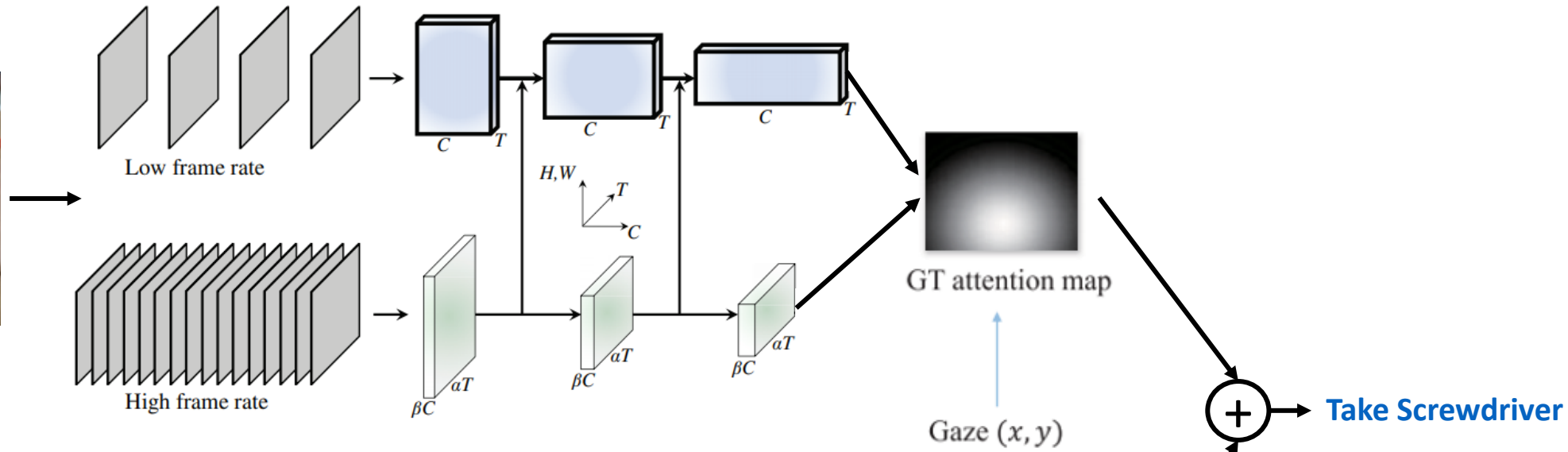
3) EHOI Interaction Detection

4) Action Anticipation

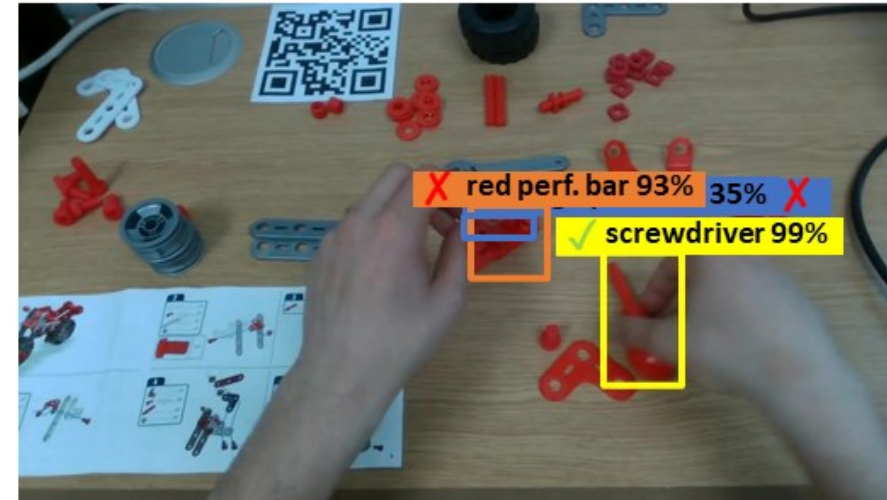
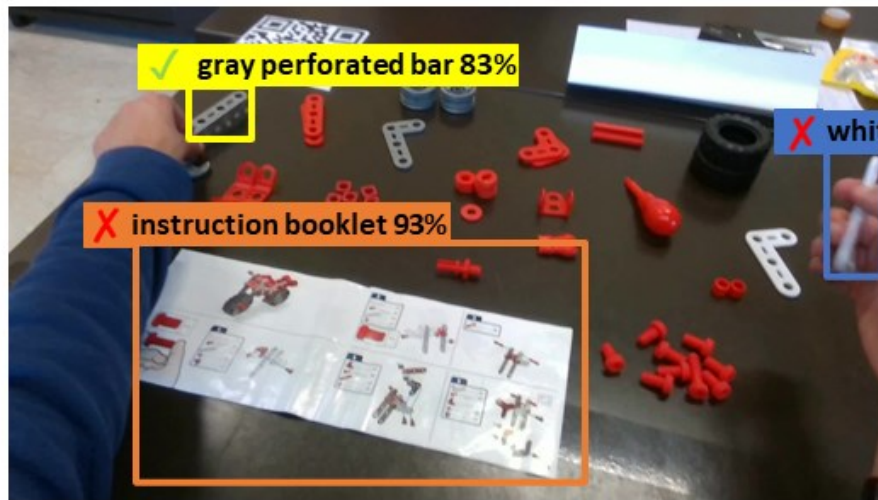
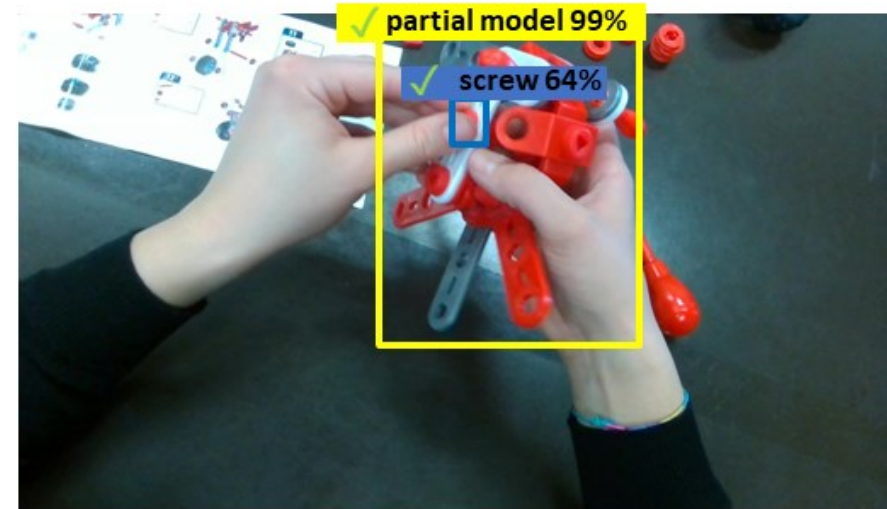
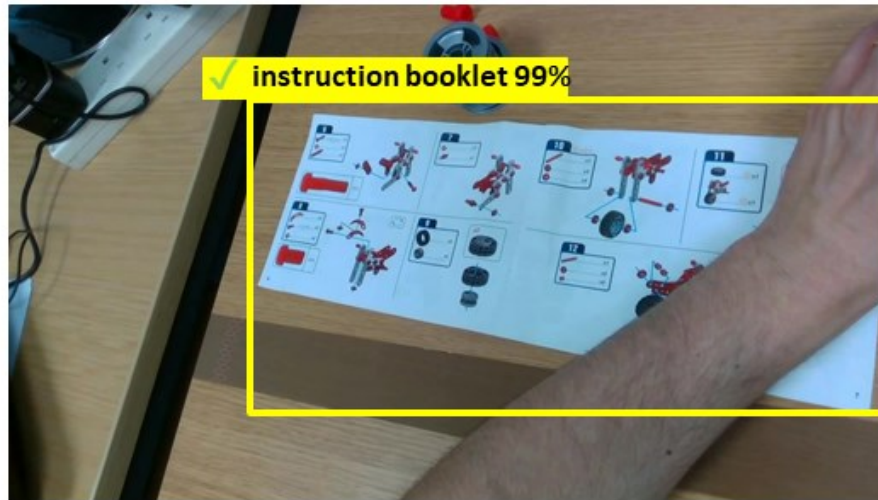
5) Next-Active Object Detection



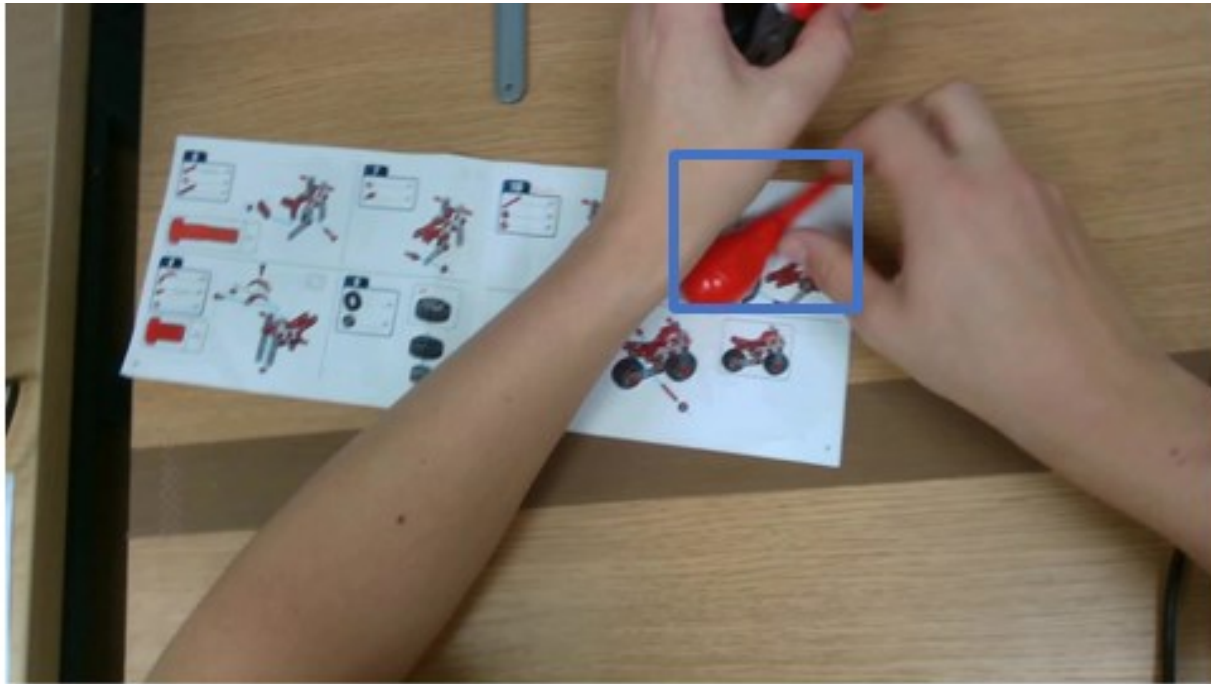
Action Recognition



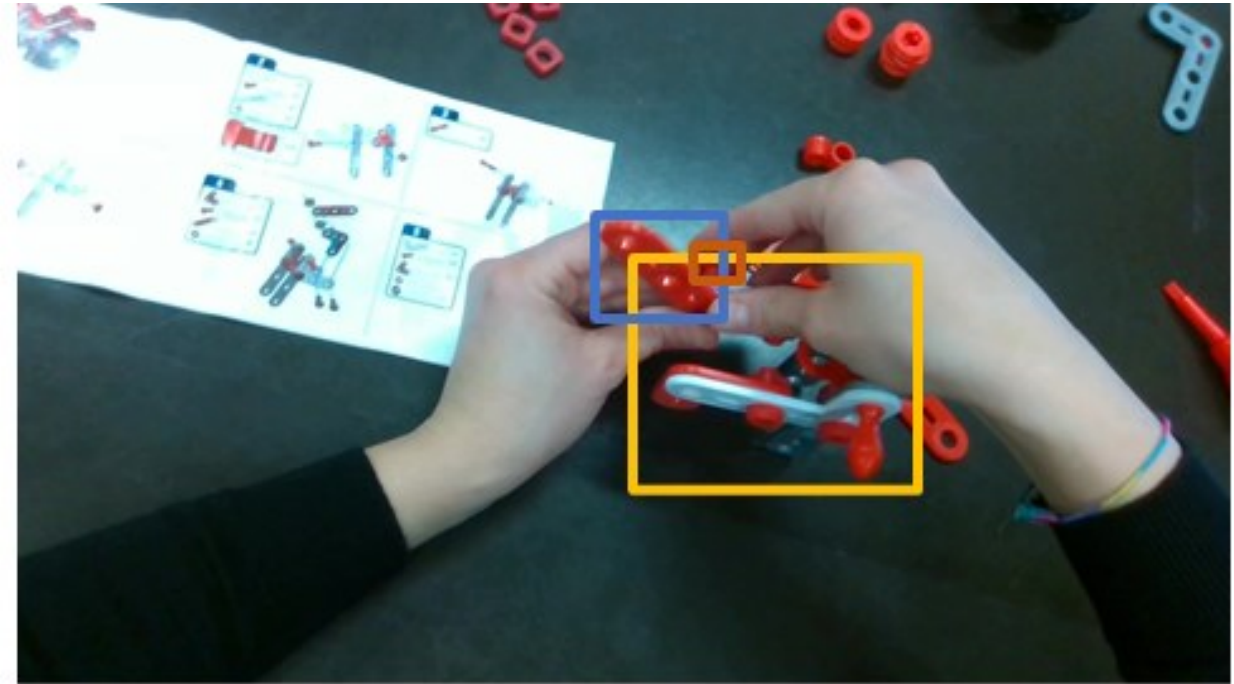
Active Object Detection and Recognition



EHOI Detection

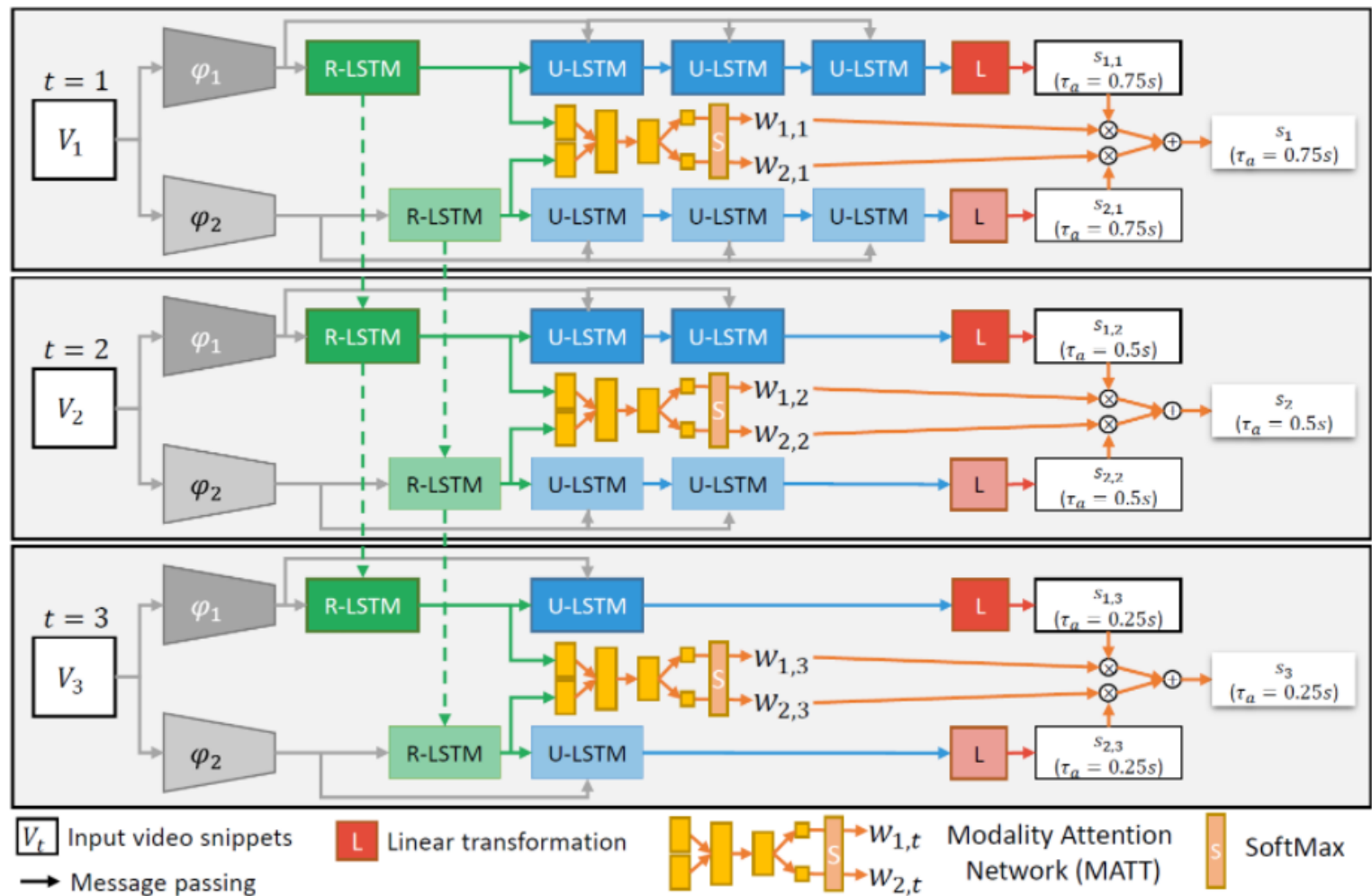


<put, screwdriver>



<plug, {red_perforated_bar, screw, partial_model}>

Action Anticipation



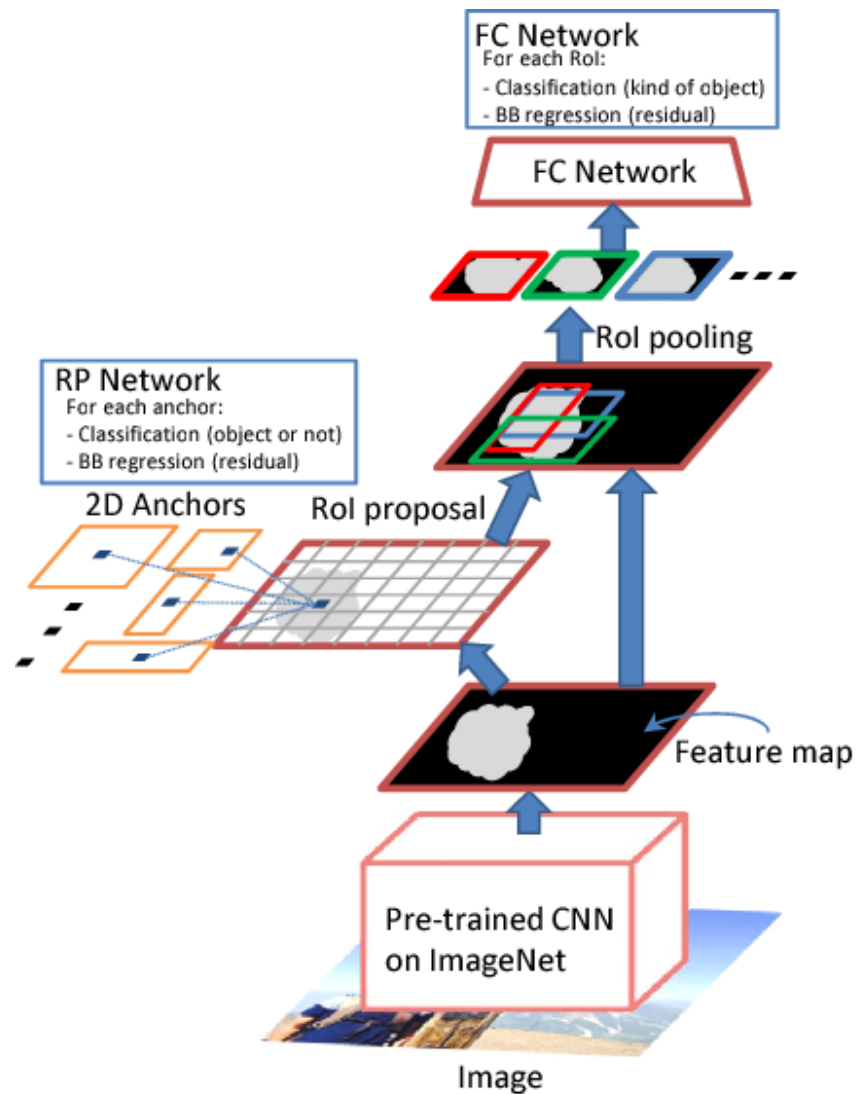
Modalities:

- RGB
- Optical Flow
- Objects

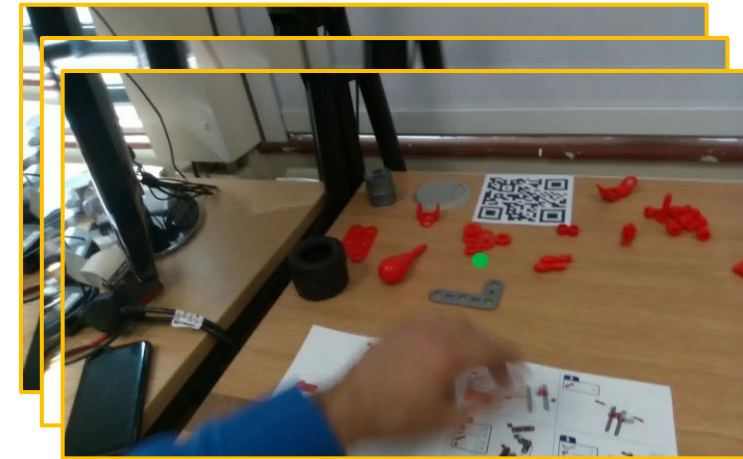
Our Modalities:

- RGB + Flow
- Depth
- Objects
- Hands
- Gaze

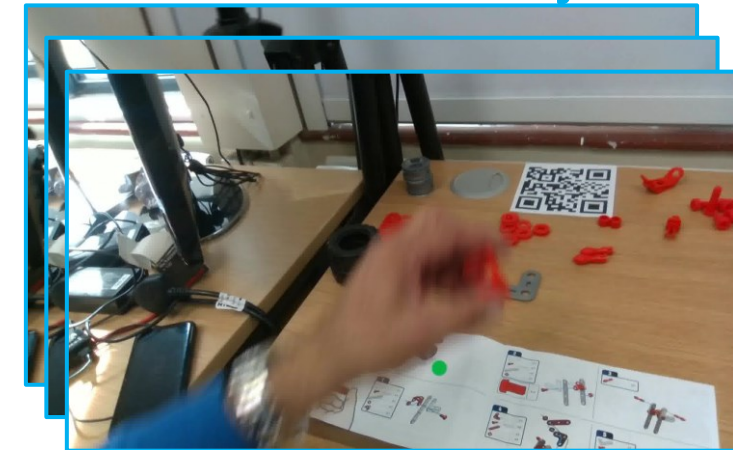
Next-Active Objects Detection



Active Objects



Next-Active Objects



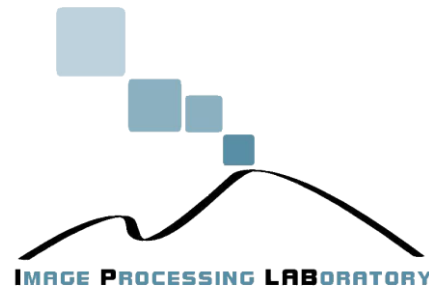
NEXT VISION N

Spin-off of the University of Catania

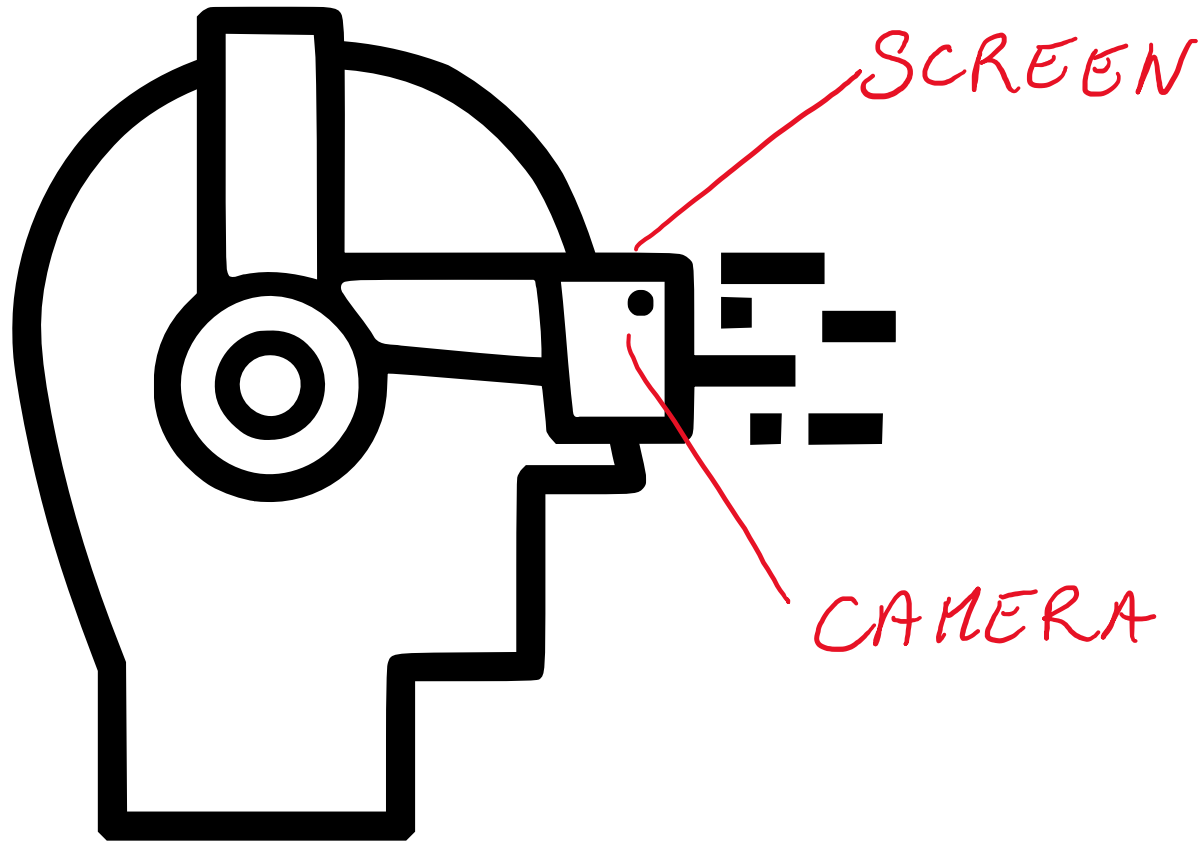
<https://www.nextvisionlab.it/>



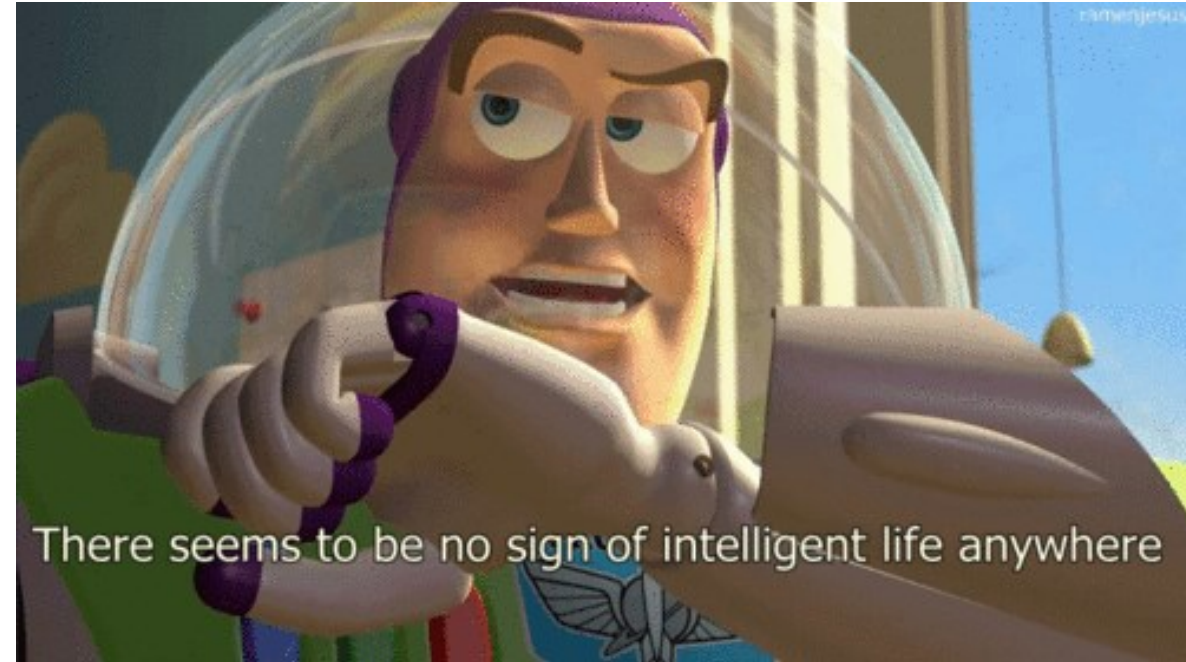
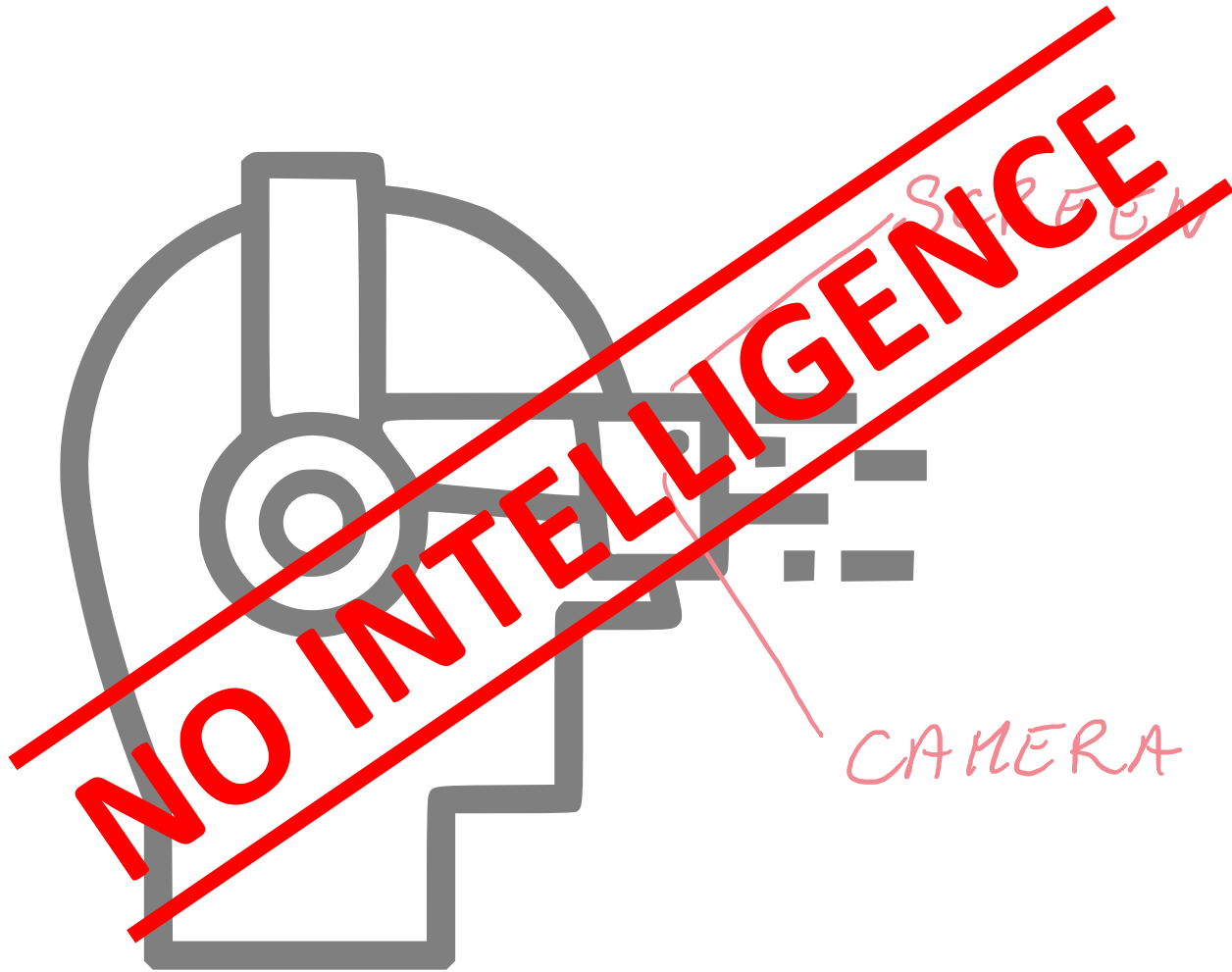
UNIVERSITÀ
degli STUDI
di CATANIA



Wearable Devices



Wearable Devices



Artificial Vision for Human Safety Prevention

Mixed Reality for Guidance and Enhanced Training on Wearable Glasses

Detection of Active Objects

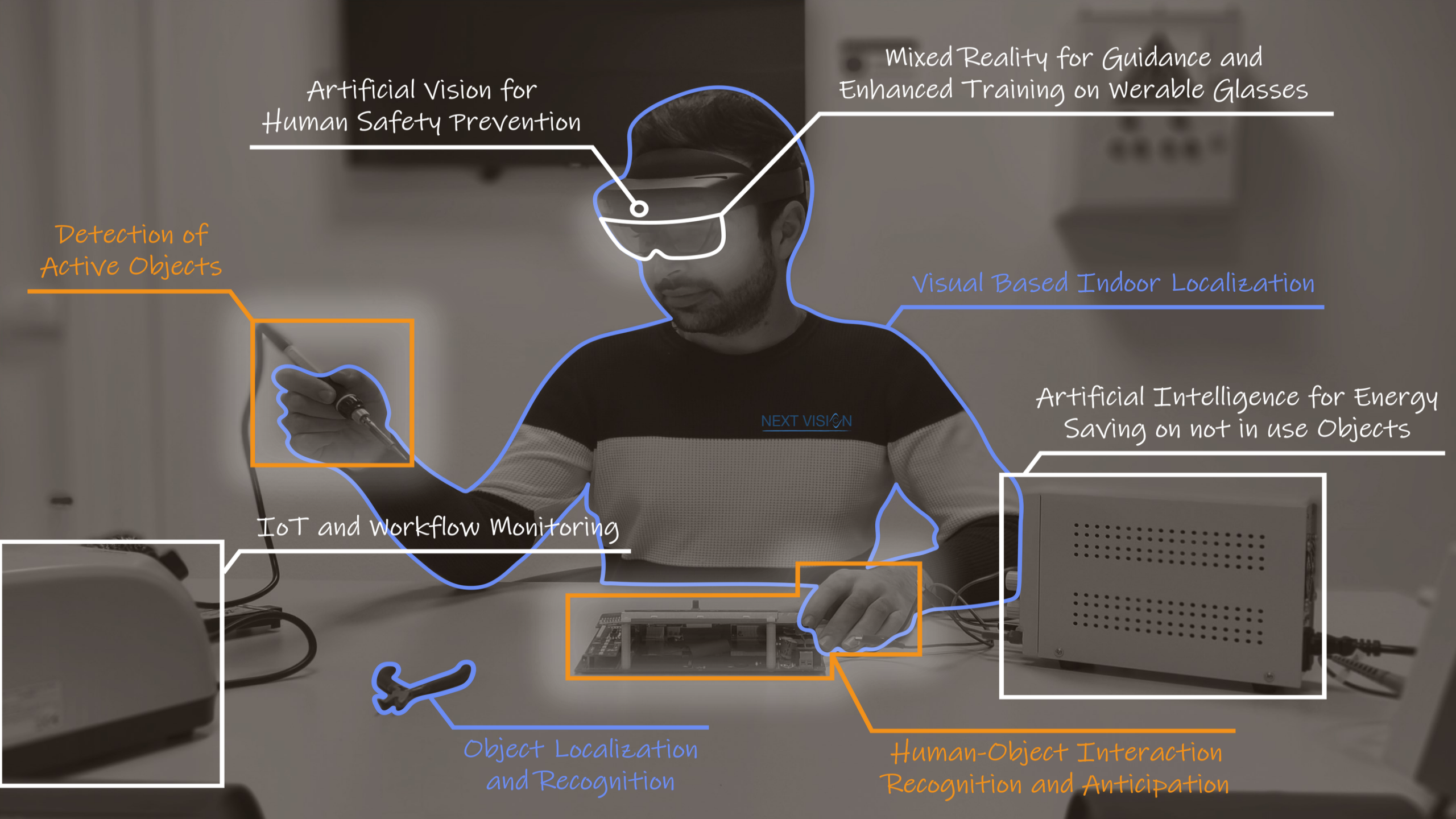
Visual Based Indoor Localization

Artificial Intelligence for Energy Saving on not in use Objects

IoT and Workflow Monitoring

Object Localization and Recognition

Human-Object Interaction Recognition and Anticipation



Solutions: NAIROBI

Solutions: Human-Object Interaction Recognition

Conclusion

- First Person Vision paves the way to a variety of user-centric applications;
- However, we are still missing solid building blocks related to fundamental problems of First Person Vision such as action recognition, object detection, action anticipation and human-object interaction detection;
- Consumer devices are starting to appear, but the near future of First Person Vision is in focused applications such as the ones in industrial scenarios.

Look for us

- **25 May 15:00-16:00 Poster Session**

- Egocentric Human-Object Interaction Detection Exploiting Synthetic Data
- Weakly Supervised Attended Object Detection Using Gaze Data as Annotations
- Panoptic Segmentation in Industrial Environments using Synthetic and Real Data

- **26 May 11:45 - 12:00 Oral session**

- Unsupervised Multi-Camera Domain Adaptation for Object Detection

- **26 May 15:30 - 16:30 Poster Session**

- Untrimmed Action Anticipation

Before we leave...

The slides of this tutorial are available online at:

<http://www.antoninofurnari.it/talks/iciap2022>



Thank you!



Antonino Furnari



Francesco Ragusa



Università
di Catania

NEXT VISION
Spin-off of the University of Catania

ICIAP 2021



First Person (Egocentric) Vision for Human-Centric Assistance: History, Building Blocks, and Applications

Antonino Furnari, Francesco Ragusa

Image Processing Laboratory - <http://iplab.dmi.unict.it/>

Department of Mathematics and Computer Science - University of Catania

Next Vision s.r.l., Italy

furnari@dmf.unict.it - <http://www.antoninofurnari.it/>

francesco.ragusa@unict.it - <https://iplab.dmi.unict.it/ragusa/>

<http://iplab.dmi.unict.it/fpv> - <https://www.nextvisionlab.it/>