



Università
di Catania

NEXT VISION

Spin-off of the University of Catania



Seeing Through the User's Eyes: Advances in Human-Centric Egocentric Vision

Francesco Ragusa

LIVE Group @ UNICT - <https://iplab.dmi.unict.it/live/>

Next Vision - <http://www.nextvisionlab.it/>

Department of Mathematics and Computer Science - University of Catania

francesco.ragusa@unict.it - <https://francescoragusa.github.io/>



VISAPP 2026

21st International Conference on Computer Vision
Theory and Applications

Marbella, Spain 9 - 11 March, 2026

- 1) Part I: History and motivations [10.30 - 12.00]
 - a) Agenda of the tutorial;
 - b) Perception and Egocentric Vision;
 - c) Seminal works in Egocentric Vision;
 - d) Differences between Third Person and First Person Vision;
 - e) First Person Vision datasets;
 - f) Wearable devices to acquire/process first person visual data;
 - g) Main research trends in First Person (Egocentric) Vision;
 - h) What's next?

Lunch [12.00 – 13.00]

- 2) **Part II: Fundamental tasks for First Person Vision systems [13.00 – 15.00]**
 - a) **Visual Localization;**
 - b) **Hand/Object Detection;**
 - c) **Hand-Object Interaction;**
 - d) **Procedural Understanding;**
 - e) **Actions and Objects anticipation;**
 - f) **Dual-Agent Language Assistance**
 - g) **Industrial Applications**

Part II

Fundamental Tasks for First Person Vision Systems

The slides of this tutorial are available online at:

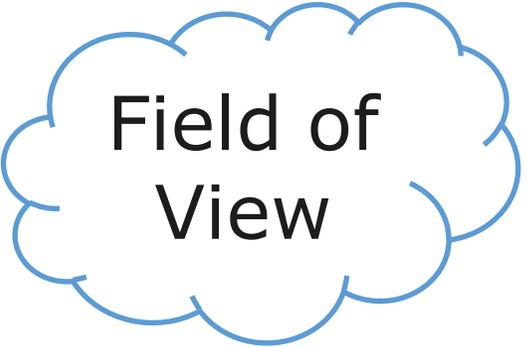
<https://francescoragusa.github.io/visapp2026>



Four things to pay attention to when collecting first person visual data



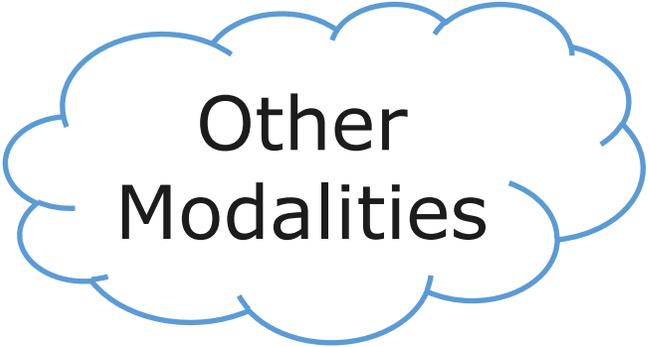
Video
Quality



Field of
View



Wearing
Modality



Other
Modalities

- 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

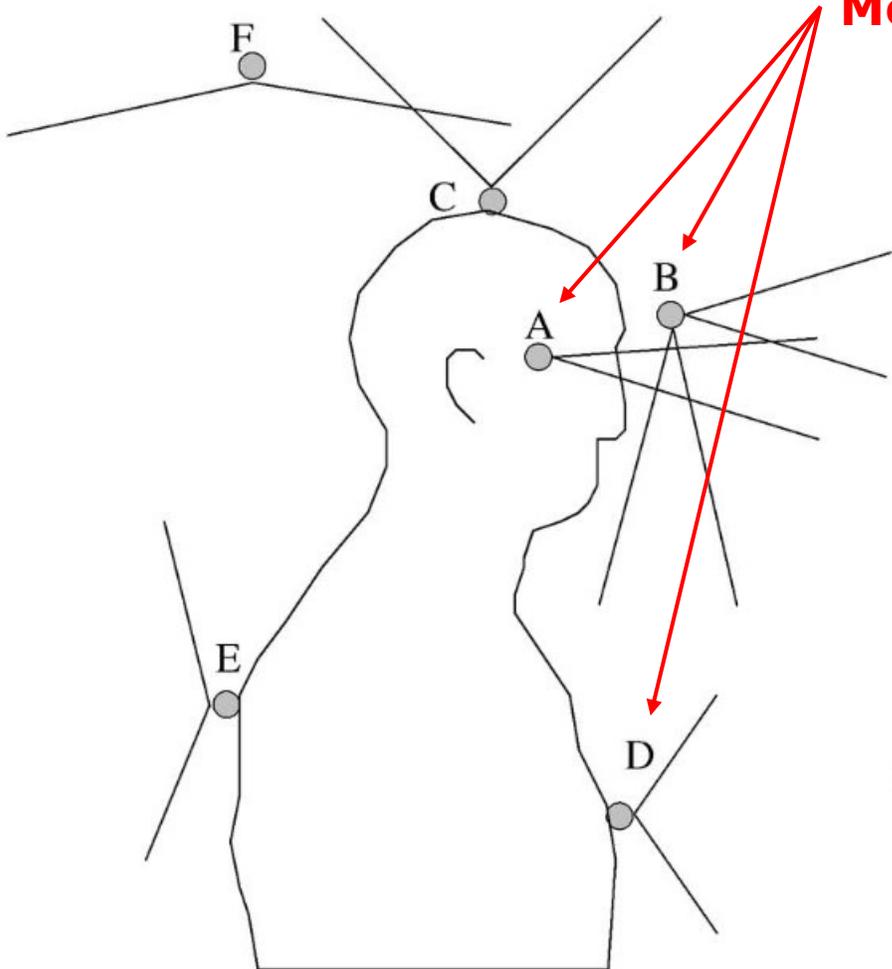


Average Quality Video



A, B: head mounted, D: chest mounted

Most Common Wearing Modalities



A



B (frontward)

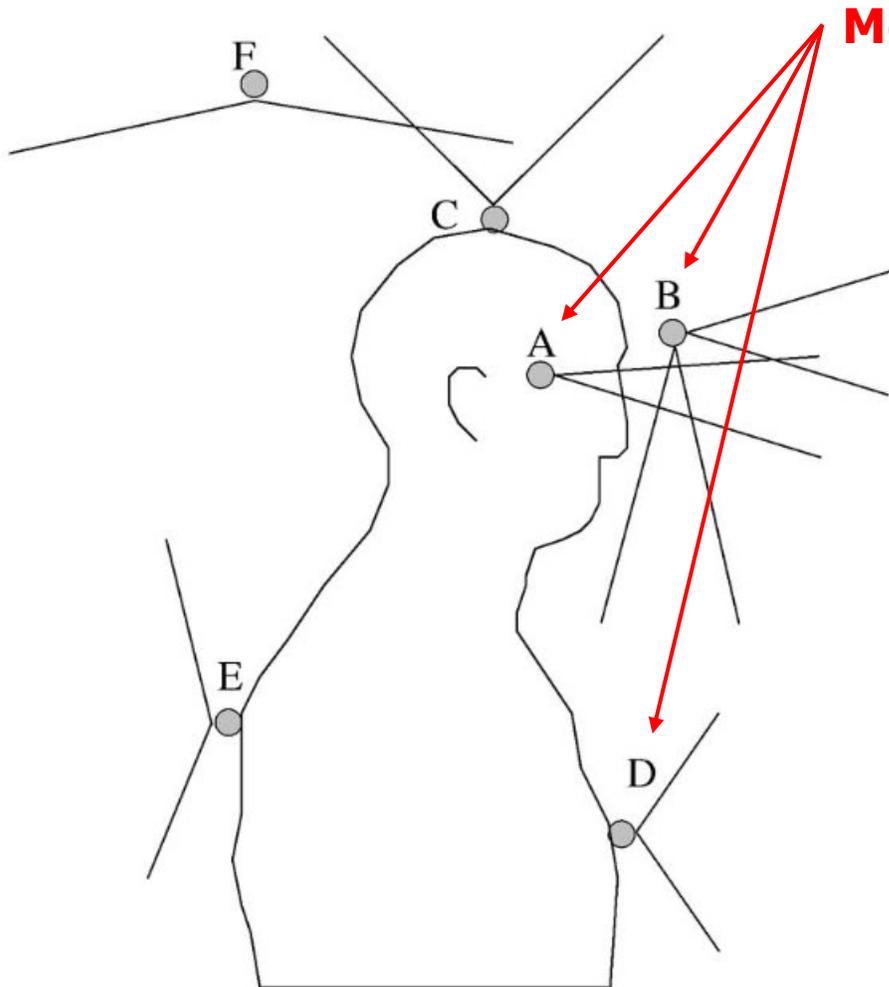


B (downward)



D





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;

A wide FOV allows to capture more scene but it may introduce distortion

Narrow Angle

Wide Angle



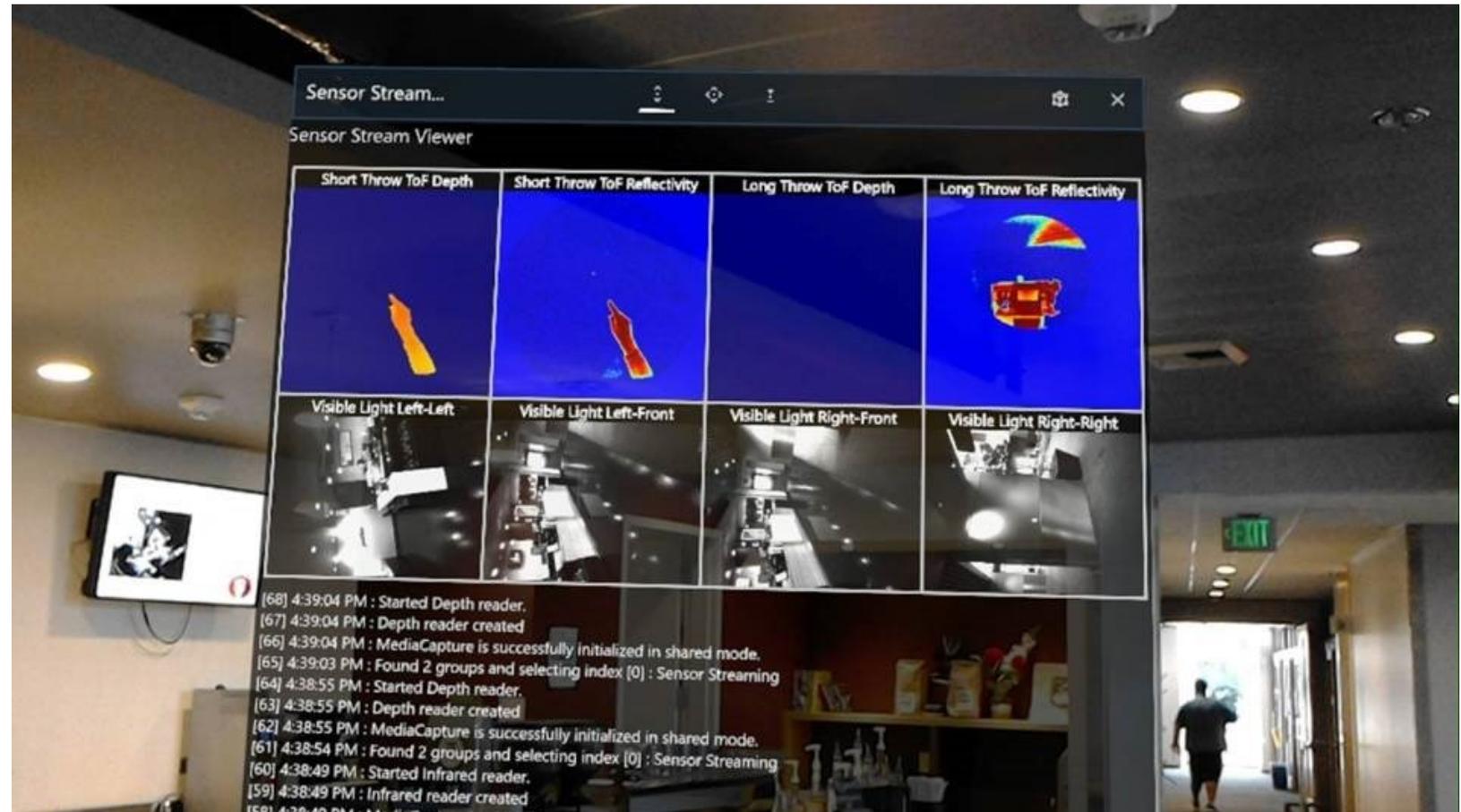
- Depth can improve scene understanding by highlighting the position of objects and hands;



<https://github.com/microsoft/HoloLensForCV>

Microsoft HoloLens Research Mode

- Microsoft HoloLens has a «Research Mode» which allows to access:
 - short-range depth
 - long-range depth;
 - IR reflectivity;



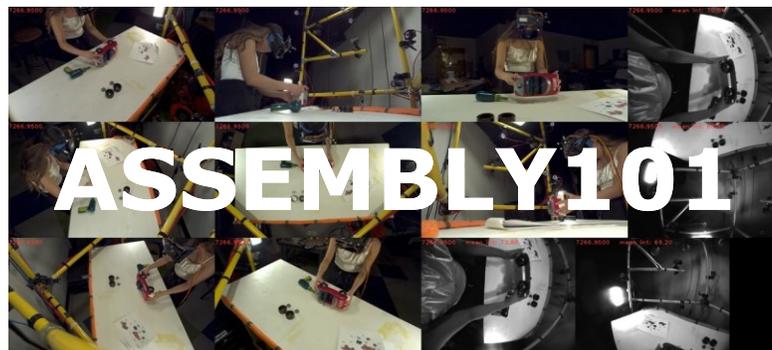
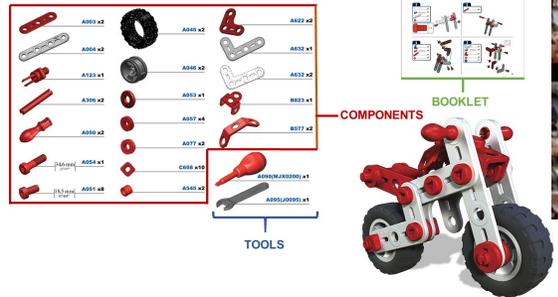
<https://docs.microsoft.com/en-us/windows/mixed-reality/research-mode>

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

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



MECCANO



Dataset	URL	Settings	Annotations	Goal
EGO-EXO4D	https://ego-exo4d-data.org/	839 participants performing procedural and physical activities.	Natural language descriptions, segmentation masks, temporal segments of keysteps, task-graphs, proficiency labels, 3D human pose	Keystep Recognition, Proficiency Estimation, Relation, Pose Estimation
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

Dataset	URL	Settings	Annotations	Goal
ENIGMA-51	https://iplab.dmi.unict.it/ENIGMA-51/	Participants performing procedural activities in the industrial domain.	Textual procedures, Hand and Object annotations, human-object interactions, next-object interactions	Untrimmed temporal annotations of human-object interactions, Egocentric Human-object interactions, short-term object interaction anticipation, NLU of intents and entities
HOLOASSIST	https://holoassist.github.io/	350 instructor-performer pairs which collaboratively complete physical manipulation tasks.	Action and conversational annotations	Action recognition and anticipation, mistake detection, intervention type prediction, 3D hand pose forecasting
ARIA Digital Twin	https://www.projectaria.com/datasets/adt/			
IndustReal	https://timschoonbeek.github.io/industreal.html	Participants performing procedural activities building a toy model of a car	Step and mistake annotations	Action recognition, assembly state detection, procedure step recognition

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	Understading daily activities, action recognition
ADL	https://www.csee.umbc.edu/~hpirs/iaj/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, sandwiche	action segments	Understading daily activities
EgoSeg	http://www.vision.huji.ac.il/egoseg/	Long term actions (walking, running, driving, etc.)	long term activity	Temporal Segmentation, Indexing

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

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

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

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/semantic-hyperlapse/cvpr2018-dataset/	80 hours/different activities	Scene/Action labels with IMU, GPS mad depth	Summarization

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



12 Egocentric Vision Research Tasks

1. Localisation
2. 3D Scene Understanding
3. Anticipation
4. Action Recognition
5. Gaze Understanding and Prediction
6. Social Behaviour Understanding
7. Full Body Pose Estimation
8. Hand and Hand-Object Interactions
9. Person Identification
10. Privacy
11. Summarisation
12. Visual Question Answering



12 Egocentric Vision Research Tasks

1. **Localisation**
2. 3D Scene Understanding
3. **Anticipation**
4. **Action Recognition**
5. Gaze Understanding and Prediction
6. Social Behaviour Understanding
7. Full Body Pose Estimation
8. **Hand and Hand-Object Interactions**
9. Person Identification
10. Privacy
11. Summarisation
12. Visual Question Answering

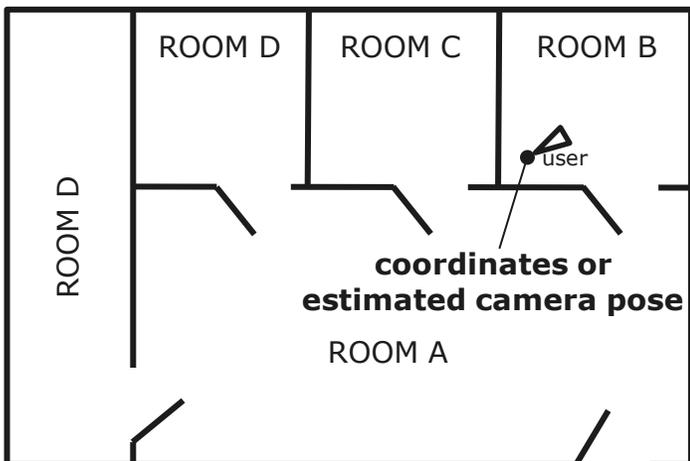
Localization

SCENE RECOGNITION

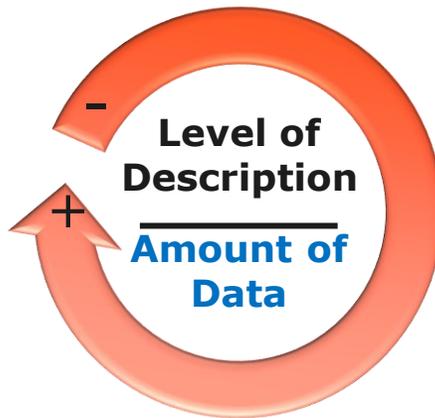


off-the-shelf detectors

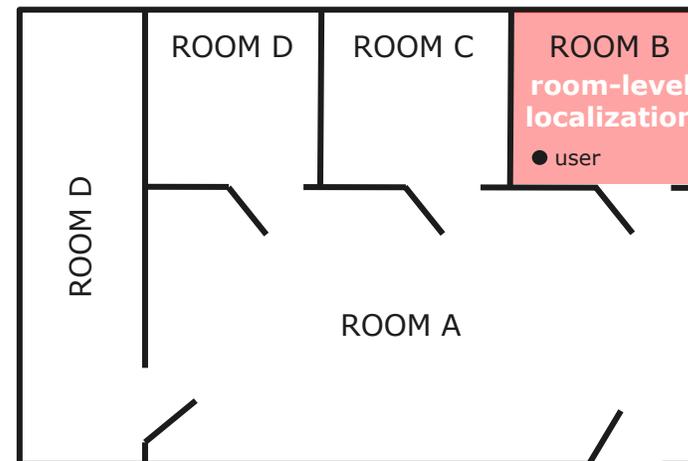
CAMERA POSE-ESTIMATION



3D reconstruction of the building

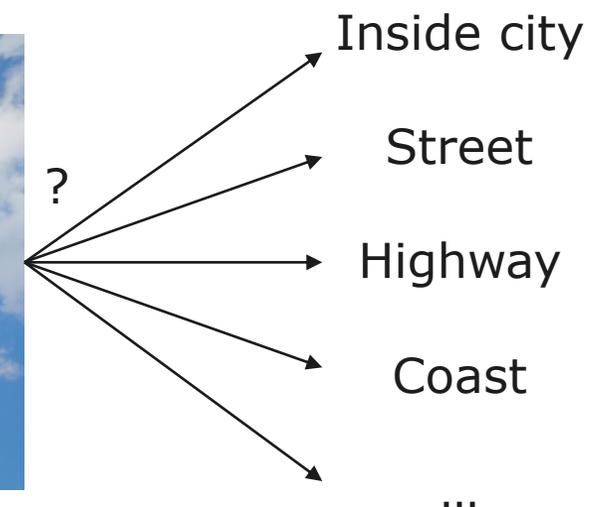
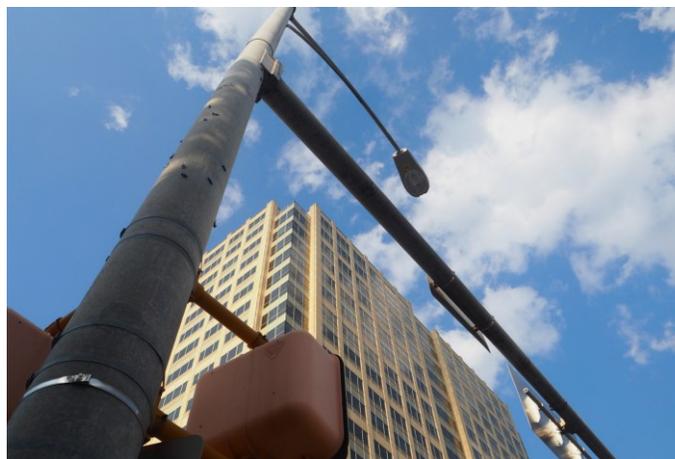


ROOM-LEVEL RECOGNITION



moderate amount of training data

- 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.



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



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)

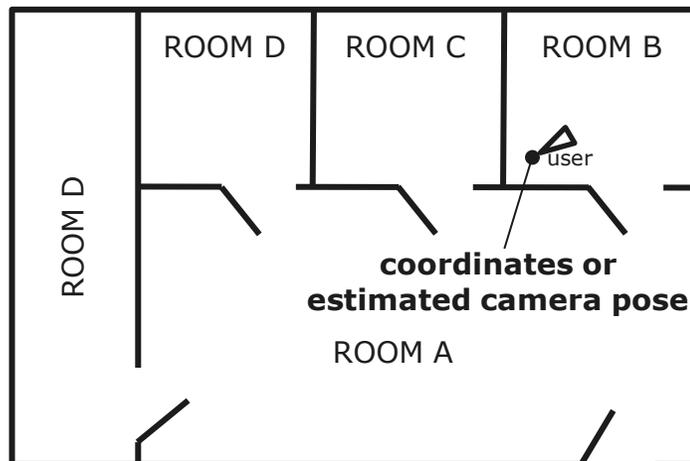
- Places is a large (10M images – 400+ classes) dataset for scene recognition;
- CNN models trained to recognize 365 scene classes available for download;
- Can be used off-the-shelf!

SCENE RECOGNITION

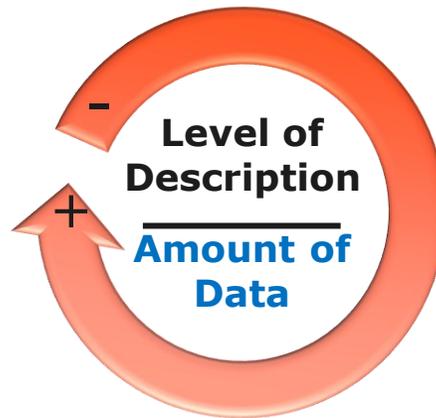


off-the-shelf detectors

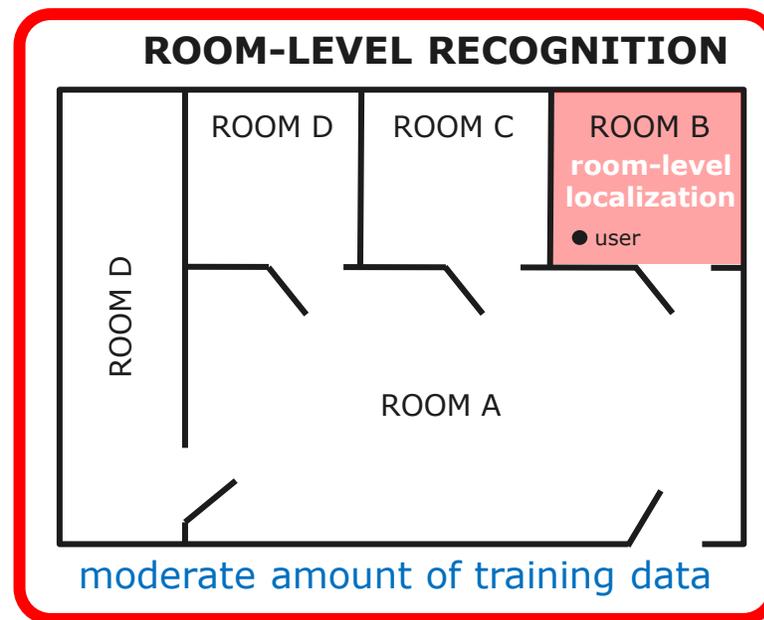
CAMERA POSE-ESTIMATION



3D reconstruction of the building



ROOM-LEVEL RECOGNITION



moderate amount of training data

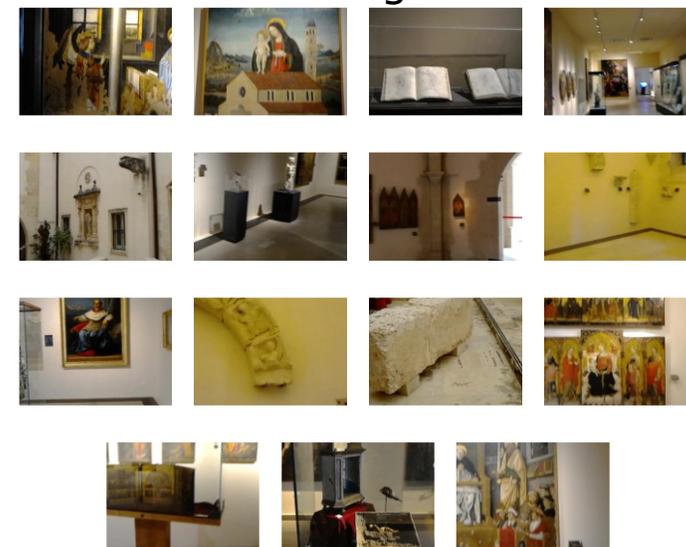
Cultural Site (e.g., museum) divided into contexts (e.g., rooms)



(videos acquired in the different contexts)



Training Set



(frames extracted from videos
acquired in the different contexts)

CODE HERE -> <https://iplab.dmi.unict.it/VEDI/>

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



Training Set (room-based images)



There are no training negatives!

1. Discrimination

estimation of $P(y_i | I_i, y_i \neq 0)$

$$\arg \max_j P(y_i = j | I_i, y_i \neq 0)$$



2. Negative Rejection

estimation of $P(y_i | I_i)$

$$\arg \max_j P(y_i = j | I_i)$$

estimation of $P(y_i = 0 | I_i)$
(variation ratio)



3. Sequential Modelling

application of HMM

$$\arg \max_L P(L | V)$$





VEDI – Vision Exploitation for Data Interpretation, PON MISE Horizon 2020

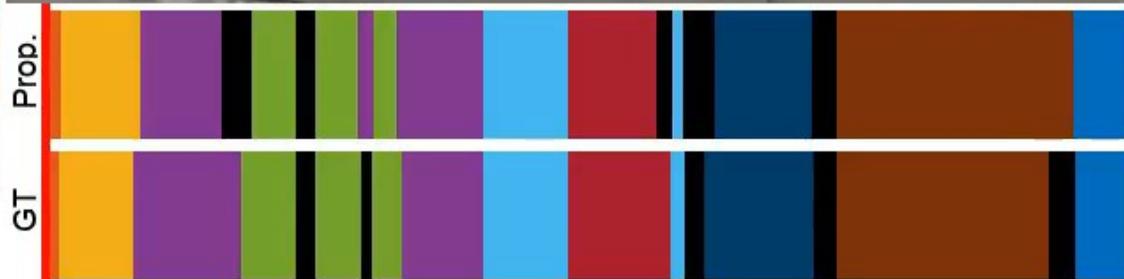
F. Ragusa, A. Furnari, S. Battiato, G. Signorello, G. M. Farinella

Time Spent at Location

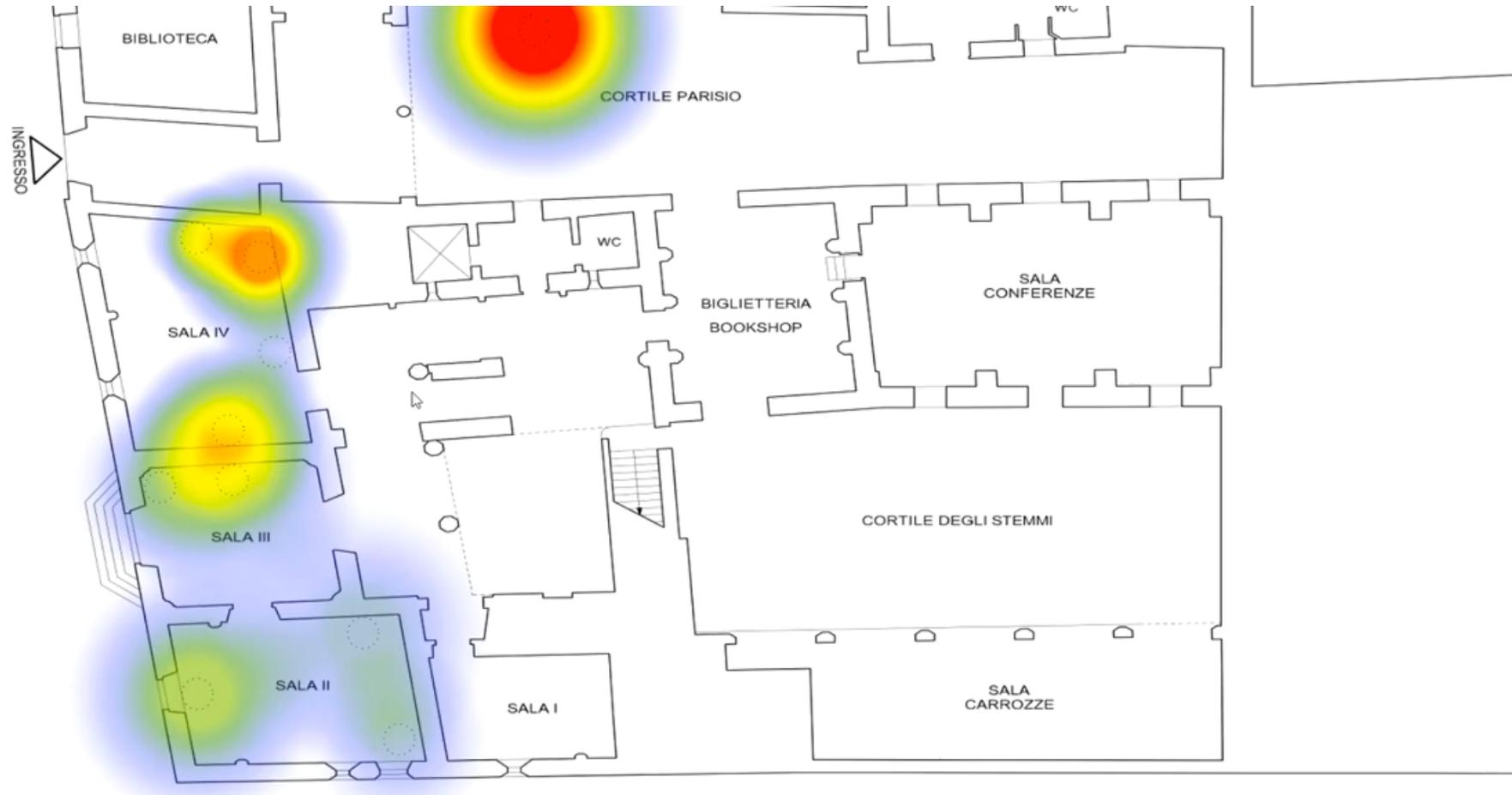
LOC	EST	GT
G. Novizi	00:00	00:00
Cortile	00:03	00:03
Scalone	00:00	00:00
Corridoi	00:00	00:00
C. Notte	00:00	00:00
Antiref.	00:00	00:00
S. Mazz.	00:00	00:00
Cucina	00:00	00:00
Ventre	00:00	00:00
Negative	00:00	00:00



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		



Ø = numero di visitatori per opera
= tempo di visione opera

Heatmap contesti

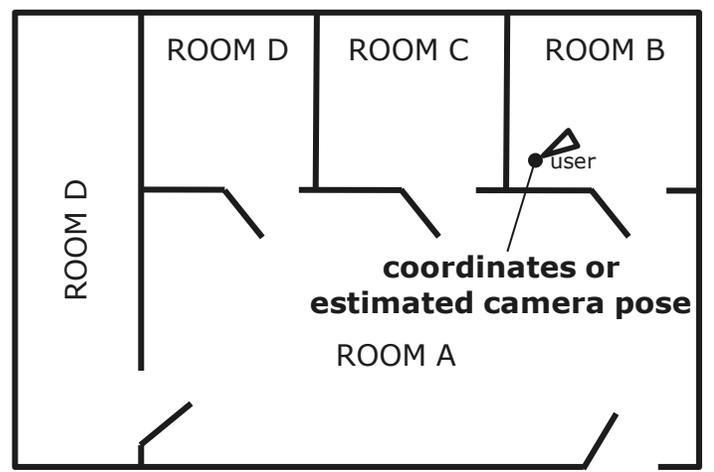
Percorsi opere

SCENE RECOGNITION

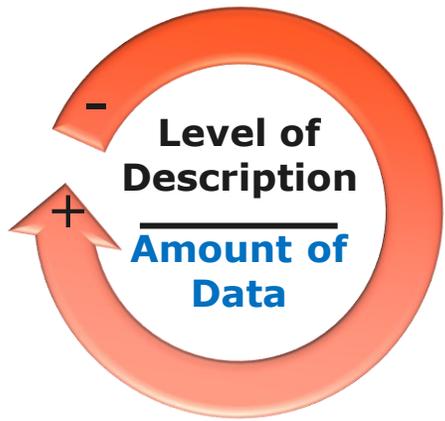


off-the-shelf detectors

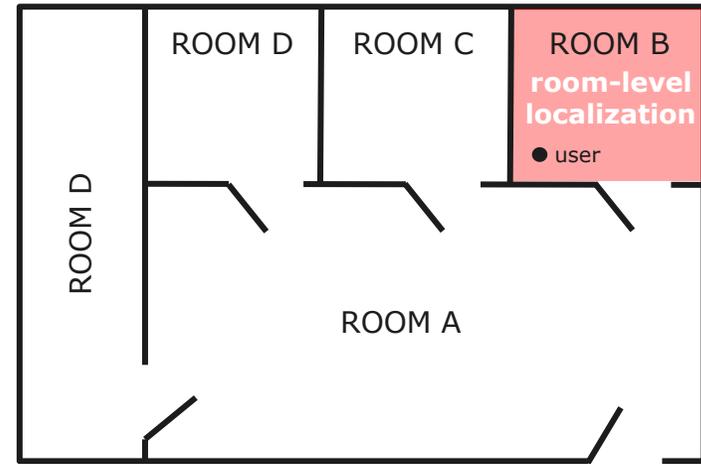
CAMERA POSE-ESTIMATION



3D reconstruction of the building

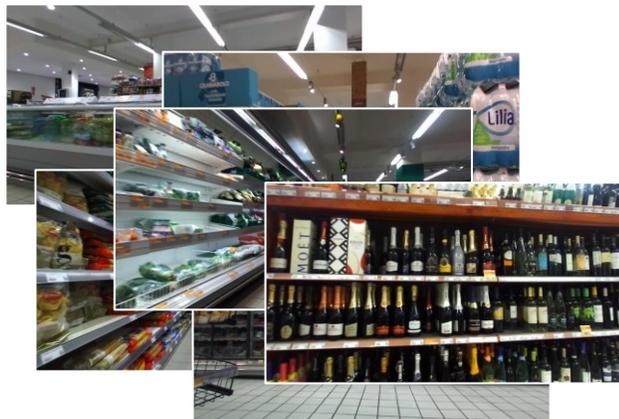


ROOM-LEVEL RECOGNITION

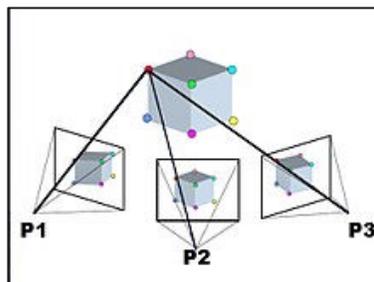


moderate amount of training data

Images



Structure from Motion (SfM)



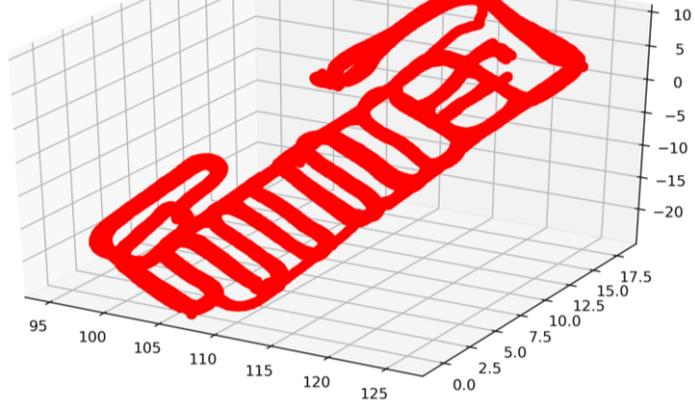
(P,Q)

Attach estimated 6DOF pose to each image

3D Model

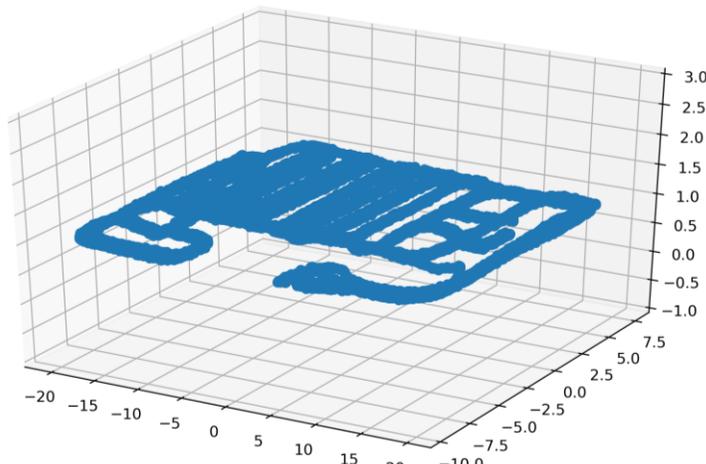


Arbitrary Coordinate System (pose/scale)

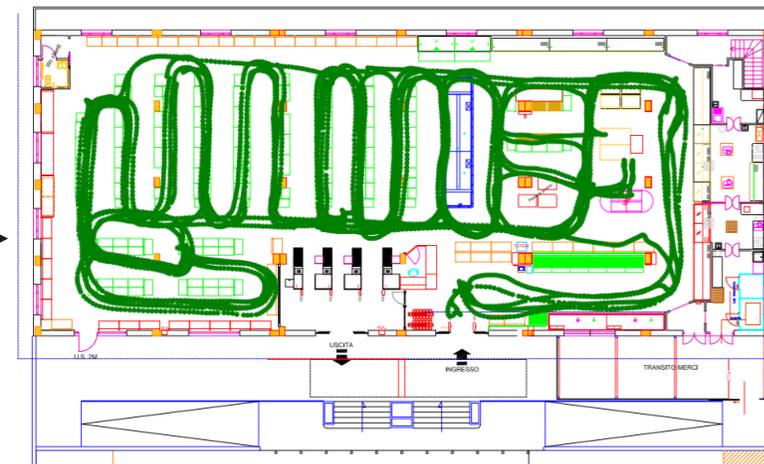


camera poses

PCA



rotated poses



scaled/aligned poses

Structure from Motion attaches every input image to a 3D model.



Many options available:

COLMAP (free)

<https://colmap.github.io/>

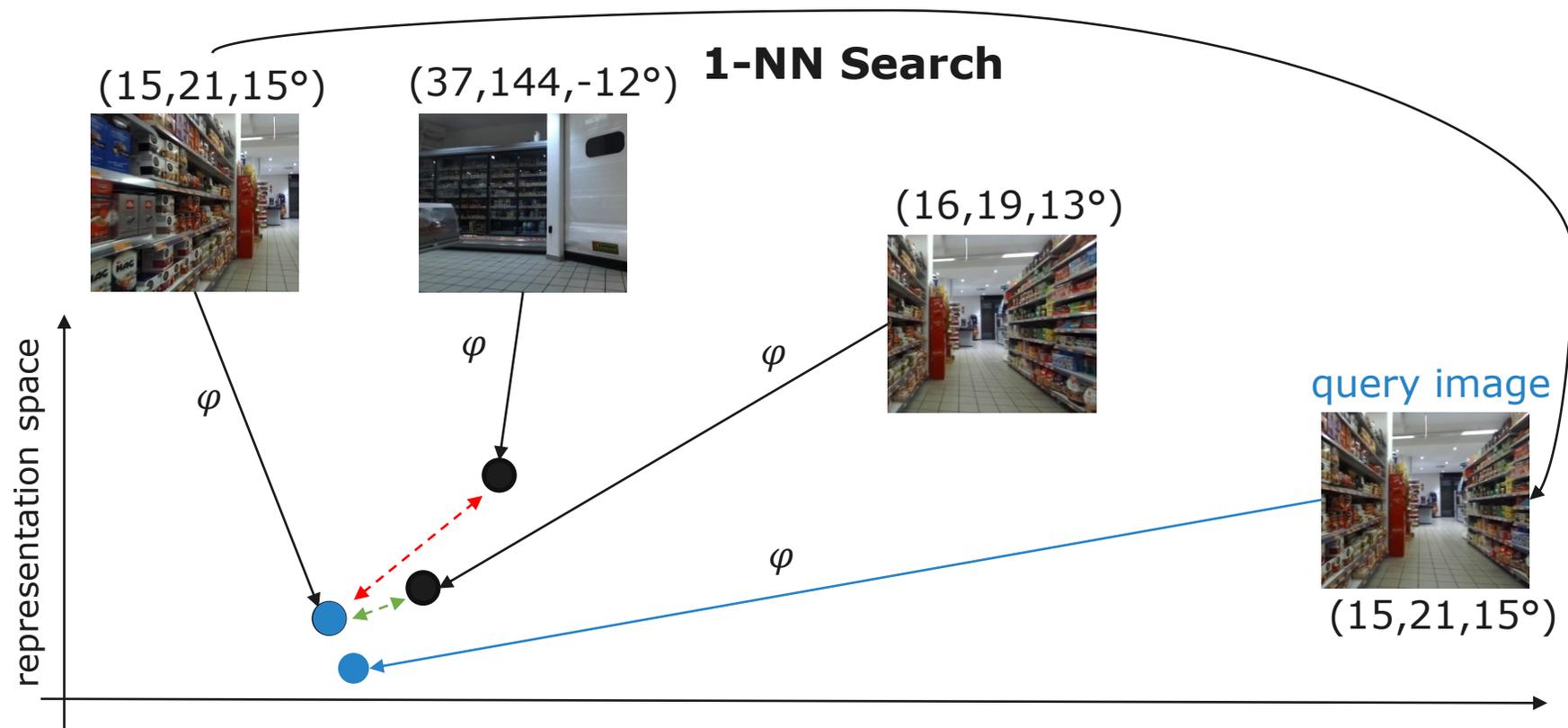
Visual SFM (free)

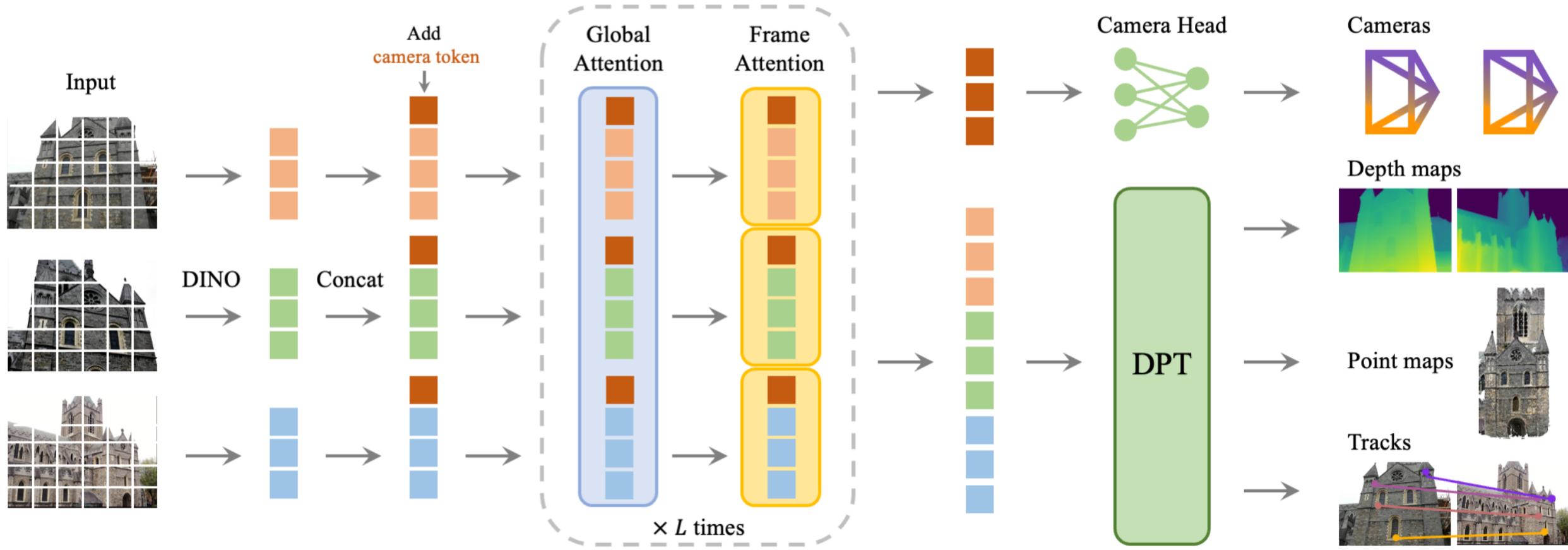
<http://ccwu.me/vsfm/>

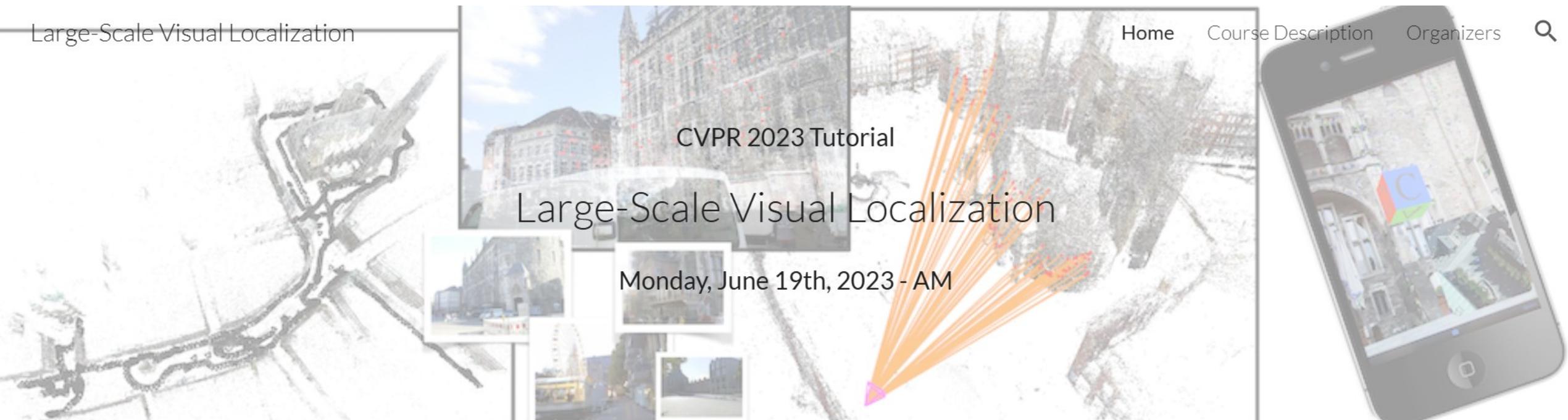
3D Zephyr (paid)

<https://www.3dflow.net/it/3df-zephyr-pro-3d-models-from-photos/>

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







Course Information

- **When:** Monday, June 19th, 2023
- **Where:** East 2
- **Time:** 8:30 AM - 12:15 PM (local time)
- **Schedule**
 - **Introduction** [5 min] [8:30 - 8:35]
 - **Part I: Instance retrieval for coarse localization** [45 min](Giorgos) [8:35 - 9:20] [[slides](#)]



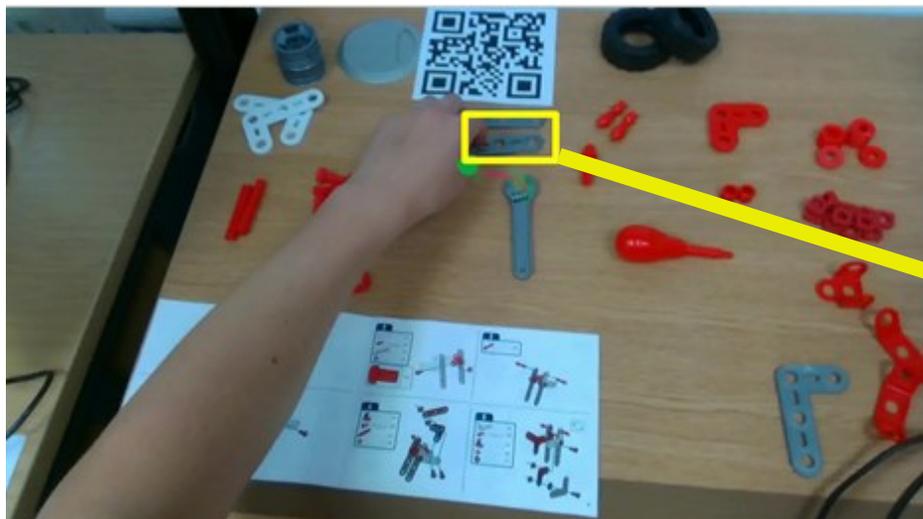
<https://sites.google.com/view/lsvl2023/home>

Object Detection/Interaction

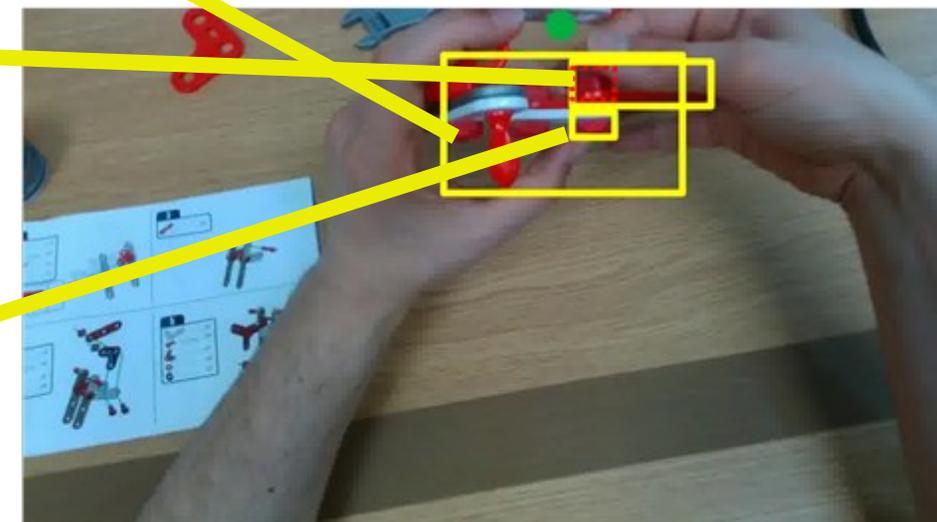
Objects and Actions are tight!

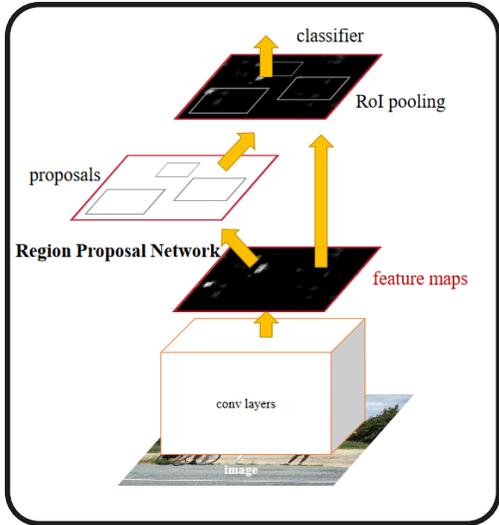
Useful to know what is in the scene

Useful to know what actions can be performed

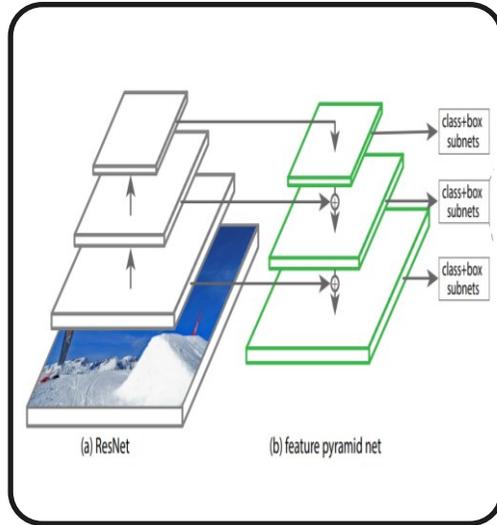


ID	Class
0	instruction booklet
1	gray_angled_perforated_bar
2	partial_model
3	white_angled_perforated_bar
4	wrench
5	screwdriver
6	gray_perforated_bar
7	wheels_axle
8	red_angled_perforated_bar
9	red_perforated_bar
10	rod
11	handlebar
12	screw
13	tire
14	rim
15	washer
16	red_perforated_junction_bar
17	red_4_perforated_junction_bar
18	bolt
19	roller

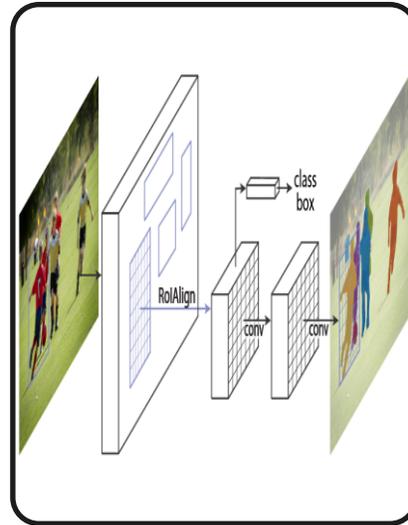




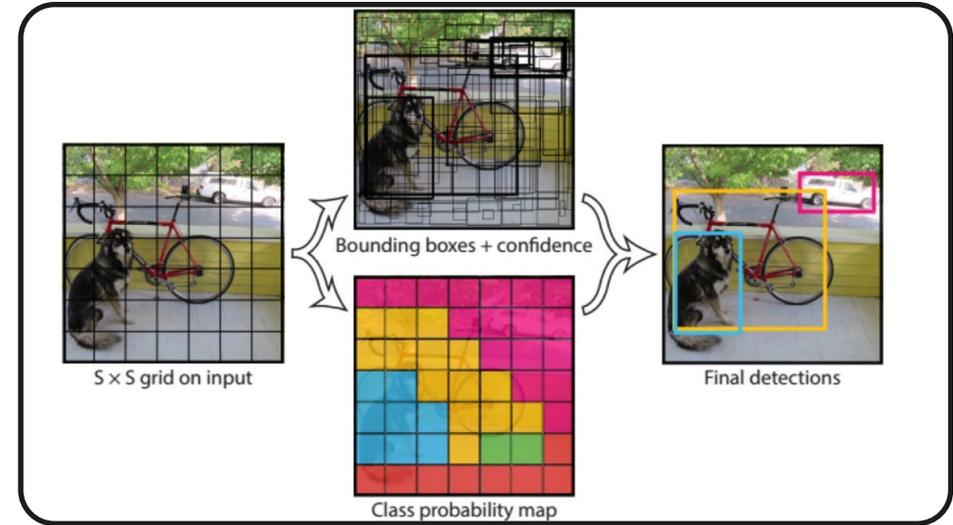
Faster-RCNN
(bounding boxes)



RetinaNet
(bounding boxes - faster)



Mask-RCNN
(boxes + segments)



YOLO
(much faster, but less accurate)

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

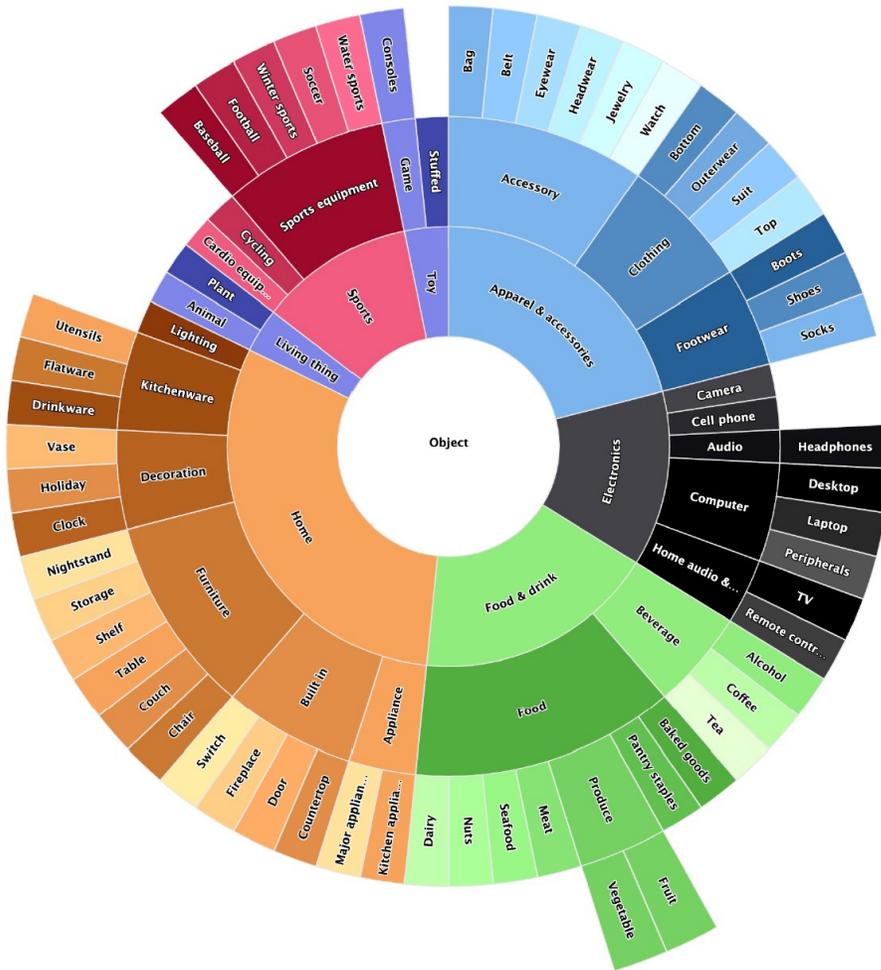
<https://pjreddie.com/darknet/yolo/>

Transformer-Based Detectors: <https://github.com/IDEA-Research/awesome-detection-transformer>

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.

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

368 categories!



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

C. Zhu et al., "EgoObjects: A Large-Scale Egocentric Dataset for Fine-Grained Object Understanding" in IEEE/CVF International Conference on Computer Vision (ICCV), 2023



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



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



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



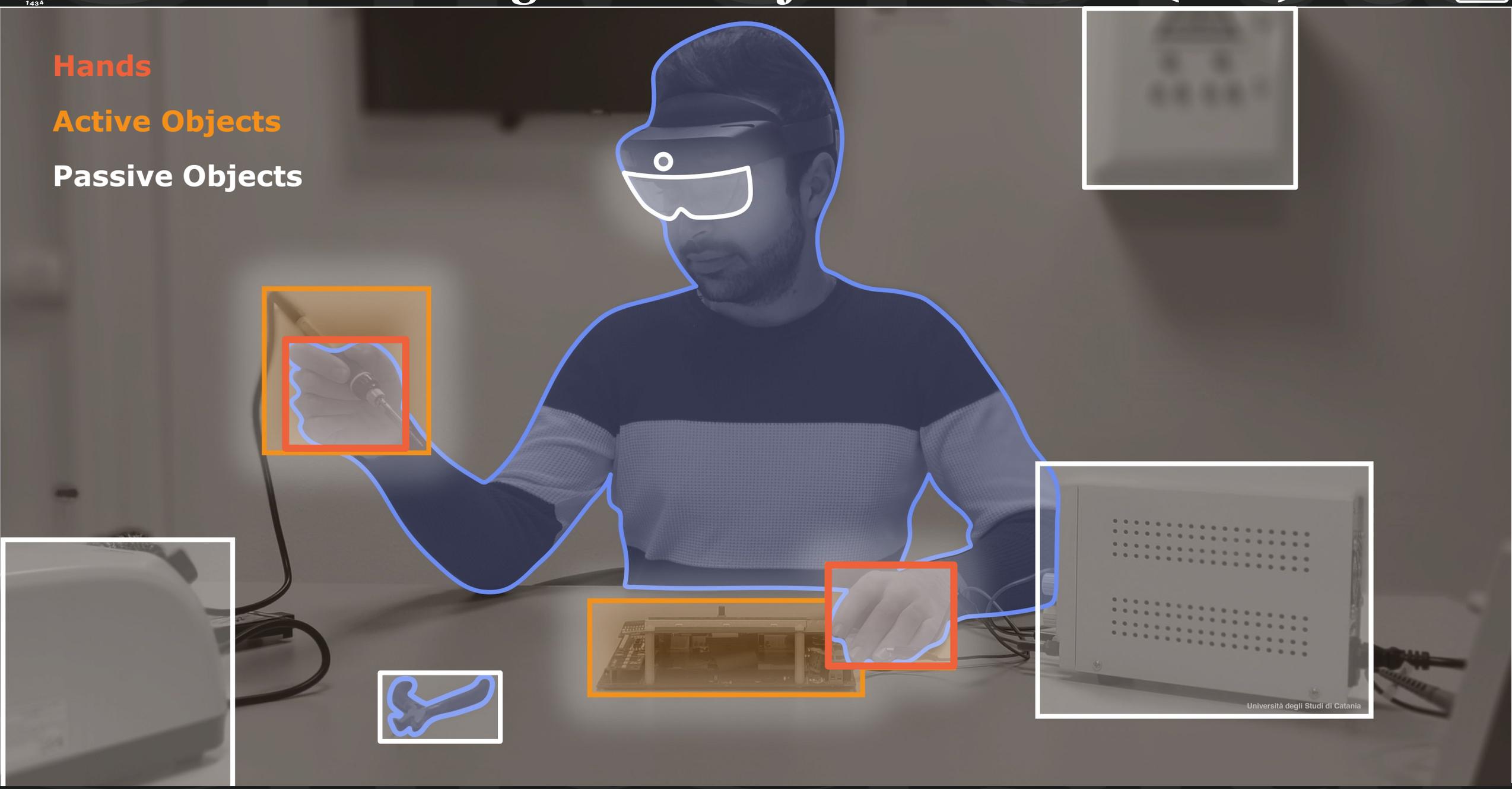
<https://iplab.dmi.unict.it/MECCANO/>

- In some scenarios, it could be necessary to fine-tune an object-detector with application-specific data.
- Main egocentric datasets providing bounding box annotations.
- EGO4D is multi- domain annotated with 295K bounding boxes.

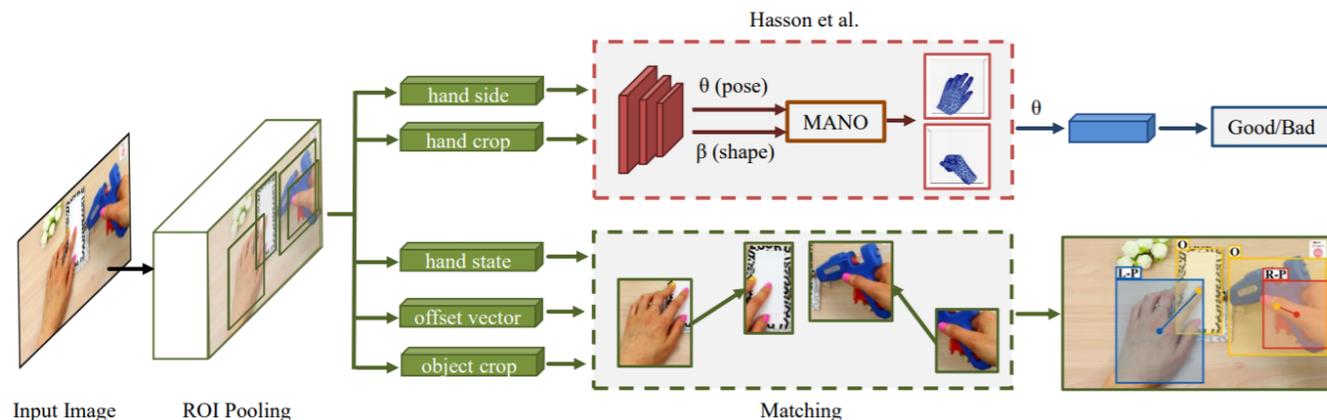
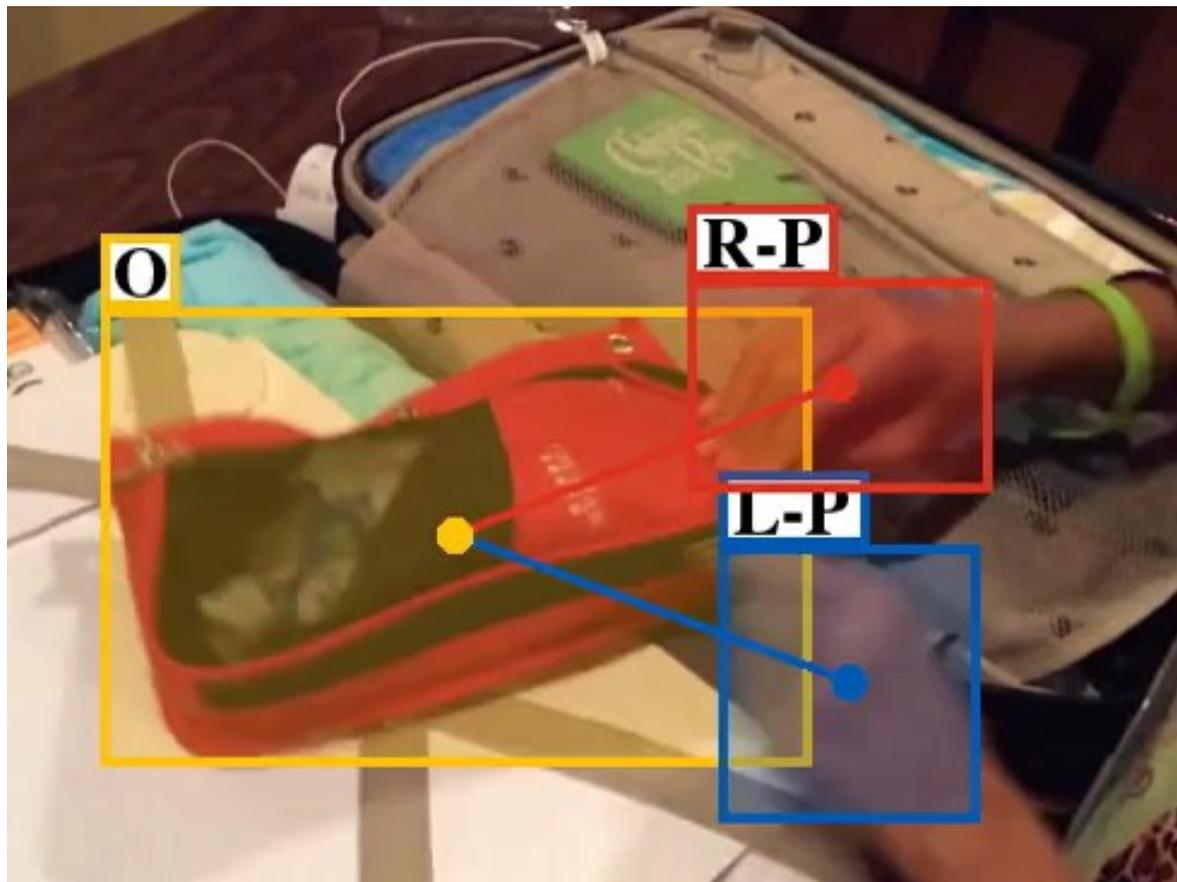
Hands

Active Objects

Passive Objects



CODE & DATA HERE -> <https://fouheylab.eecs.umich.edu/~dandans/projects/100DOH/>



An «augmented» detector which recognizes:

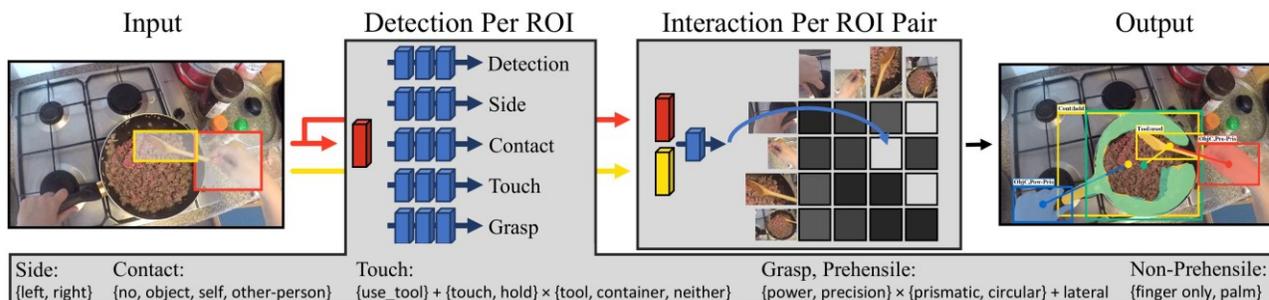
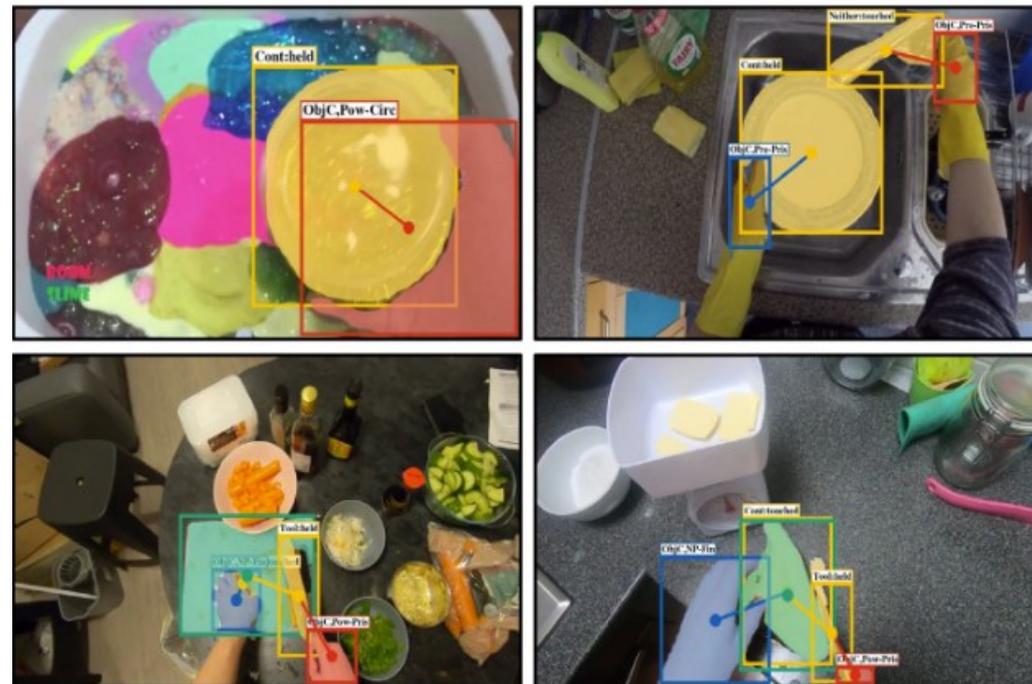
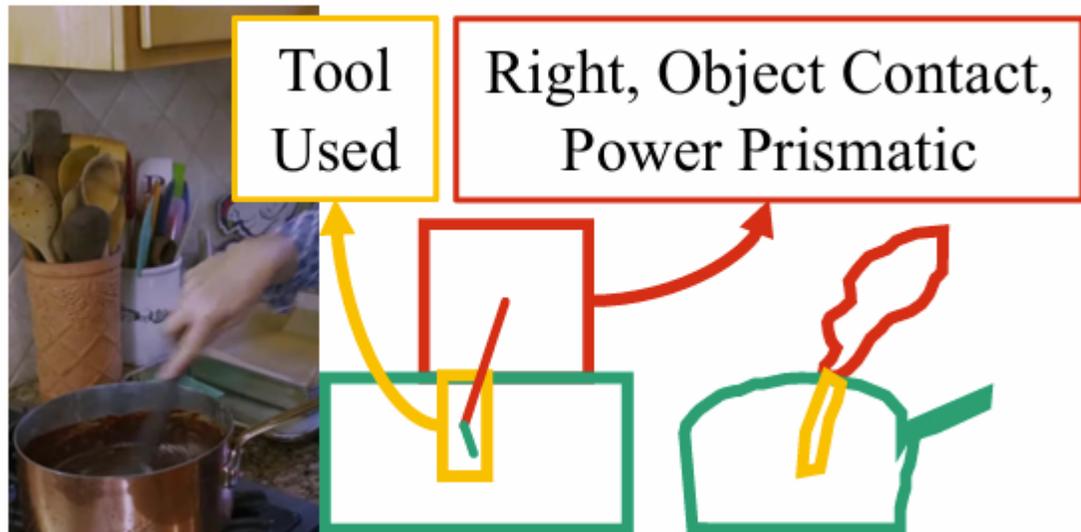
- The left hand;
- The right hand;
- The interacted object.

VISOR DATASET

Darkhalil, Ahmad, et al. "Epic-kitchens visor benchmark: Video segmentations and object relations." *Advances in Neural Information Processing Systems* 35 (2022): 13745-13758.

Shan, D., Geng, J., Shu, M., & Fouhey, D. F. (2020). Understanding human hands in contact at internet scale. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (pp. 9869-9878).

CODE & DATA HERE -> <https://fouheylab.eecs.umich.edu/~dandans/projects/hands23/>



An «augmented» detector which recognizes:

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



Standard approach:

- Collect a lot of images and videos of construction sites;
- Label the data with domain-specific annotations;
- Train and test deep learning algorithms.

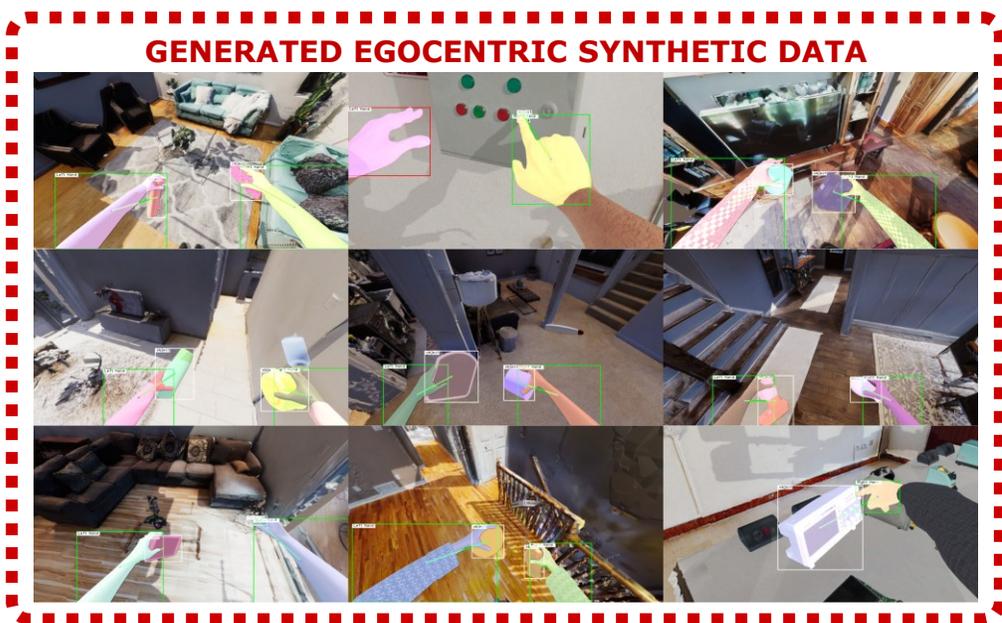


**What if we could learn the
«real thing» in simulation?**

DATA HERE -> <https://iplab.dmi.unict.it/HOI-Synth/>

Are Synthetic Data Useful for Egocentric Hand-Object Interaction Detection?





- **Epic-Kitchens VISOR**
 - 32,857 real + 30,259 synthetic images
- **EgoHOS**
 - 8,107 real + 8,107 synthetic images
- **ENIGMA-51**
 - 3,479 real images
 - In-Domain 16,773 + Out-domain 20,321 synthetic images

a) Unsupervised Setting

% Real Labeled Data	Approach	Overall	H	H+S	H+C	O
0%	Synthetic-Only	09.88	28.41	24.89	08.64	01.23
	UDA	33.33	80.16	65.98	33.47	08.35
Absolute Improvement		+23.45	+51.75	+41.09	+24.83	+7.12

b) Semi-supervised Setting

% Real Labeled Data	Approach	Overall	H	H+S	H+C	O
25% (8,215 images)	Real-Only	37.90	90.14	85.66	53.99	17.85
	Synthetic+Real	38.19	89.98	84.67	55.88	18.49
	SSDA	45.55	90.37	84.42	52.59	22.15
Absolute Improvement		+7.65	+0.23	-0.99	+1.89	+4.30

C) Fully-supervised Setting

% Real Labeled Data	Approach	Overall	H	H+S	H+C	O
100% (32,857 images)	Real-Only	45.33	92.25	88.54	59.24	24.23
	Synthetic+Real	44.52	91.45	88.94	56.55	27.77
	FSDA	46.48	91.83	87.65	57.63	24.03
Absolute Improvement		+1.15	-0.42	+0.40	-1.61	+3.54

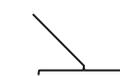
Action Recognition



Model

VERB

NOUN



Open

$v = 3$

- Box

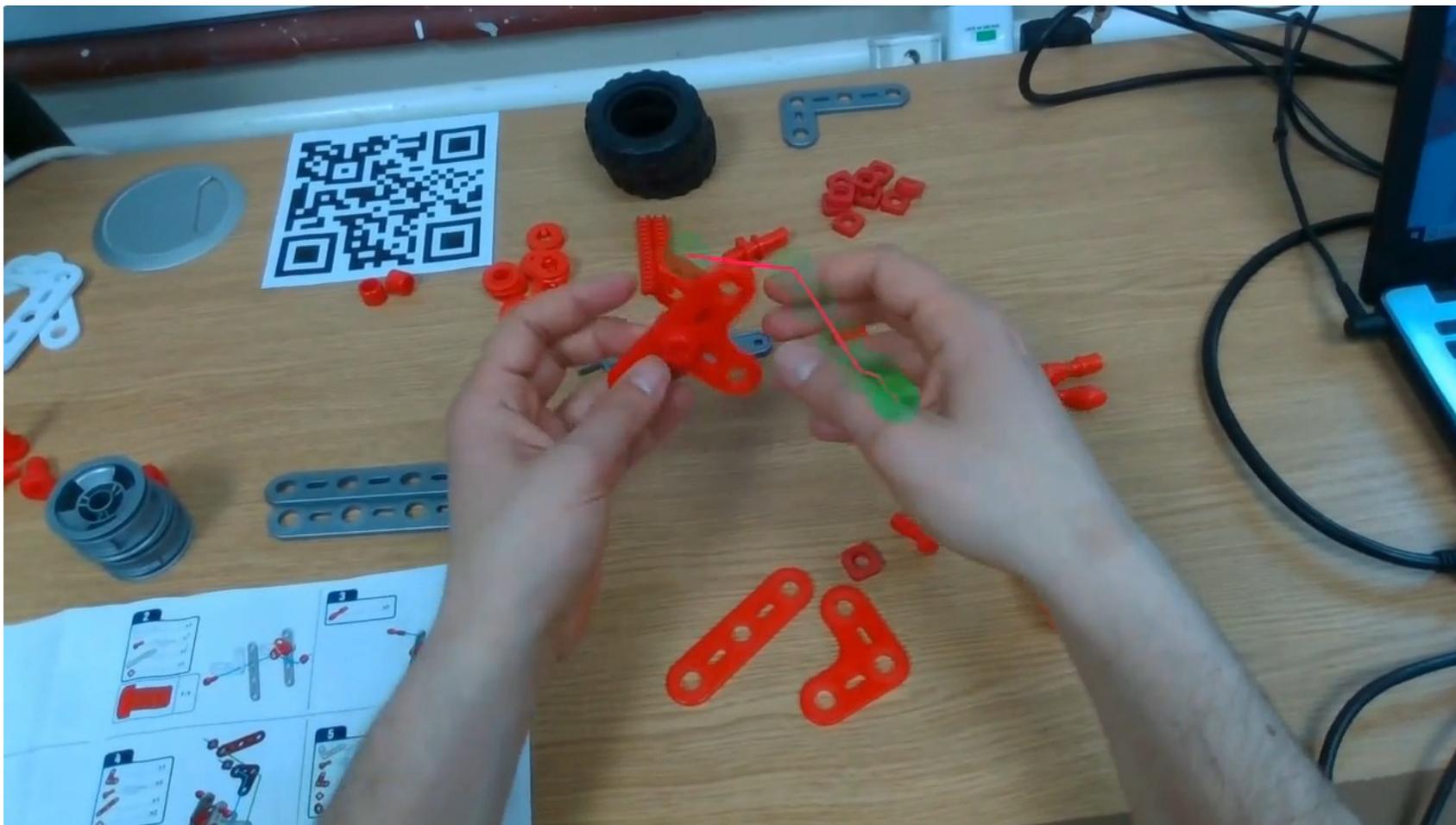
$n = 23$



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

As defined in EPIC-KITCHENS-2020

TAKE SCREWDRIVER



TAKE SCREWDRIVER



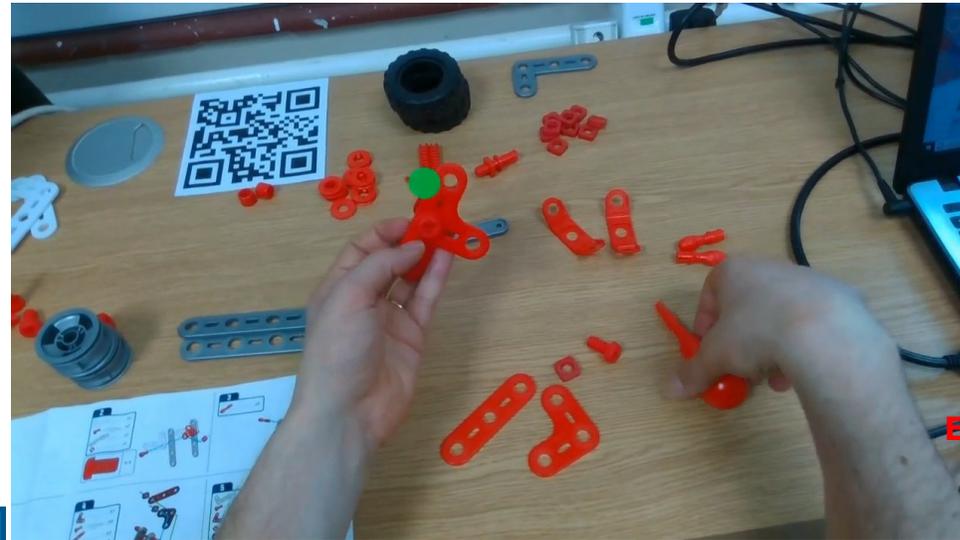
Start Action

Start Interaction (H-O)



Frame of Contact

TAKE SCREWDRIVER



Start Action

Start Interaction (H-O)

End Action



Frame of Contact

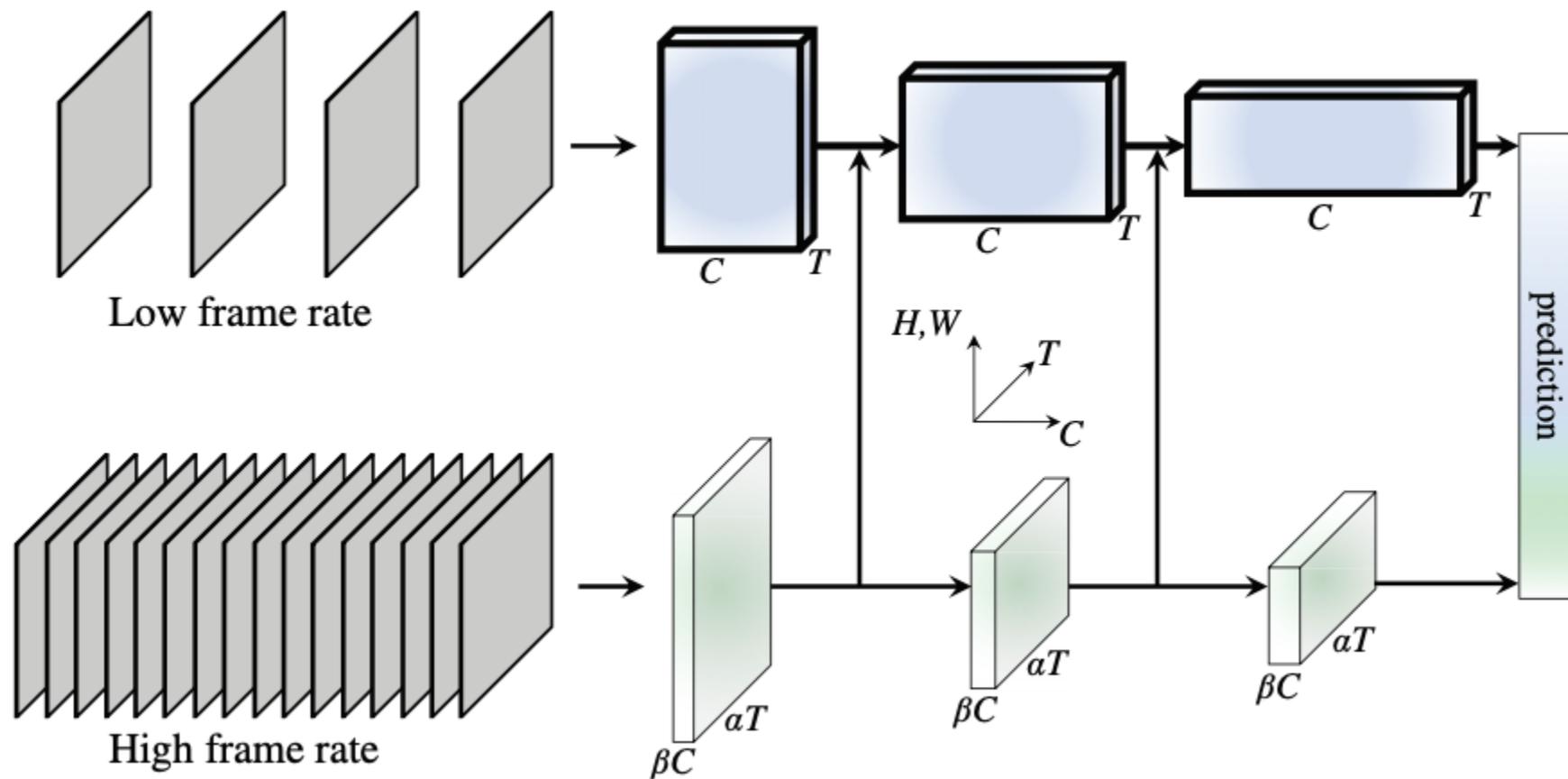
Frame of Decontact



F. Ragusa, A. Furnari, G. M. Farinella. MECCANO: A Multimodal Egocentric Dataset for Humans Behavior Understanding in the Industrial-like Domain. Computer Vision and Image Understanding (CVIU), 2023.

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

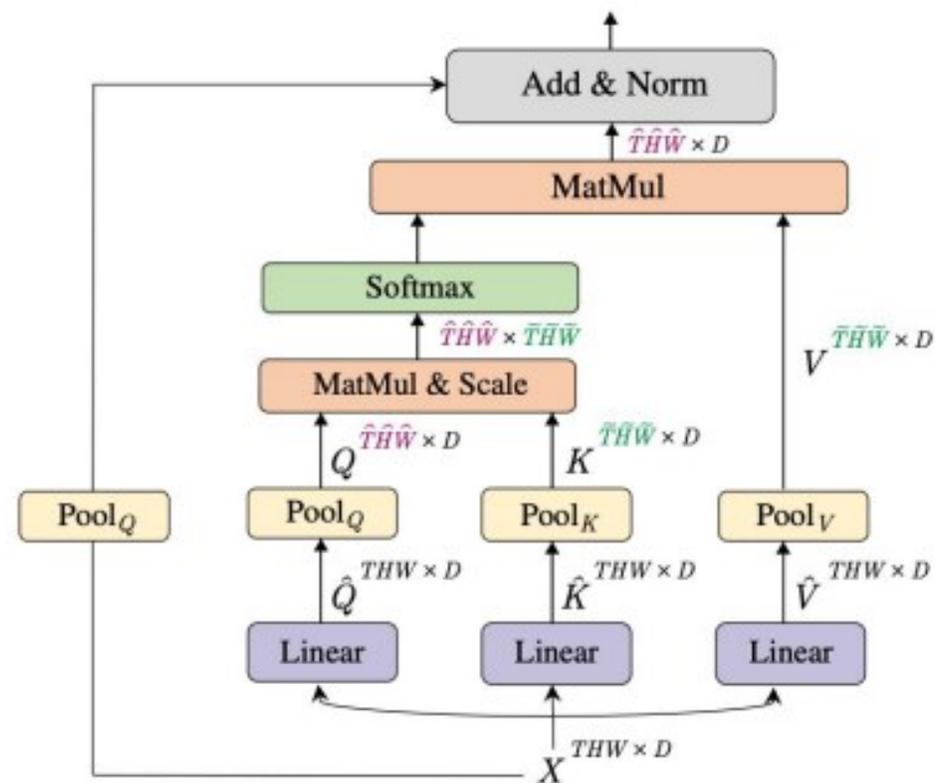
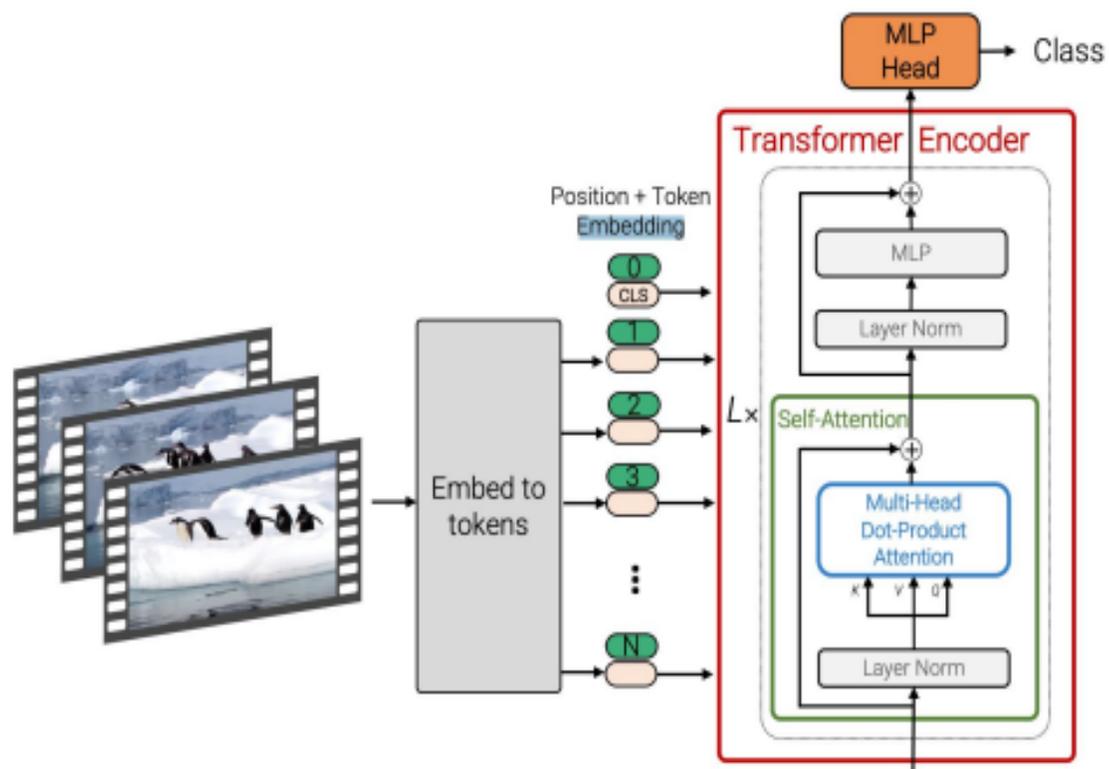
CODE HERE -> <https://github.com/facebookresearch/SlowFast>



CODE HERE -> <https://github.com/facebookresearch/SlowFast>

Factorized attention: Attend over space / time

Pooling module: Reduce number of tokens



Bertasius et al, "Is Space-Time Attention All You Need for Video Understanding?", ICML 2021

Arnab et al, "ViViT: A Video Vision Transformer", ICCV 2021

Neimark et al, "Video Transformer Network", ICCV 2021

Fan et al, "Multiscale Vision Transformers", ICCV 2021

Li et al, "MViTv2: Improved Multiscale Vision Transformers for Classification and Detection", CVPR 2022

☰ README.md

PySlowFast

PySlowFast is an open source video understanding codebase from FAIR that provides state-of-the-art video classification models with efficient training. This repository includes implementations of the following methods:

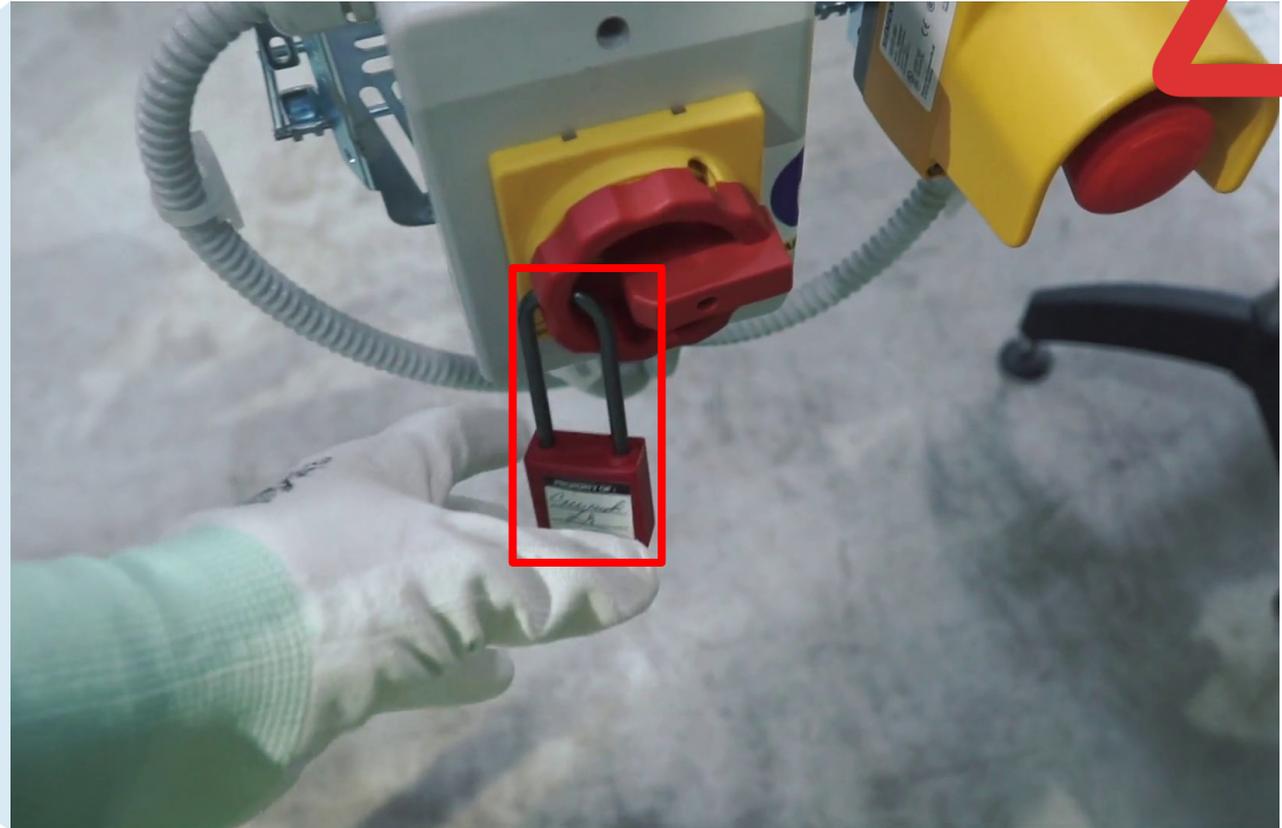
- [SlowFast Networks for Video Recognition](#)
- [Non-local Neural Networks](#)
- [A Multigrid Method for Efficiently Training Video Models](#)
- [X3D: Progressive Network Expansion for Efficient Video Recognition](#)
- [Multiscale Vision Transformers](#)
- [A Large-Scale Study on Unsupervised Spatiotemporal Representation Learning](#)
- [MViTv2: Improved Multiscale Vision Transformers for Classification and Detection](#)
- [Masked Feature Prediction for Self-Supervised Visual Pre-Training](#)
- [Masked Autoencoders As Spatiotemporal Learners](#)
- [Reversible Vision Transformers](#)

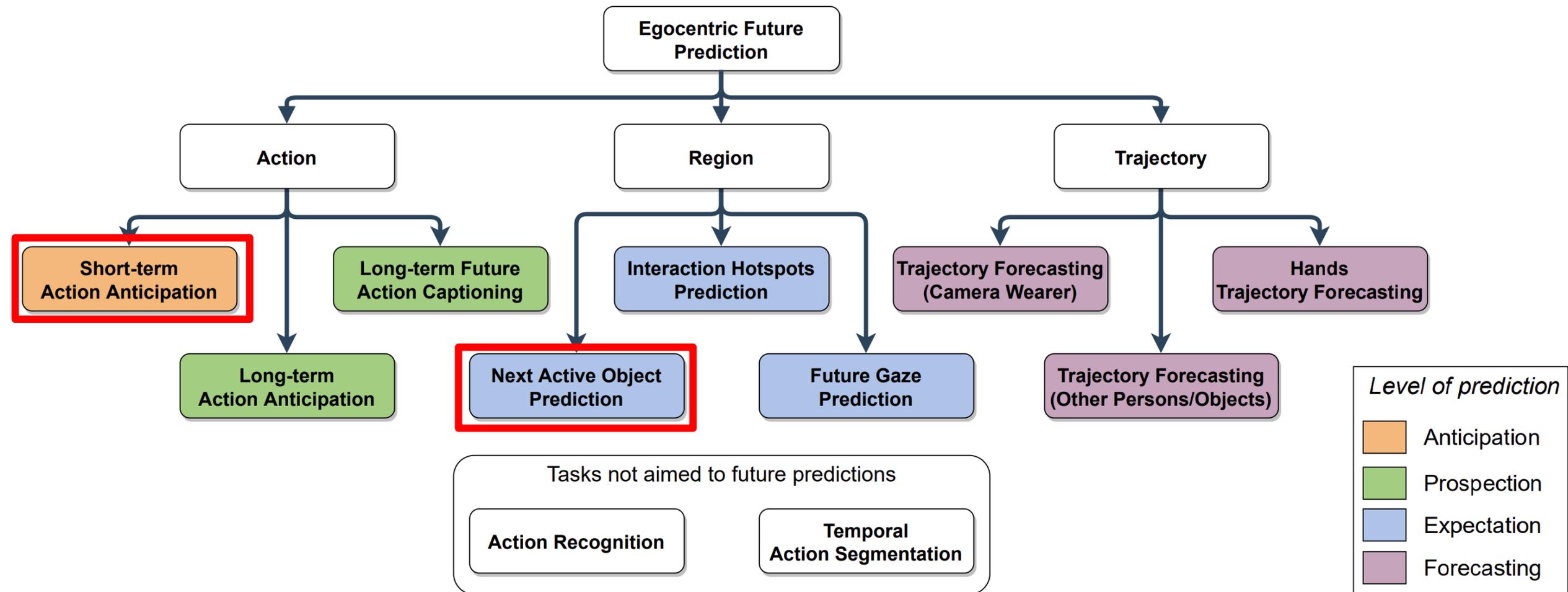
<https://github.com/facebookresearch/SlowFast>

Anticipation

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**



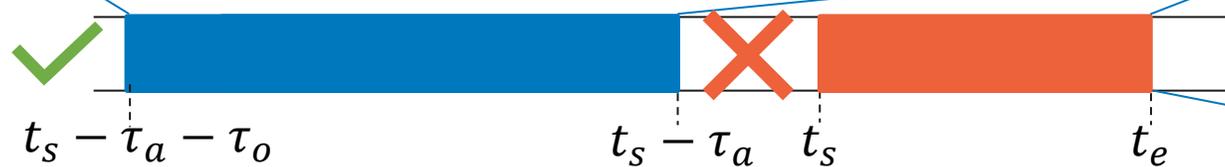


(observed video)



Model

Take - Plate



τ_0 arbitrary

$\tau_a = 1s$;



(unobserved)

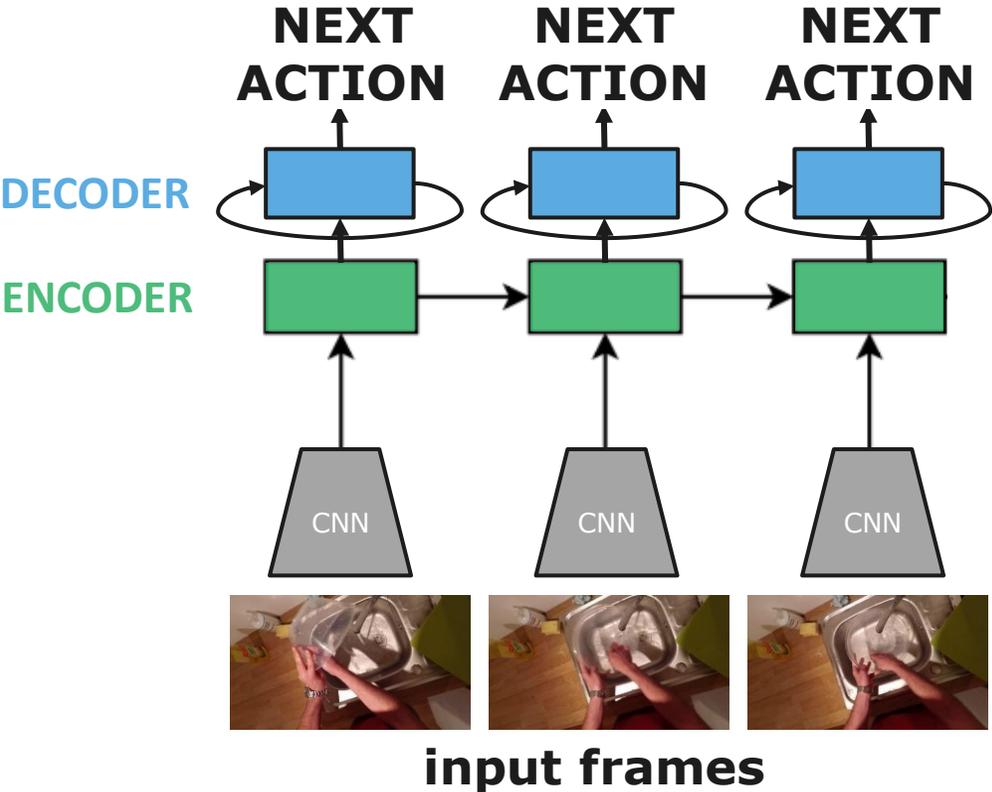
Rolling-LSTM

ENCODING

Unrolling-LSTM

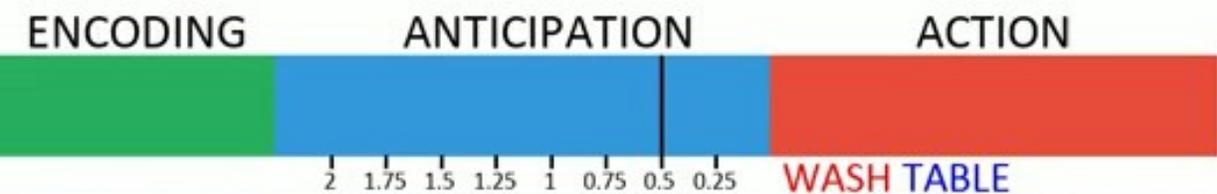
INFERENCE

We take inspiration from sequence to sequence models.



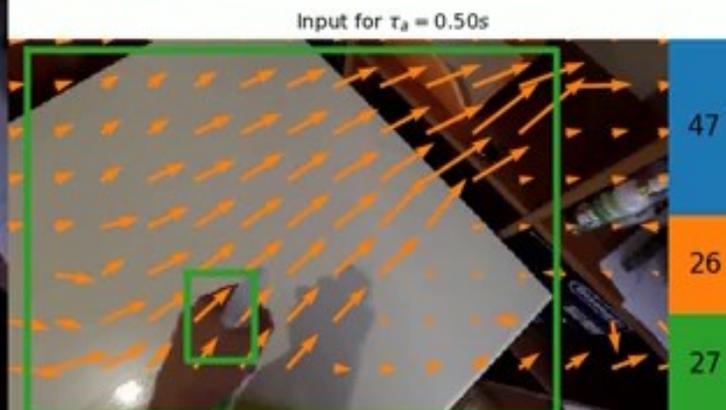
A. Furnari, G. M. Farinella, What Would You Expect? Anticipating Egocentric Actions with Rolling-Unrolling LSTMs and Modality Attention. ICCV 2019 (ORAL).

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



Anticipated Actions (in 0.50s)

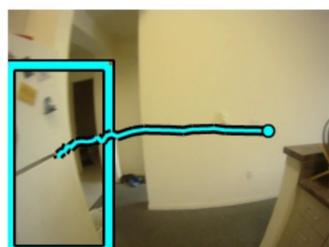
- WASH TABLE**
- SPRAY LIQUID:WASHING
- TAKE SHEETS
- MOVE BOTTLE
- PUT LIQUID:WASHING
- PUT SHEETS
- WASH TOP
- OPEN TAP
- CLOSE CUPBOARD
- TAKE BAG
- WASH SINK
- MOVE BREAD



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

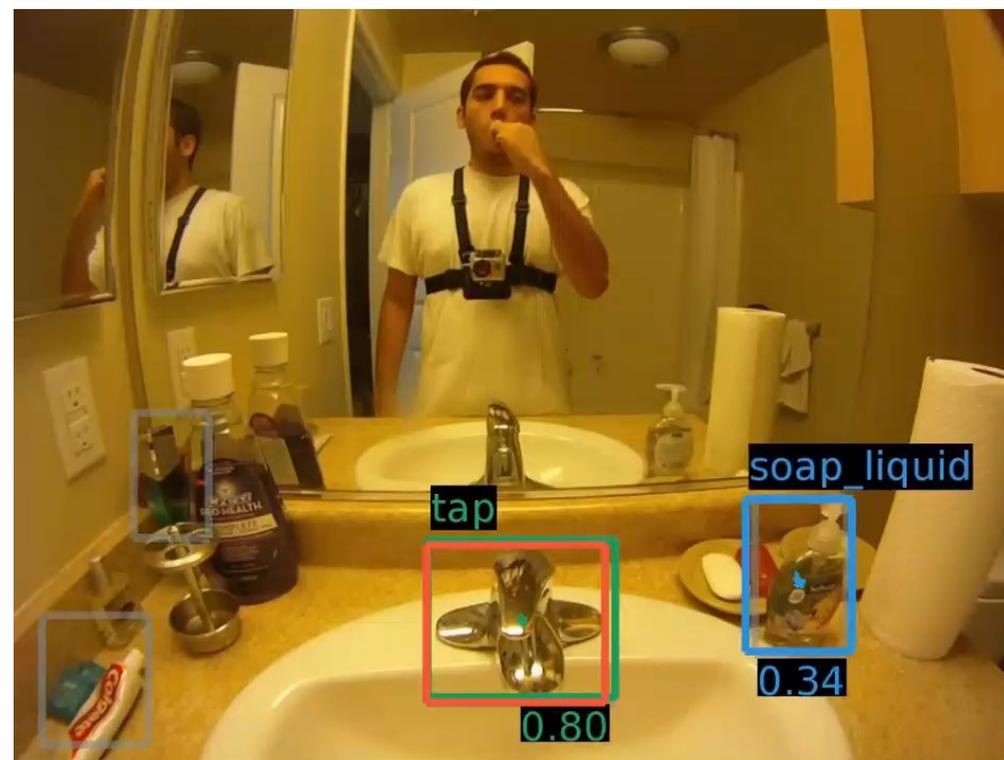
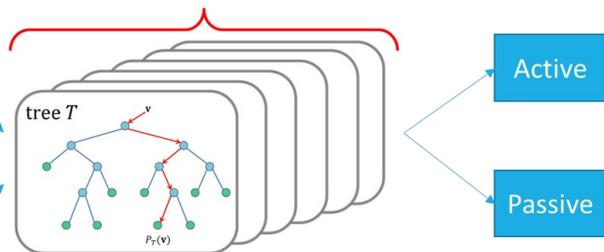
Use egocentric object trajectories to distinguish passive from next-active-objects (i.e., those which will be used soon by the user).

Active Trajectory



Passive Trajectory

Random Decision Forest



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

prediction

bbox =
[1391,101,531,713]

noun = *wooden block*

verb = *take*

ttc = 0.75s

score = 0.83

Last observed frame (V_t)

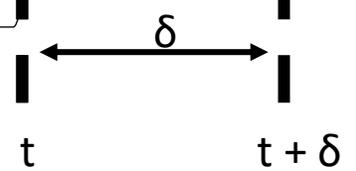


Unobserved future frame ($V_{t+\delta}$)

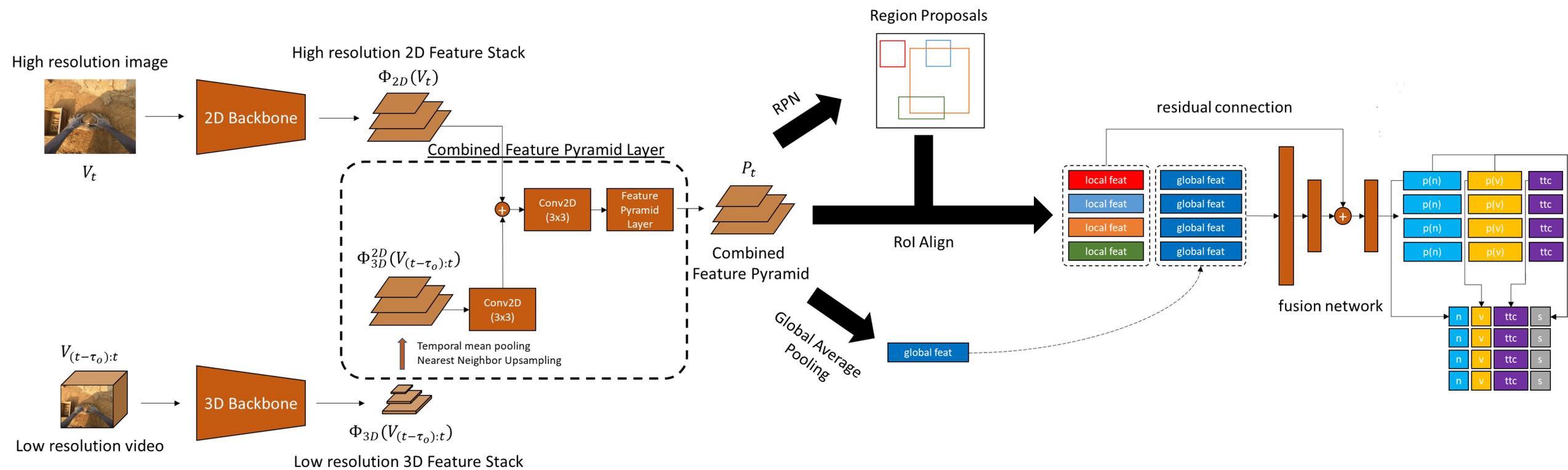


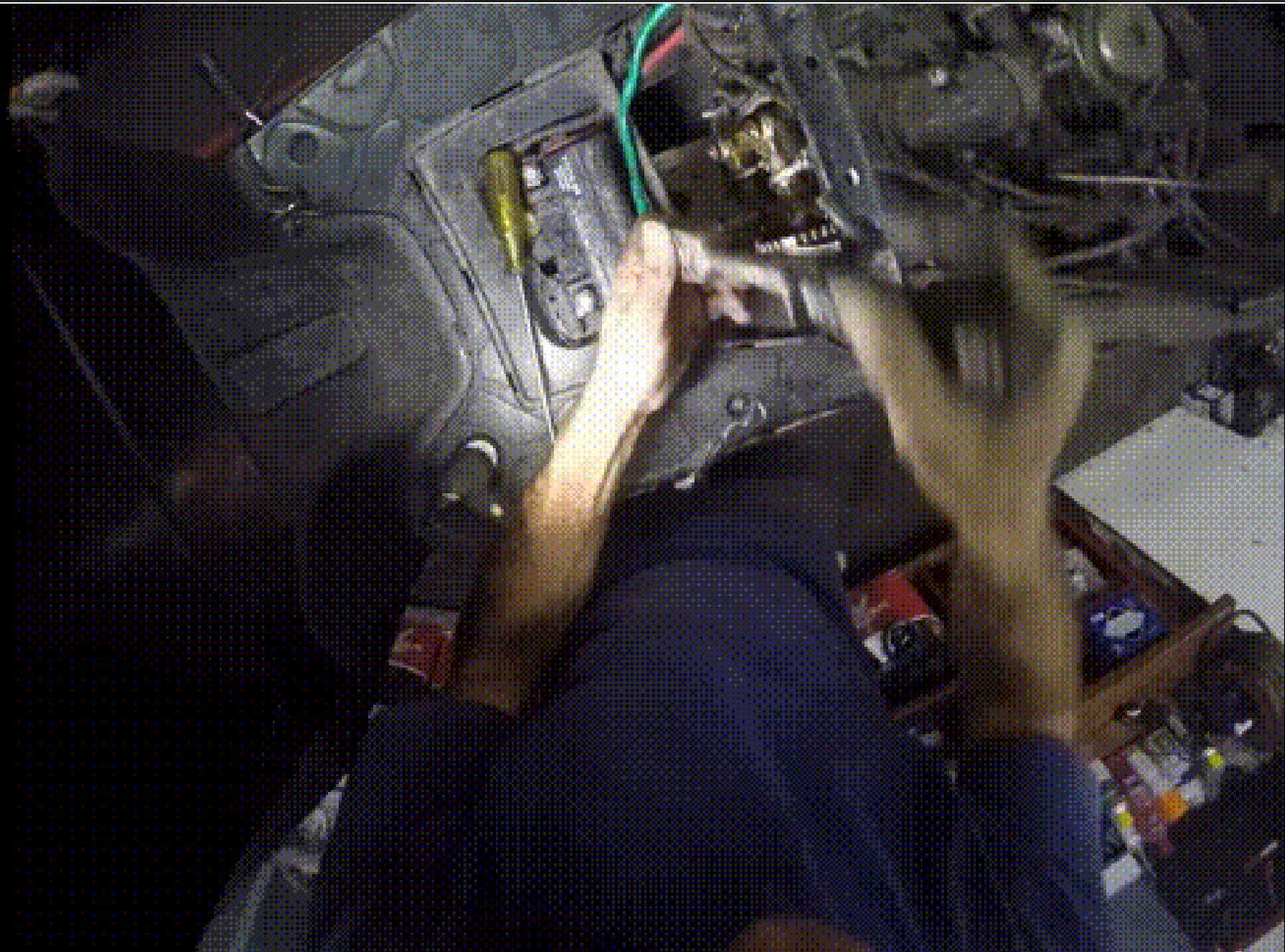
frame of contact

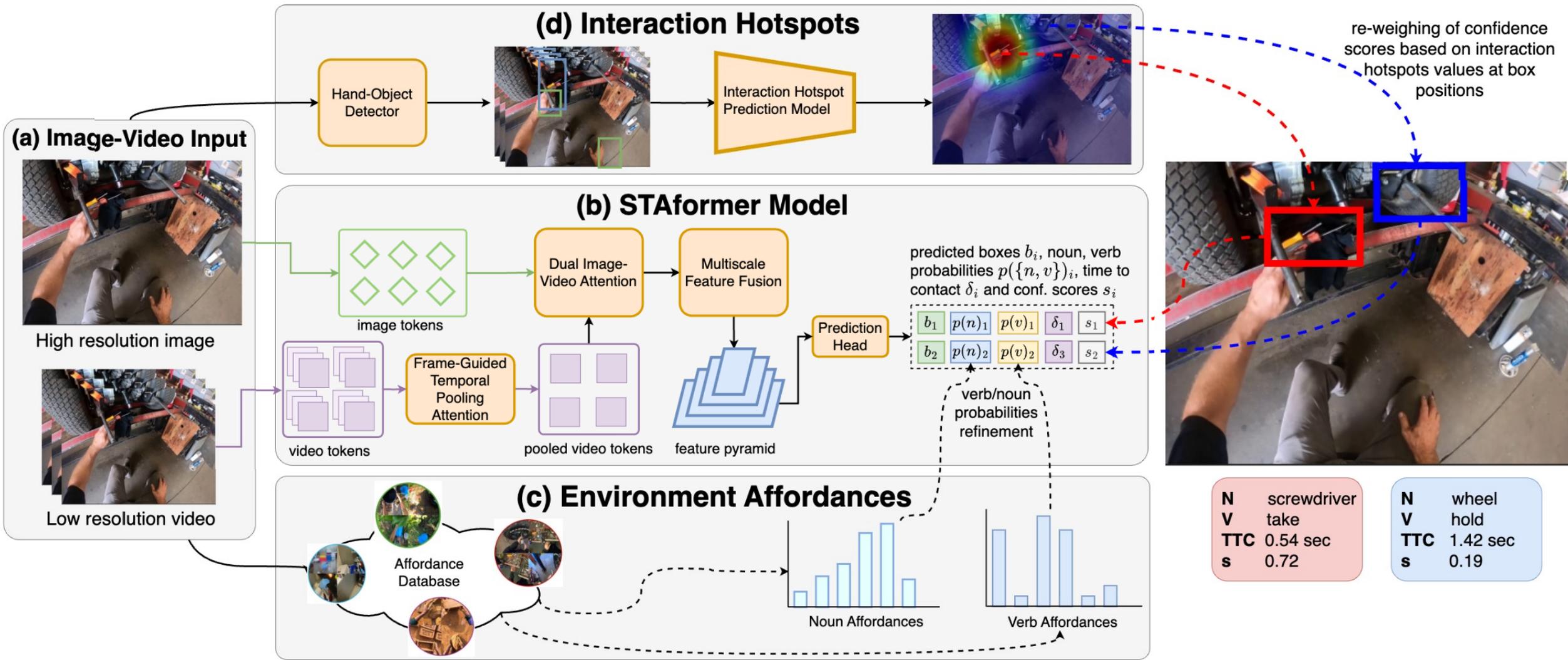
Input video: $V_{:t}$



An end-to-end approach for predicting next-active-objects based on an 2D-3D backbone taking as input a high resolution image and a video clip.







Setup:

- Model: LLaVA-OneVision
- Anticipation time: 0.25s

Standard answer: **The washing soup.**

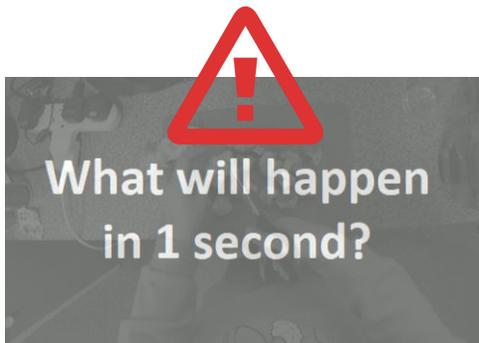
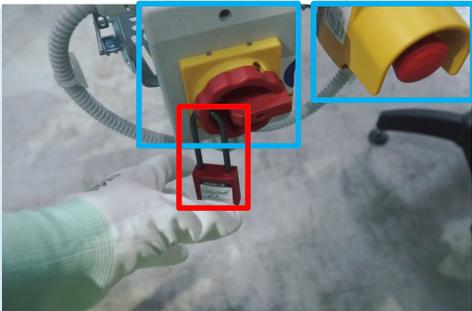
SoM answer: **The spatula.**

Gaze answer: **The small bowl.**

SoM_Gaze answer: **The spatula.**



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;

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

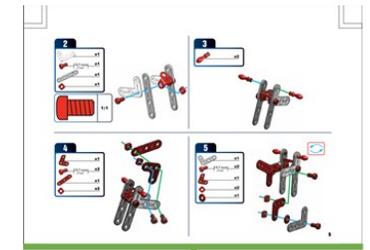
We asked subjects to record egocentric videos while assembling a toy motorbike.

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



	A003 x2		A045 x2		A622 x2
	A004 x2		A046 x2		A632 x1
	A123 x1		A053 x1		A632 x2
	A306 x2		A057 x4		B823 x1
	A050 x2		A077 x2		B577 x2
	A054 x1 34,6 mm 1 3/8"		C658 x10		A090(MJX0200) x1
	A051 x8 18,5 mm 47/64"		A545 x2		A095(J0095) x1

COMPONENTS

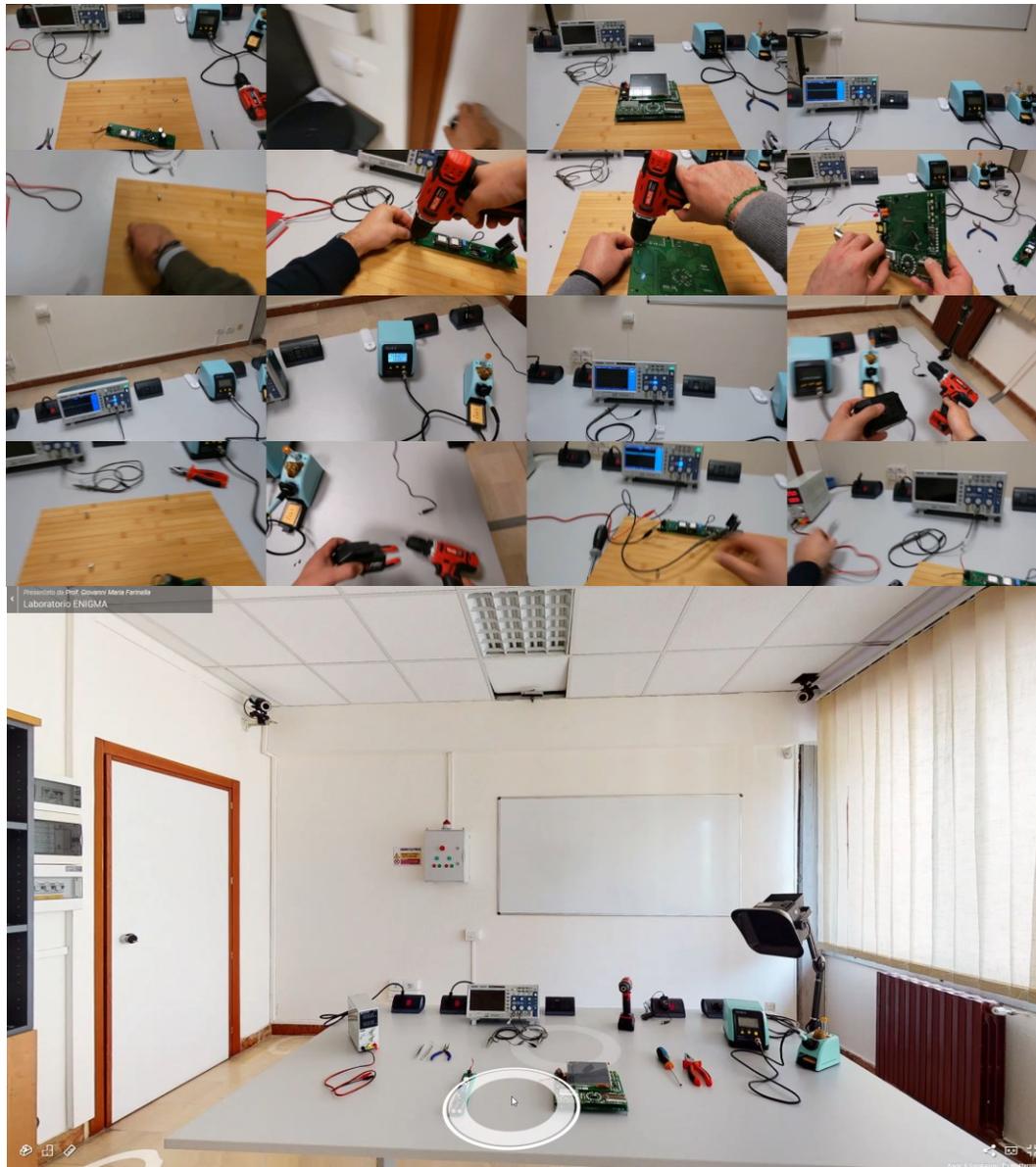


BOOKLET



TOOLS

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

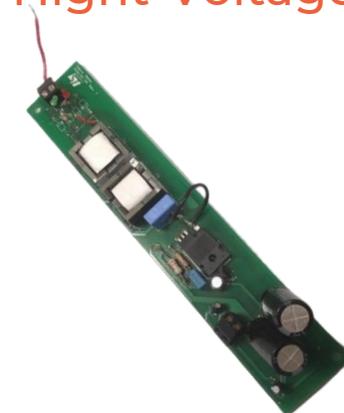


We designed two procedures consisting of instructions that involve humans interacting with the objects present in the laboratory to achieve the goal of repairing two electrical boards

Low-Voltage



High-Voltage



Industrial Applications

NEXT VISION

Spin-off of the University of Catania

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





Intelligent Navigation



Image-based Localization



Augmented Reality



Multi-platform

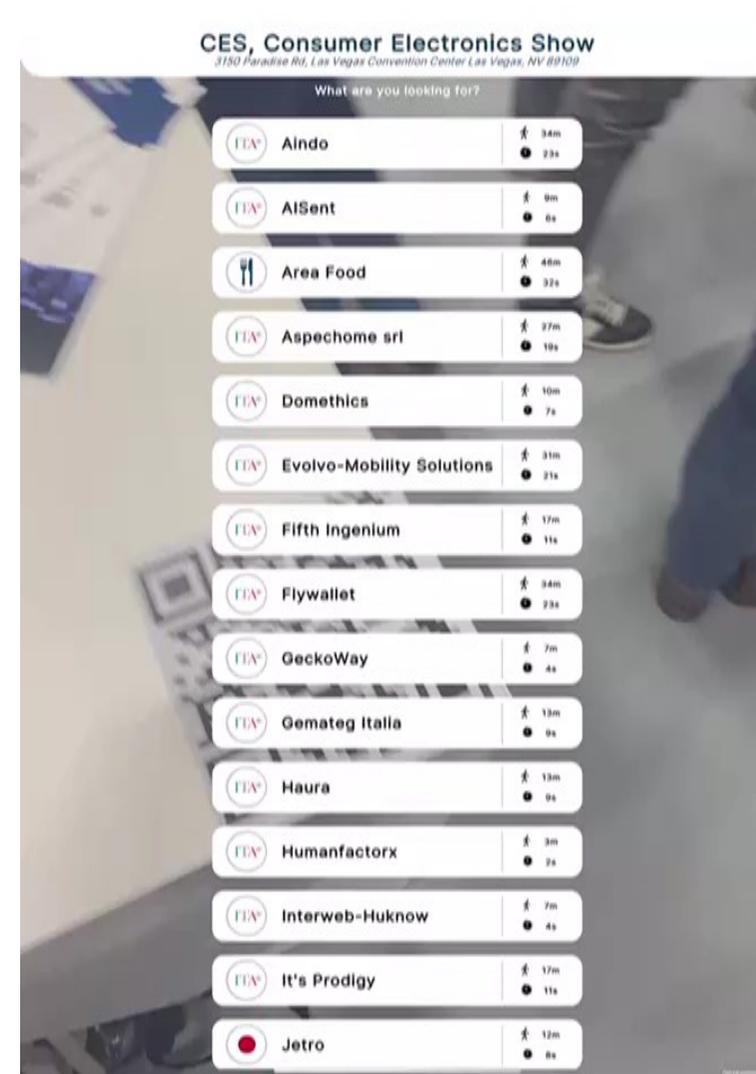


Founders of Next Vision are authors of patents related to the developed technologies





https://drive.google.com/file/d/1Ile4yF6b1kLp9P3ywqKOi77koTvn5OuE/view?usp=share_link



https://drive.google.com/file/d/1FAkLceBzwCkDCsAJFq-nYBwFPZVciQV/view?usp=drive_link



Università di Catania



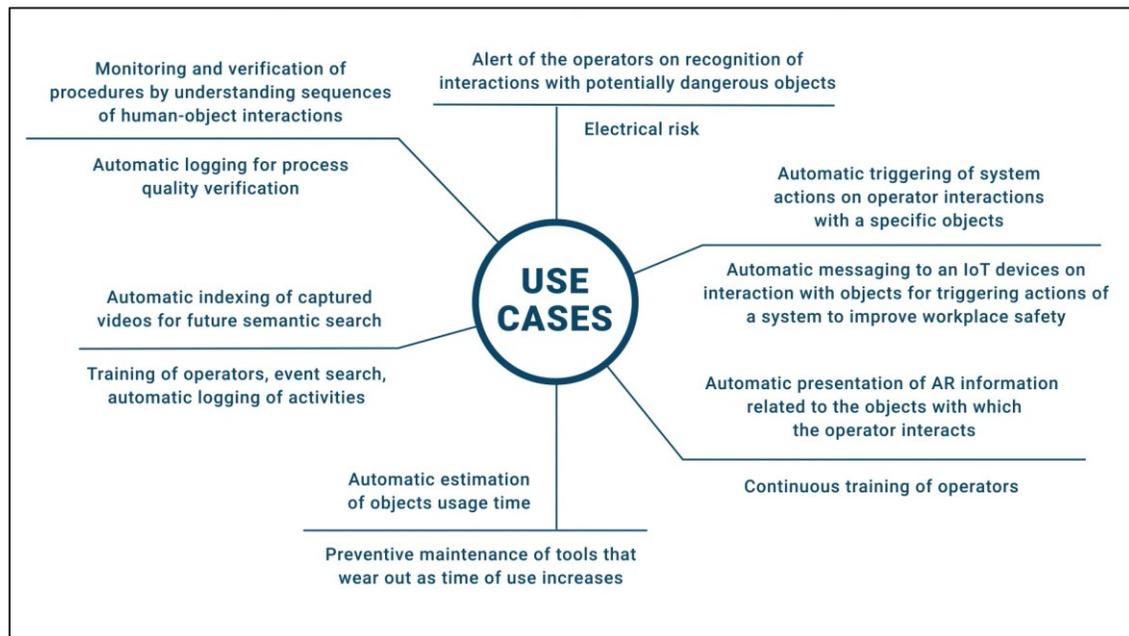




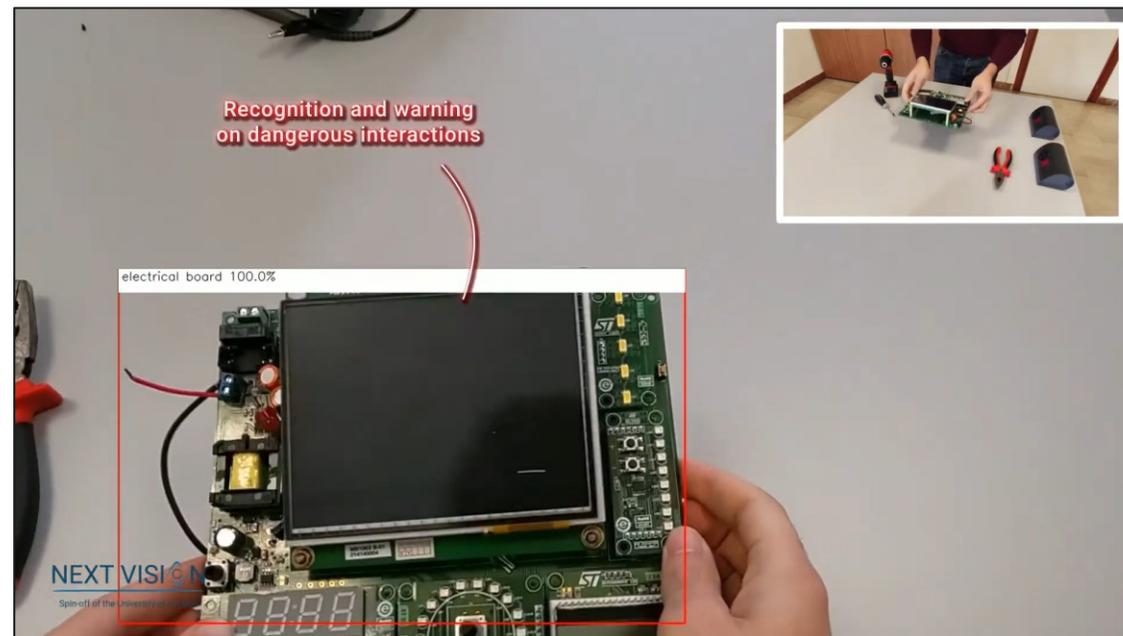


Michele Mazzamuto, Francesco Ragusa, Antonino Furnari, Irene D'Ambra, Antonia Guarriera, Armando Sorbello, Giovanni Maria Farinella (2024). A Mixed Reality Application to Help Impaired People Rehabilitate Outside Clinical Environments. In IEEE International Conference on Metrology for eXtended Reality, Artificial Intelligence and Neural Engineering (MetroXRINE).

- **NAOMI** is an AI Assistant able to support humans to monitor interactions, predict/anticipate next interactions, verify correctness in a sequence of interactions.



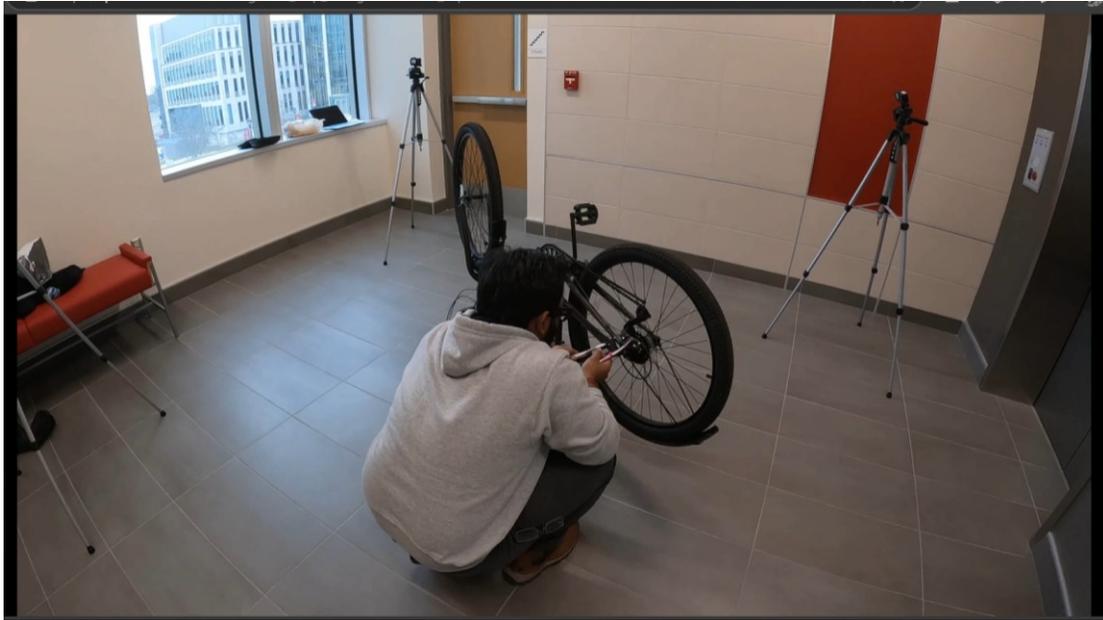
Use cases



The video shows an example of object interaction monitoring. The operator is notified on an interaction with a dangerous object.

https://drive.google.com/file/d/1oOvhVbbyR7AZ35I-V90Zy7RyRTR7IkD4/view?usp=drive_link

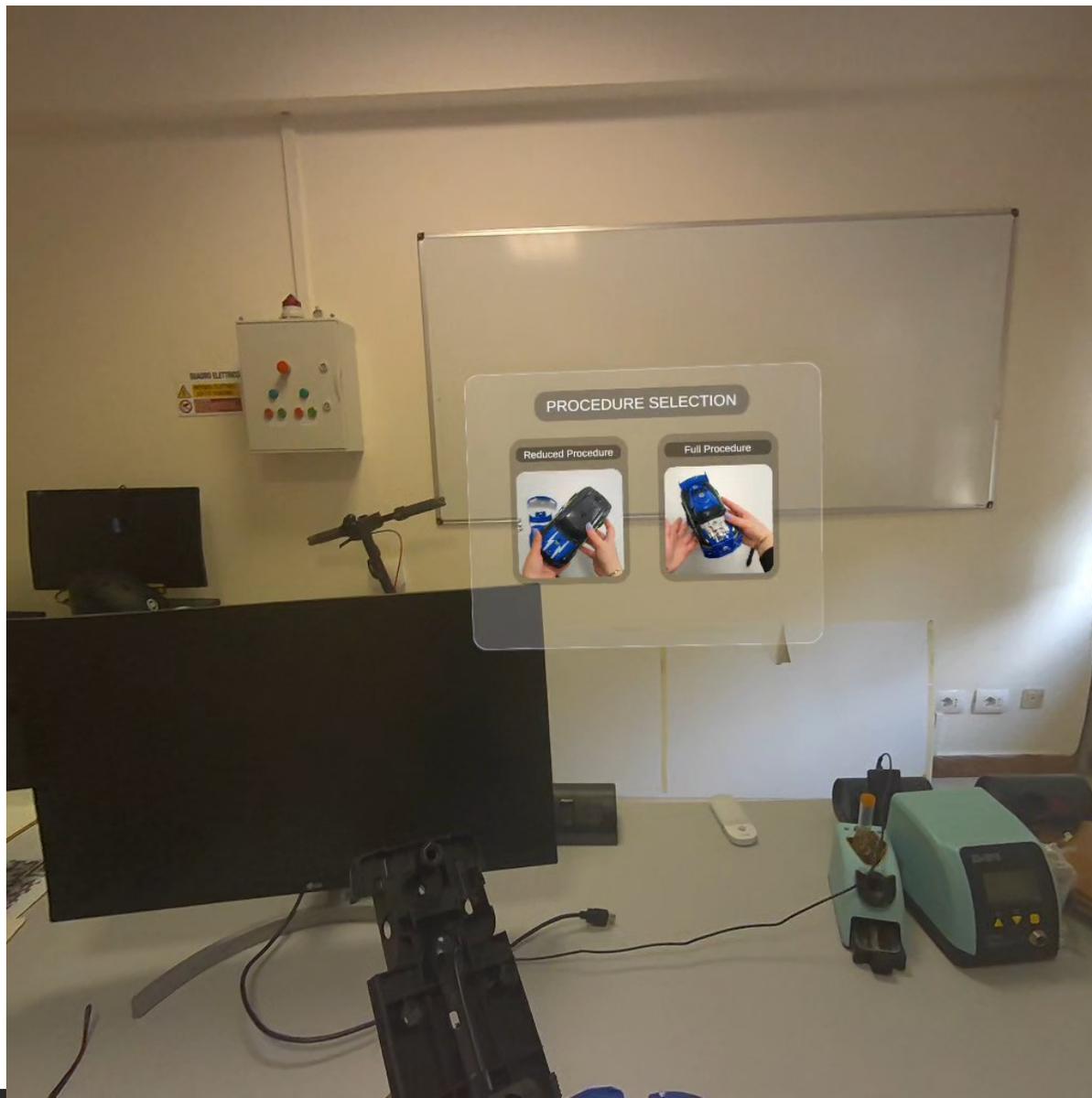
Skill Assessment



Beginner



Expert



The screenshot displays a skill assessment interface. On the left, a sidebar lists two assembly steps:

- Engine assembly**
 - Duration: 00:09
 - Cumulative: 00:09
- Body assembly**
 - Duration: 00:14
 - Cumulative: 00:23

The main video player shows a top-down view of a blue and black toy car on a white surface. To the left of the car are several blue plastic parts and a small metal component. A hand with blue nail polish is visible at the bottom right of the frame. At the bottom of the video player, there are icons for a refresh button, an information button, and a 'Replay' button.

Step

Video



Engine assembly
Duration: 00:09
Cumulative: 00:09

Body assembly
Duration: 00:14
Cumulative: 00:23

Replay

Step

Video



Engine assembly
Duration: 00:09
Cumulative: 00:09

Body assembly
Duration: 00:14
Cumulative: 00:23

Replay

Step

Video



Engine assembly
Duration: 00:09
Cumulative: 00:09

Body assembly
Duration: 00:14
Cumulative: 00:23

Replay

Step

Video

The interface displays the following information:

- Left Sidebar:**
 - Engine assembly**
 - Duration: 00:09
 - Cumulative: 00:09
 - Body assembly**
 - Duration: 00:14
 - Cumulative: 00:23
- Central Video Player:** Shows a top-down view of a blue and black toy car on a white surface. To the left of the car are several blue plastic parts and a small metal component. A hand is visible at the bottom right corner of the video frame. Below the video are icons for a refresh button, an information button, and a 'Replay' button.
- Right Sidebar:**
 - Task name: **Wheel assembly**
 - Step indicator: **3/3 STEP** (with left and right navigation arrows)
 - Timer: **03:06**

Step

Video

Procedure

The interface is divided into three main sections:

- Left Panel (Task List):**
 - Engine assembly:** Duration: 00:09, Cumulative: 00:09
 - Body assembly:** Duration: 00:14, Cumulative: 00:23
- Center Panel (Video):** Shows a top-down view of a car chassis with blue body panels and a black steering wheel. A hand is visible at the bottom right. Below the video are icons for refresh, info, and a 'Replay' button.
- Right Panel (Procedure):** Features a blue circle around the title 'Wheel assembly', a large circular progress indicator showing '3/3 STEP', and a timer at the bottom displaying '03:06'.

Step

Video

Procedure

Engine assembly
Duration: 00:09
Cumulative: 00:09

Body assembly
Duration: 00:14
Cumulative: 00:23

Wheel assembly

3/3
STEP

03:06

Replay

Step

Video

Procedure

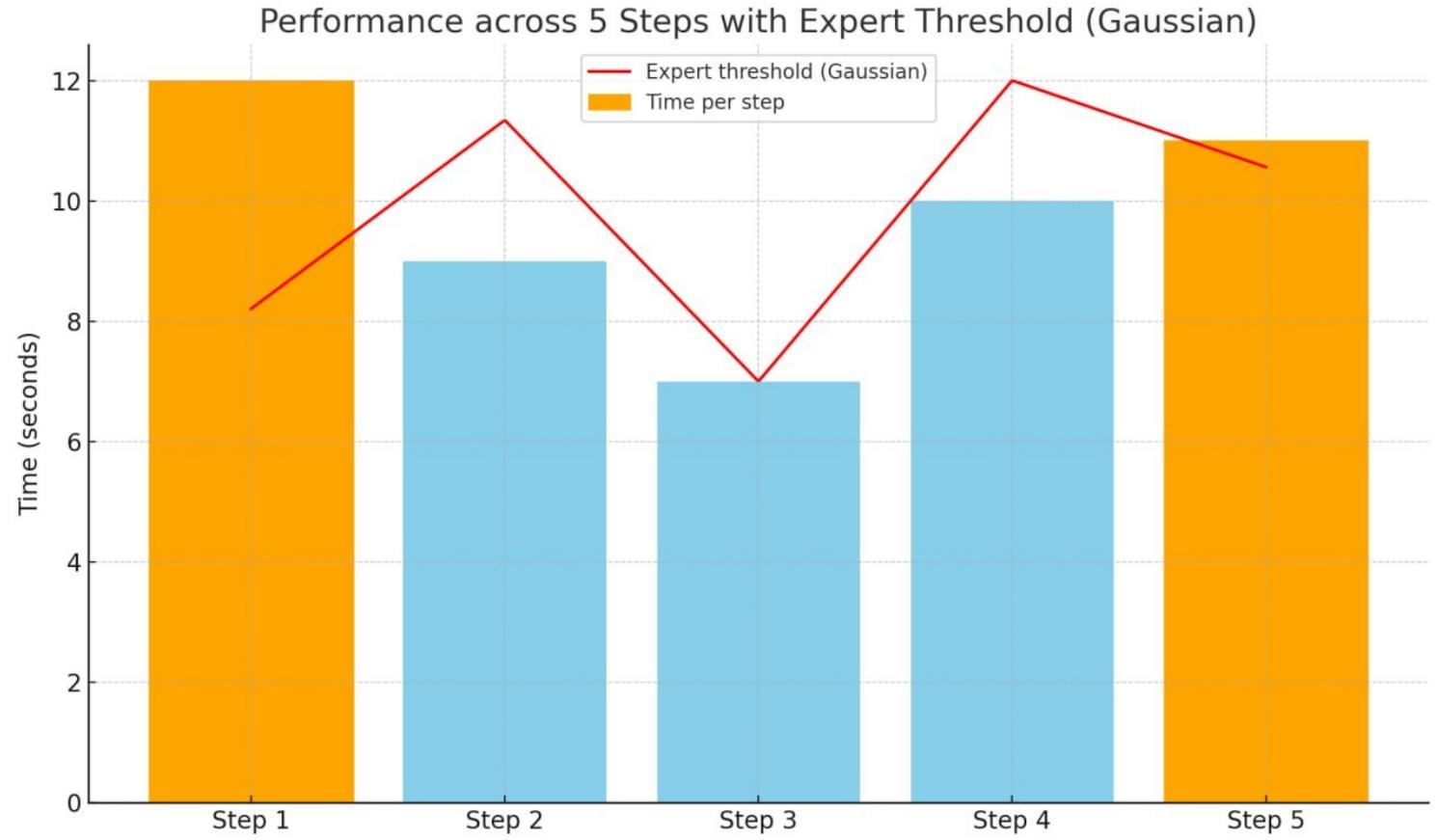
The interface displays the following information:

- Left Sidebar:**
 - Engine assembly**
 - Duration: 00:09
 - Cumulative: 00:09
 - Body assembly**
 - Duration: 00:14
 - Cumulative: 00:23
- Central Video Player:** Shows a top-down view of a blue and black toy car with various parts (blue plastic pieces, a metal component, a black tool) laid out on a white surface. A hand is visible at the bottom right. Below the video are icons for refresh, info, and a 'Replay' button.
- Right Sidebar:**
 - Task name: **Wheel assembly**
 - Step indicator: **3/3 STEP** (with left and right navigation arrows)
 - Timer: **03:06** (circled in blue)

Step

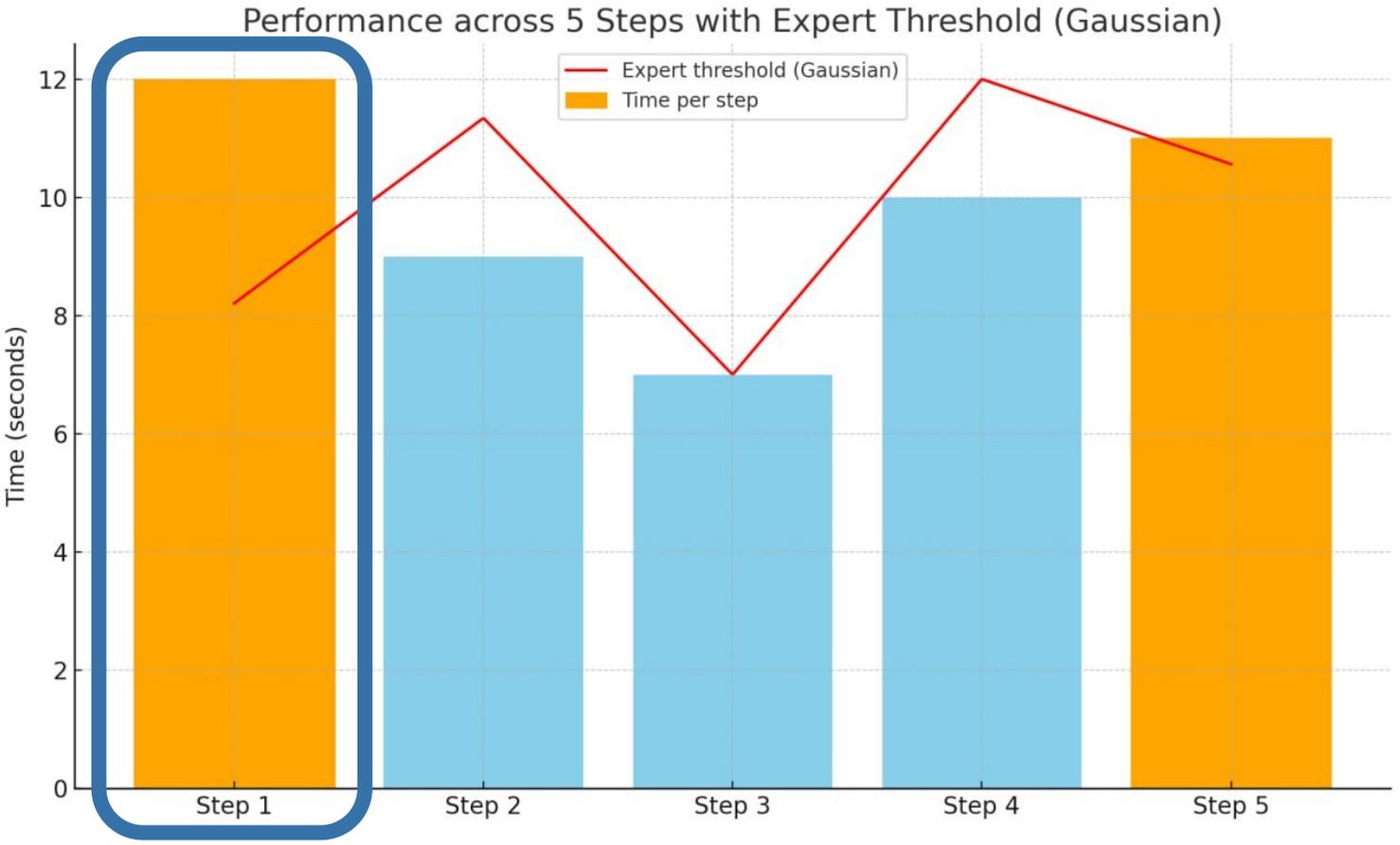
Video

Procedure



Step

Procedure



Step

Procedure

Wheel assembly

3/3 STEP

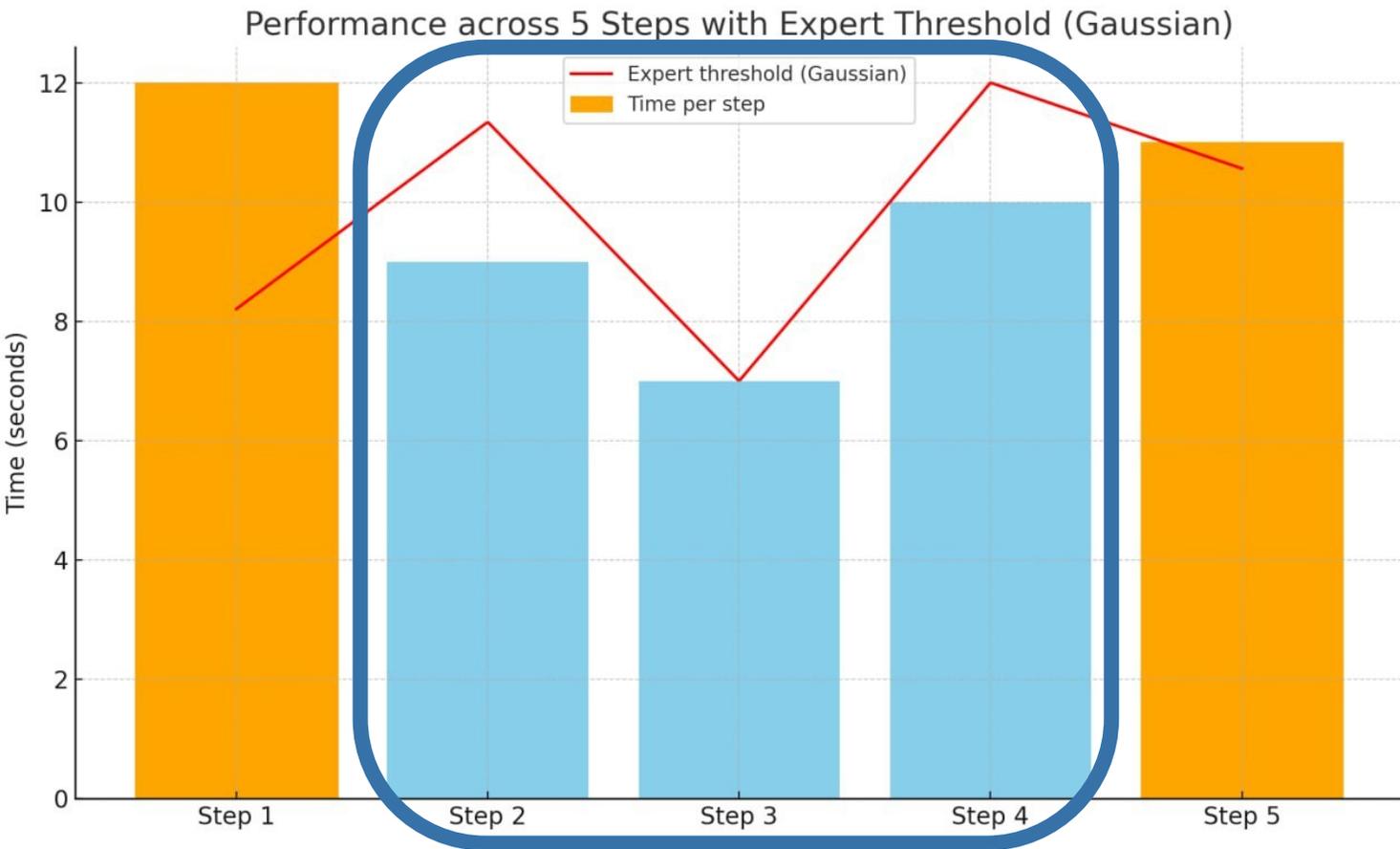
03:06

Engine assembly

Duration: 00:09
Cumulative: 00:09

Body assembly

Duration: 00:14
Cumulative: 00:23



Step

Procedure

Wheel assembly

3/3 STEP

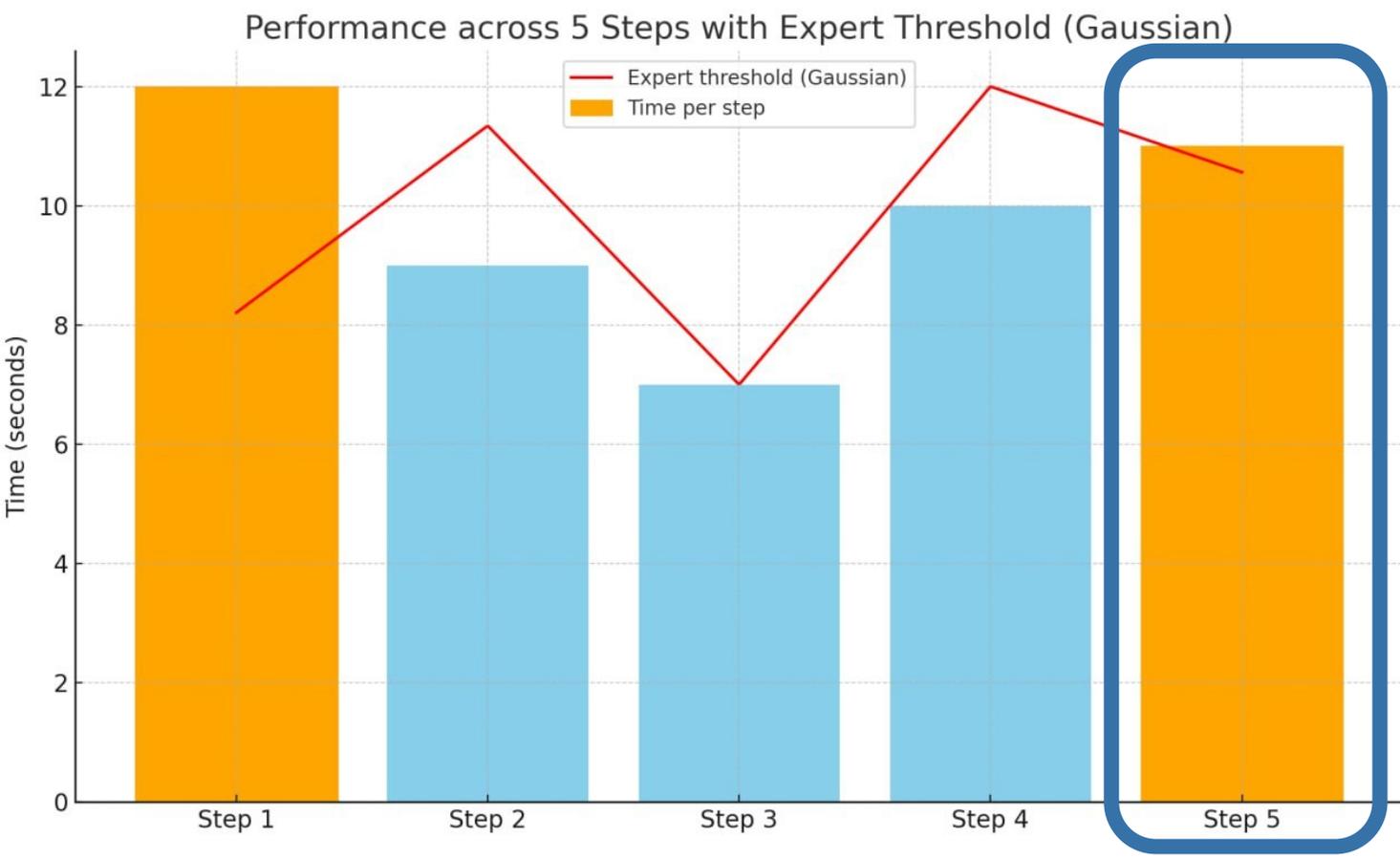
03:06

Engine assembly

Duration: 00:09
Cumulative: 00:09

Body assembly

Duration: 00:14
Cumulative: 00:23



Step

Procedure



Doing Research in Egocentric Vision: Where to start?



Download only certain data types

We provide videos, RGB/optical flow frames, GoPro's metadata (for the extension only) and object detection frames (for EPIC KITCHENS-55's videos only). You can also download the consent form templates.

If you want to download only one (or a subset) of the above, you can do so with the following self-explanatory arguments:

- `--videos`
- `--rgb-frames`
- `--flow-frames`
- `--object-detection-images`
- `--masks`
- `--metadata`
- `--consent-forms`

If you want to download only videos, then:

```
python epic_downloader.py --videos
```

Note that these arguments can be **combined** to download multiple things. For example:

```
python epic_downloader.py --rgb-frames --flow-frames
```

Will download both RGB and optical flow frames.

Specifying participants

You can use the argument `--participants` if you want to download data for only a subset of the participants. Participants can be specified with their numerical or string ID.

You can specify a single participant, e.g. `--participants 1` or `--participants P01` for participant P01, or a comma-separated list of them, e.g. `--participants 1,2,3` or `--participants P01,P02,P03` for participants P01, P02 and P03

This argument can also be combined with the aforementioned arguments. For example:

```
python epic_downloader.py --videos --participants 1,2,3
```

Will download only videos from P01, P02 and P03.

Data download

Canonical videos and annotations can be downloaded using the following command:

```
python -m ego4d.cli.cli --output_directory=~/.ego4d_data" --datasets full_scale annotations --benchmarks FH0
```

v2.0 annotations can be downloaded with:

```
python -m ego4d.cli.cli --output_directory=~/.ego4d_data" --datasets annotations --version v2
```

Detailed Flags

Flag Name	Description
<code>--dataset</code>	[Required] A list of identifiers to download: [annotations, full_scale, clips] Each dataset will be stored in folders in the output directory with the name of the dataset (e.g. <code>output_dir/v2/full_scale/</code>) and manifest.
<code>--output_directory</code>	[Required] A local path where the downloaded files and metadata will be stored
<code>--metadata</code>	[Optional] Download the primary <code>ego4d.json</code> metadata at the top level (Default: True)
<code>--benchmarks</code>	[Optional] A list of benchmarks to filter dataset downloads by - e.g. Narrations/EM/FHO/AV
<code>-y --yes</code>	[Optional] If this flag is set, then the CLI will not show a prompt asking the user to confirm the download. This is so that the tool can be used as part of shell scripts.
<code>--aws_profile_name</code>	[Optional] Defaults to "default". Specifies the AWS profile name from <code>~/.aws/credentials</code> to use for the download
<code>--video_uids</code>	[Optional] List of video or clip UIDs to be downloaded. If not specified, all relevant UIDs will be downloaded.
<code>--video_uid_file</code>	[Optional] Path to a whitespace delimited file that contains a list of UIDs. Mutually exclusive with the <code>video_uids</code> flag.
<code>--universities</code>	[Optional] List of university IDs. If specified, only UIDs from the S3 buckets belonging to the listed universities will be downloaded.
<code>--version</code>	[Optional] A version identifier - e.g. "v1" or "v2" (default)
<code>--no-metadata</code>	[Optional] Bypass the <code>ego4d.json</code> metadata download
<code>--config</code>	[Optional] Local path to a config JSON file. If specified, the flags will be read from this file instead of the command line

Modern datasets are HUGE!

- EPIC-KITCHENS ~ 796 GB
- EGO4D ~ 30+ TB

Datasets

The following datasets are available (not exhaustive):

Dataset	Description
annotations	The full set of annotations for the majority of benchmarks.
full_scale	The full scale version of all videos. (Provide <code>benchmarks</code> or <code>video_uids</code> filters to reduce the 5TB download size.)
clips	Clips available for benchmark training tasks. (Provide <code>benchmarks</code> or <code>video_uids</code> filters to reduce the download size.)
video_540ps	The downsampled version of all videos - rescaled to 540px on the short side. (Provide <code>benchmarks</code> or <code>video_uids</code> filters to reduce the 5TB download size.)
annotations_540ps	The annotations corresponding to the downsampled <code>video_540ps</code> videos - primarily differing only in spatial annotations (e.g. bounding boxes).
3d	Annotations for the 3D VQ benchmark.
3d_scans	3D location scans for the 3D VQ benchmark.
3d_scan_keypoints	3D location scan keypoints for the 3D VQ benchmark.
imu	IMU data for the subset of videos available
slowfast8x8_r101_M400	Precomputed action features for the Slowfast 8x8 (R101) model
omnivore_video_swini	Precomputed action features for the Omnivore Video model
omnivore_image_swini	Precomputed action features for the Omnivore Image model
fut_loc	Images and annotations for the future locomotion benchmark.
av_models	Model checkpoints for the AV/Social benchmark.
lta_models	Model checkpoints for the Long Term Anticipation benchmark.
moments_models	Model checkpoints for the Moments benchmark.
nlq_models	Model checkpoints for the NLQ benchmark.
sta_models	Model checkpoints for the Short Term Anticipation benchmark.
vq2d_models	Model checkpoints for the 2D VQ benchmark.





EPIC-KITCHENS-100 2025 CHALLENGES

Challenge Details with links to ★NEW★ Codalab Leaderboards

leaderboards are now open for the challenge phase from Jan 2025.

In 2024, we have 7 open challenges. These are

- [Semi-Supervised Video Object Segmentation Challenge](#)
- [EPIC-SOUNDS Audio-Based Interaction Recognition](#)
- [EPIC-SOUNDS Audio-Based Interaction Detection](#)
- [Action Recognition](#)
- [Action Detection](#)
- [UDA for Action Recognition](#)
- [Multi-Instance Retrieval](#)

EPIC-Kitchens 2025 Challenges

Feb 1st 2025,

May 19th 2025,

May 23rd 2025,

June 17 2024,

All leaderboards are open

Server Submission Deadline at 00:00:00 UTC

Deadline for Submission of Technical Reports on CMT [HERE](#)

Results announced at 2nd EgoVis workshop in Nashville [EgoVis@CVPR2025 workshop](#)

Challenges Guidelines

The eight challenges below and their test sets and evaluation servers are available via CodaLab. The leaderboards will decide the winners for each individual challenge. For each challenge, the CodaLab server page details submission format and evaluation metrics.

This year, we offer four new challenges in: Semi-Supervised Video Object Segmentation using the [VISOR](#) annotations, Hand-object-segmentations using the [VISOR](#) annotations, single-object tracking and audio-based action recognition using the [epic-sounds](#) dataset.

To enter any of the nine competitions, you need to register an account for that challenge using a valid institute (university/company) email address. To enable your account, [fill this form with your team's details](#). A single registration per research team is allowed. We perform a manual check for

<https://epic-kitchens.github.io/2025>

Ego4D and EgoExo4D Challenge 2024

Overview

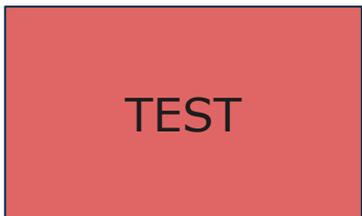
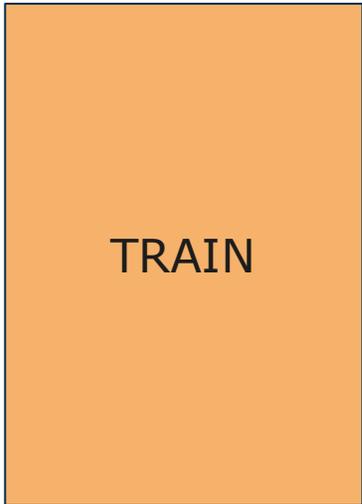
At CVPR 2024, we will host **16** challenges including 2 new challenges (Goal Step and Ego Schema), representing each of Ego4D's five benchmarks. Included in the 16 challenges hosted at CVPR are two teaser Ego-Exo4D challenges (Ego-Pose Body and Ego-Pose Hands). Please find details below on the challenges:

Ego4D challenges

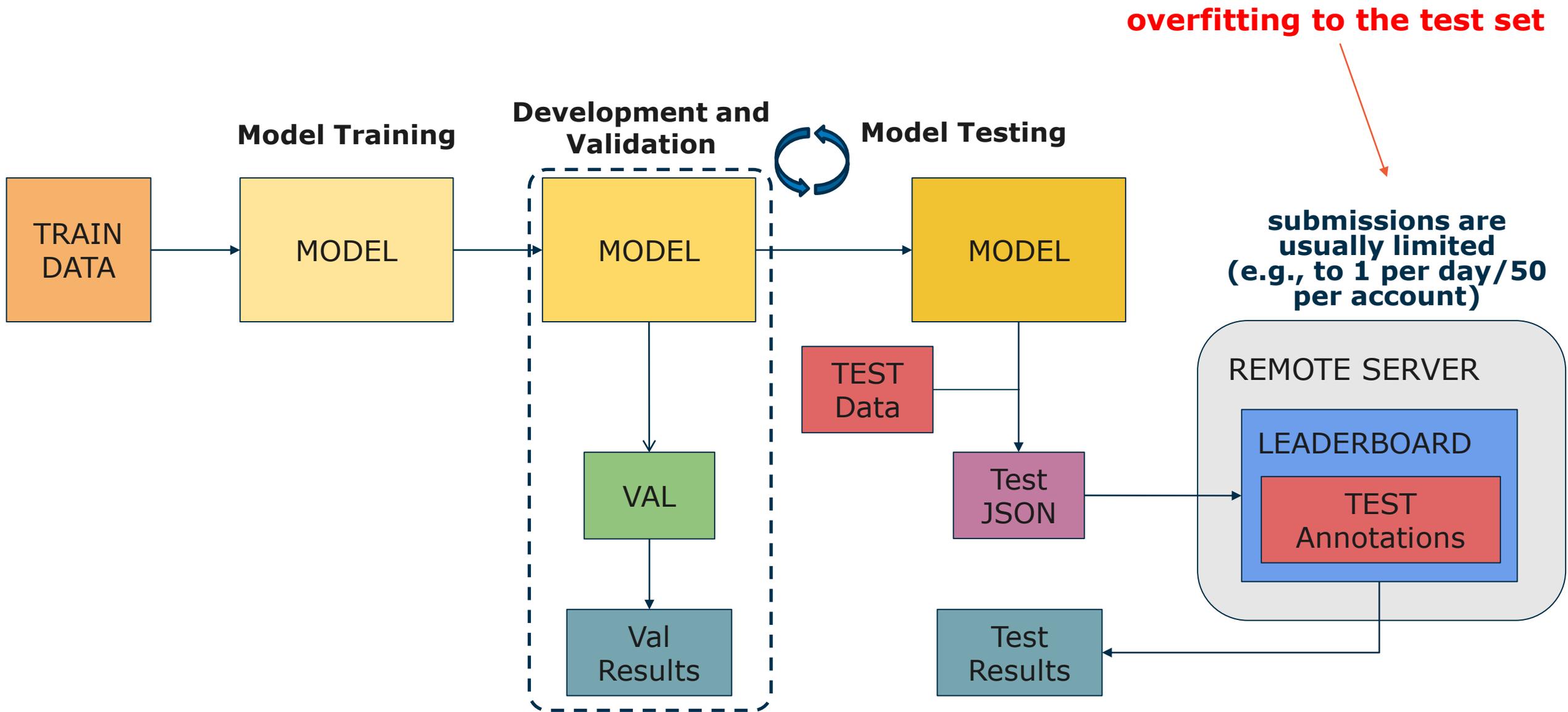
Episodic memory:

- [Visual queries with 2D localization \(VQ2D\)](#) and [Visual Queries 3D localization \(VQ3D\)](#): Given an egocentric video clip and an image crop depicting the query object, return the most recent occurrence of the object in the input video, in terms of contiguous bounding boxes (2D + temporal localization) or the 3D

<https://ego4d-data.org/docs/challenge/>



- Datasets are usually divided into train/val/test splits;
- All videos are publicly released;
- Train annotations are publicly released and meant for training models for the different challenges;
- Val annotations are publicly released and meant for model development and hyperparameter search;
- Test annotations are private and meant for assessing the performance of models avoiding bias in model design and optimization;
- Hence, the only way to obtain results on the test set is to send model predictions to an evaluation server.





EPIC-KITCHENS-100 Action Detection

Organized by antonino - Current server time: Feb. 23, 2025, 7:28 p.m. UTC

First phase

End

2024 Open Testing

Competition Ends

Oct. 1, 2024, 8 a.m. UTC

Jan. 31, 2025, midnight UTC

Test Set (Mean Average Precision - mAP)

#	User	Entries	Date of Last Entry	Team Name	mAP@0.1 (%)			mAP@0.2 (%)			mAP@0.3 (%)			mAP@0.4 (%)			mAP@0.5 (%)			Avg. mAP (%)					
					PT	TL	TD	Verb	Noun	Action	Verb	Noun	Action	Verb	Noun	Action									
1	shuming	12	05/02/24	KAUST-4Paradigm-MoonshotAI-Nvidia	2.0	3.0	4.0	34.11	40.66	36.09	32.75	38.63	34.70	30.48	36.32	32.67	28.03	32.55	29.91	24.73	27.98	26.50	30.02	35.23	31.97
2	xyx	23	05/30/24	dg_team(deepglint)	3.0	3.0	4.0	30.28	33.73	29.29	29.27	32.39	28.33	27.32	30.44	26.80	25.32	27.48	24.77	22.14	23.74	22.06	26.87	29.56	26.25
3	TIM_method	1	04/06/24	Oxford+Bristol	2.0	3.0	3.0	32.14	34.88	28.13	30.01	32.99	26.74	27.84	30.57	25.01	25.24	26.60	22.29	20.37	21.78	18.86	27.12	29.36	24.21
4	mzs	5	05/27/23	mzs	2.0	3.0	4.0	31.01	30.32	25.54	30.04	28.76	24.54	28.01	27.20	23.16	25.44	24.28	21.04	22.32	20.74	18.35	27.36	26.26	22.52
5	lijun	18	06/01/22		2.0	3.0	4.0	30.67	30.96	24.57	29.40	29.36	23.50	26.81	26.78	21.94	24.34	23.27	19.65	20.51	18.80	16.74	26.35	25.83	21.28
6	tzzcl	19	06/01/22	4Paradigm-UWMadison-NJU	2.0	3.0	4.0	26.97	28.61	23.90	25.91	27.14	22.98	24.21	24.92	21.37	21.77	22.14	19.57	18.47	18.69	16.94	23.47	24.30	20.95
7	Haniel	44	06/01/23	Bristol-MaVi	2.0	3.0	3.0	27.57	24.18	19.64	26.16	22.92	18.66	24.02	20.85	17.29	21.48	18.24	15.55	18.42	15.21	13.49	23.53	20.28	16.93
8	Alibaba-MMAL-Research	1	12/15/21	CVPR 2021 Challenges	2.0	3.0	3.0	22.77	26.44	18.76	22.01	24.55	17.73	19.63	22.30	16.26	17.81	19.82	14.91	14.65	16.25	12.87	19.37	21.87	16.11
9	ctai	7	05/30/23	ctai	2.0	3.0	3.0	23.37	21.21	17.09	22.57	20.22	16.49	21.67	19.00	15.71	19.51	17.01	14.24	16.87	13.94	12.38	20.80	18.28	15.18



Ego4D Short Term Object Interaction Anticipation Challenge

★ 15

Organized by: Ego4D

Starts on: Oct 25, 2022 2:00:00 AM CET (GMT + 1:00)

Ends on: Jun 1, 2099 1:59:59 AM CET (GMT + 1:00)

- Overview
- Evaluation
- Phases
- Participate
- Leaderboard
- Discuss

Leaderboard

Overall Top-5 mAP

Phase: Test Phase, Split: Test Split

Order by metric

B - Baseline * - Private V - Verified

Visible Metrics

Rank	Participant team	Noun (↑)	Noun_Verb (↑)	Noun_TTC (↑)	Overall (↑)	Last submission at	Meta Attributes
1	123456ABCD (j)	31.08	16.18	12.41	7.21	8 months ago	View
2	Zarrio (IH_new)	33.50	17.26	11.77	6.75	8 months ago	View
3	ICL@SNU (YOLO + CLIP)	34.89	17.61	10.91	6.22	9 months ago	View
4	Language NAO (TransFusion, Ego4D v2)	30.09	13.58	10.39	5.41	1 year ago	View
5	PAVIS (GANO_v2)	25.67	13.60	9.02	5.16	2 years ago	View

- 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.

VISAPP 26-1B

Monday, March 9th

15:00 - 16:00

proposed data generation framework and GlovEgo-Net.

<https://nextvisionlab.github.io/GlovEgo-HOI/>

Alfio Spoto, Rosario Leonardi, Francesco Ragusa, Giovanni Maria Farinella (2026). GlovEgo-HOI: Bridging the Synthetic-to-Real Gap for Industrial Egocentric Human-Object Interaction Detection. In International Conference on Computer Vision Theory and Applications (VISAPP).

SignIT

VISAPP 26-40

Tuesday, March 10th

11:00 - 12:00

 Paper

 Dataset

<https://fpv-iplab.github.io/SignIT/>

Alessia Micieli, Giovanni Maria Farinella, Francesco Ragusa (2026). SignIT: A Comprehensive Dataset and Multimodal Analysis for Italian Sign Language Recognition. In International Conference on Computer Vision Theory and Applications (VISAPP).

EGO

EXO

VISAPP 26-70

Wednesday, March 11th

10:45 - 12:00

Alessandro Passanisi, Giovanni Maria Farinella, Francesco Ragusa (2026). Ego-Exo Temporal Action Segmentation in Industrial Environments. In International Conference on Computer Vision Theory and Applications (VISAPP).

francesco.ragusa@unict.it – fragusa@nextvisionab.it





Università
di Catania

NEXT VISION

Spin-off of the University of Catania



THANK YOU!

Seeing Through the User's Eyes: Advances in Human-Centric Egocentric Vision

Francesco Ragusa

LIVE Group @ UNICT - <https://iplab.dmi.unict.it/live/>

Next Vision - <http://www.nextvisionlab.it/>

Department of Mathematics and Computer Science - University of Catania

francesco.ragusa@unict.it - <https://francescoragusa.github.io/>



VISAPP 2026

21st International Conference on Computer Vision
Theory and Applications

Marbella, Spain 9 - 11 March, 2026