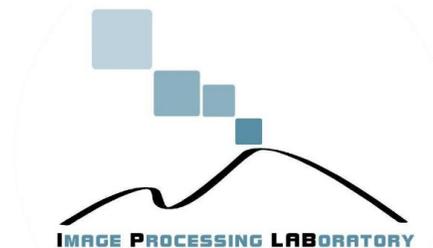




Università
di Catania

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/>



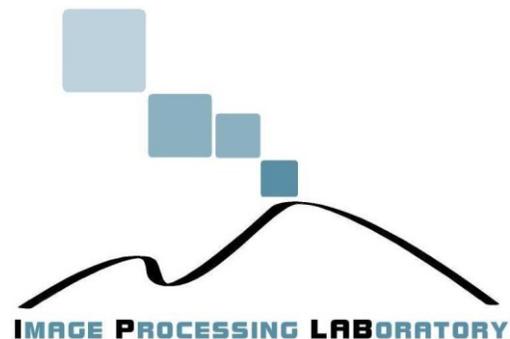
VISAPP 2026

21st International Conference on Computer Vision
Theory and Applications

Marbella, Spain 9 - 11 March, 2026



Università di Catania



LIVE Group @ UNICT



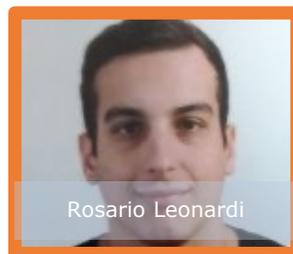
Giovanni Maria Farinella



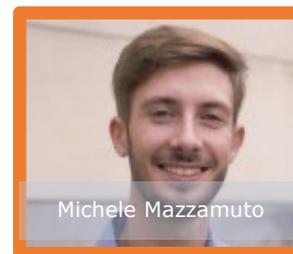
Francesco Ragusa



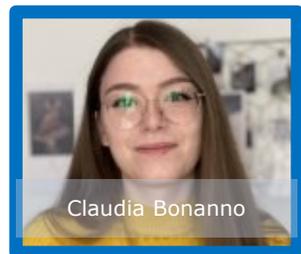
Daniele Di Mauro



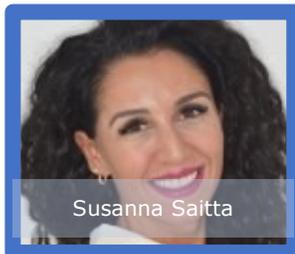
Rosario Leonardi



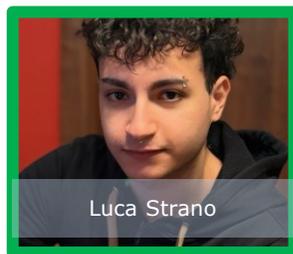
Michele Mazzamuto



Claudia Bonanno



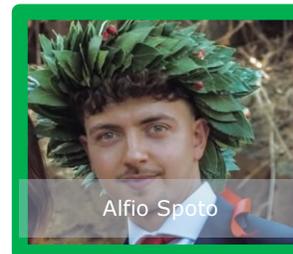
Susanna Saitta



Luca Strano



Giovanni Maria Manduca



Alfio Spoto



Alessia Micieli



Salvatore Carota



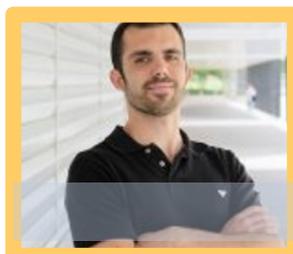
Alessandro Passanisi



Daniele Materia



Irene D'Ambra



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

NEXT VISION

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

19 Members

1 Full Professor

1 Assistant Professor

3 Post Docs

2 PhD Students

7 Master Students

1 Lab Assistant

4 Visiting PhD Students

The slides of this tutorial are available online at:
<https://francescoragusa.github.io/visapp2026>



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 1

History and Motivations



Perception and Egocentric Vision

Perception is the process of **receiving**, **organize** and **interpret** information in order to give meaning to the surrounding world.

Perception is the process of **receiving**, **organize** and **interpret** information in order to give meaning to the surrounding world.

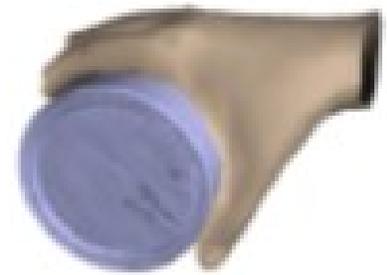
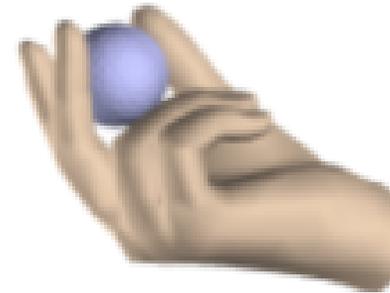
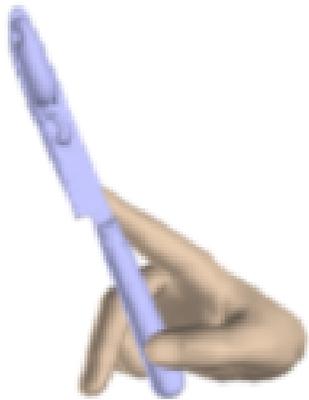


Perception is the process of **receiving**, **organize** and **interpret** information in order to give meaning to the surrounding world.

I'm in the
kitchen!



Perception is the process of **receiving**, **organize** and **interpret** information in order to give meaning to the surrounding world.



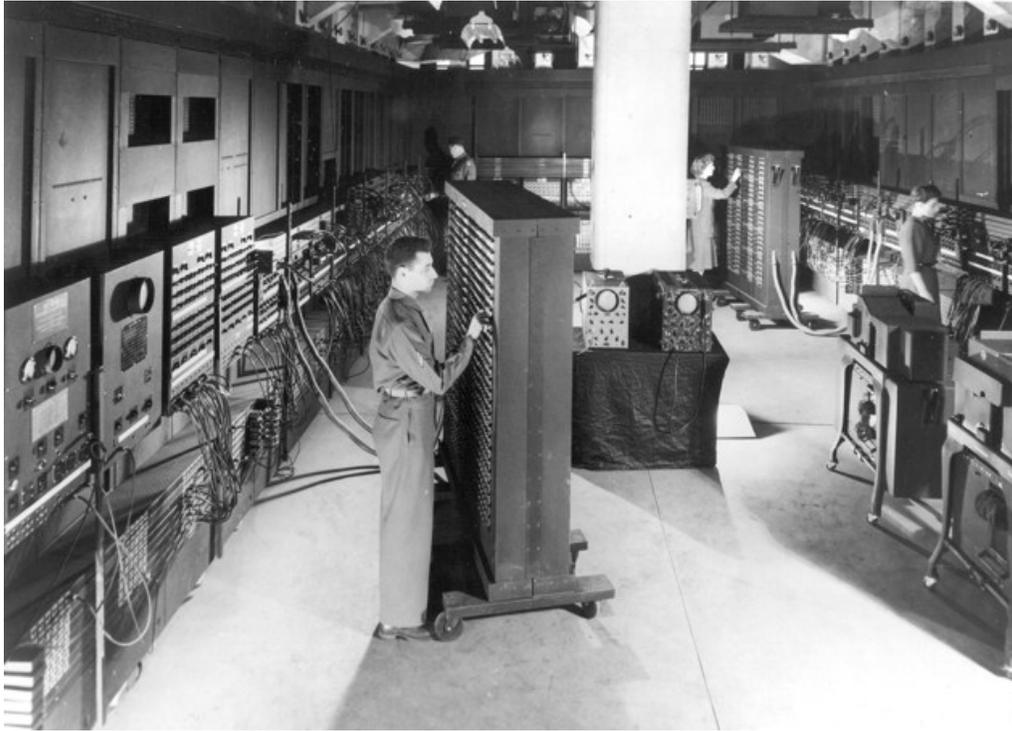
Perception is the process of **receiving**, **organize** and **interpret** information in order to give meaning to the surrounding world.



Computer vision enables computers to **acquire, process, analyze** and **understand** digital images, and extract of high-dimensional data from the real world in order to produce numerical or symbolic information

Computer vision enables computers to **acquire**, **process**, **analyze** and **understand** digital images, and extract of high-dimensional data from the real world in order to produce numerical or symbolic information, e.g. in the forms of decisions

Perception is the process of **receiving**, **organize** and **interpret** information in order to give meaning to the surrounding world.



Mainframe Era (1950s–1970s)

Centralized, inaccessible, institutional



Personal Computer Era (1980s–1990s)

Desktop computing enters homes and offices



Laptop Era (1990s–2000s)

*Computing for the mass, but not mobile
and not context aware - dedicated
access to computing*

Smartphone Era

*Computing is
context aware - forces to switch*

Short News

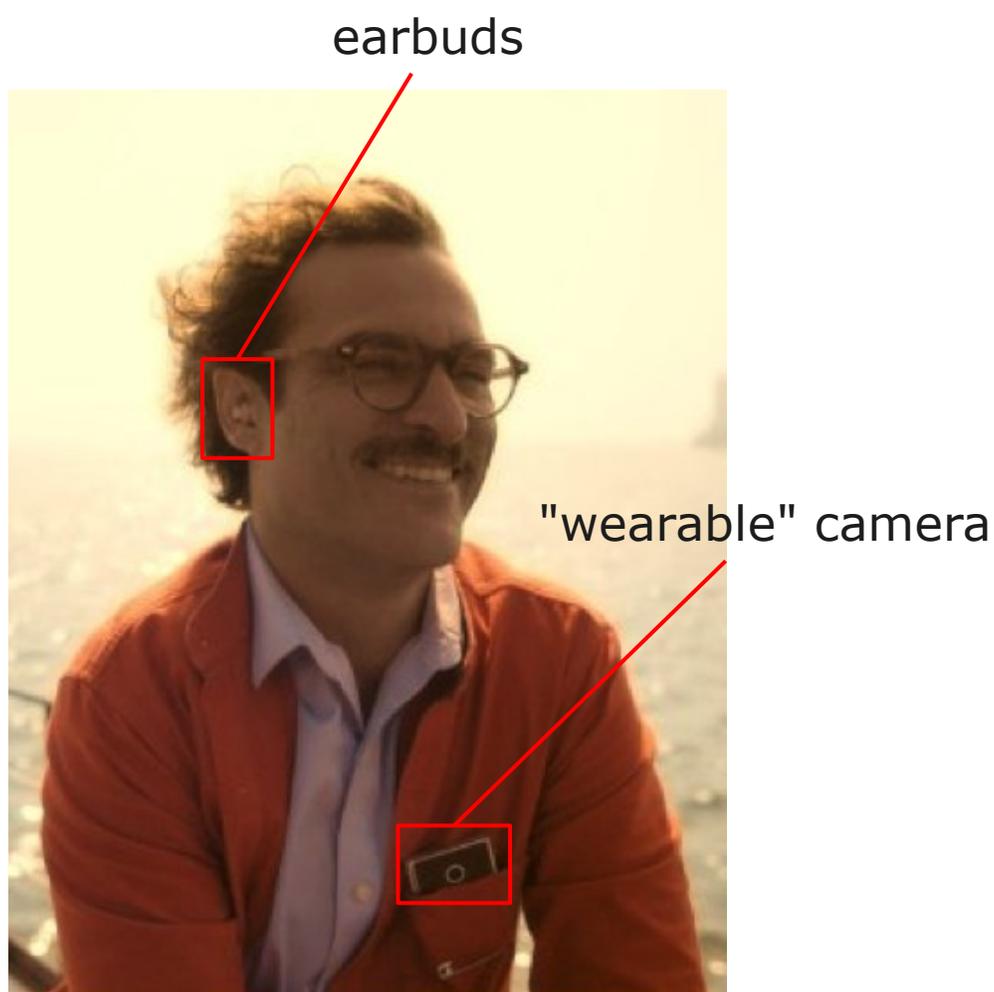


Smartphone Era (2007–present)

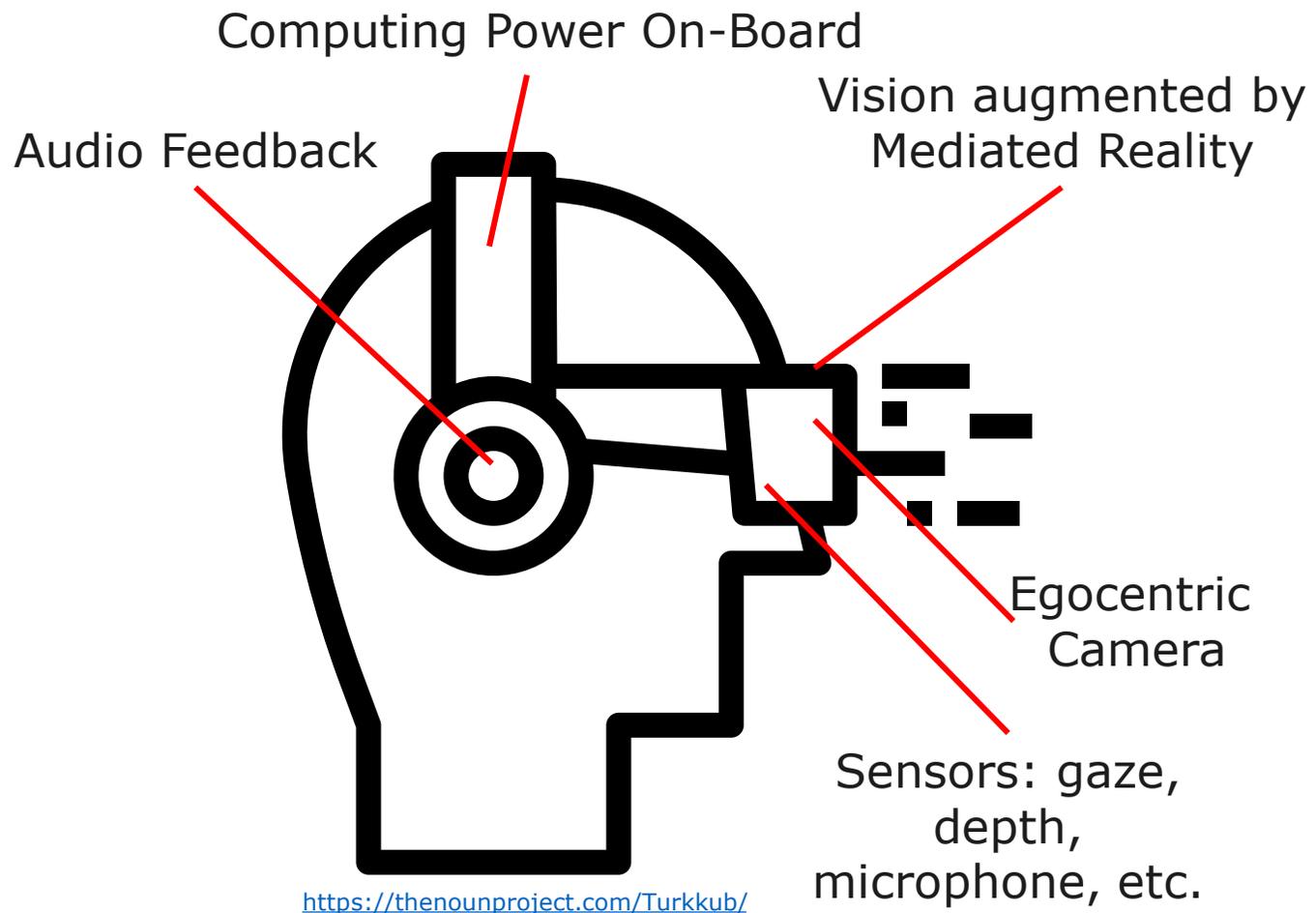
*Computing in your pocket.
Computing is always accessible, but
forces to switch between the digital and
real world*

Smartglasses Era (Now and Future)

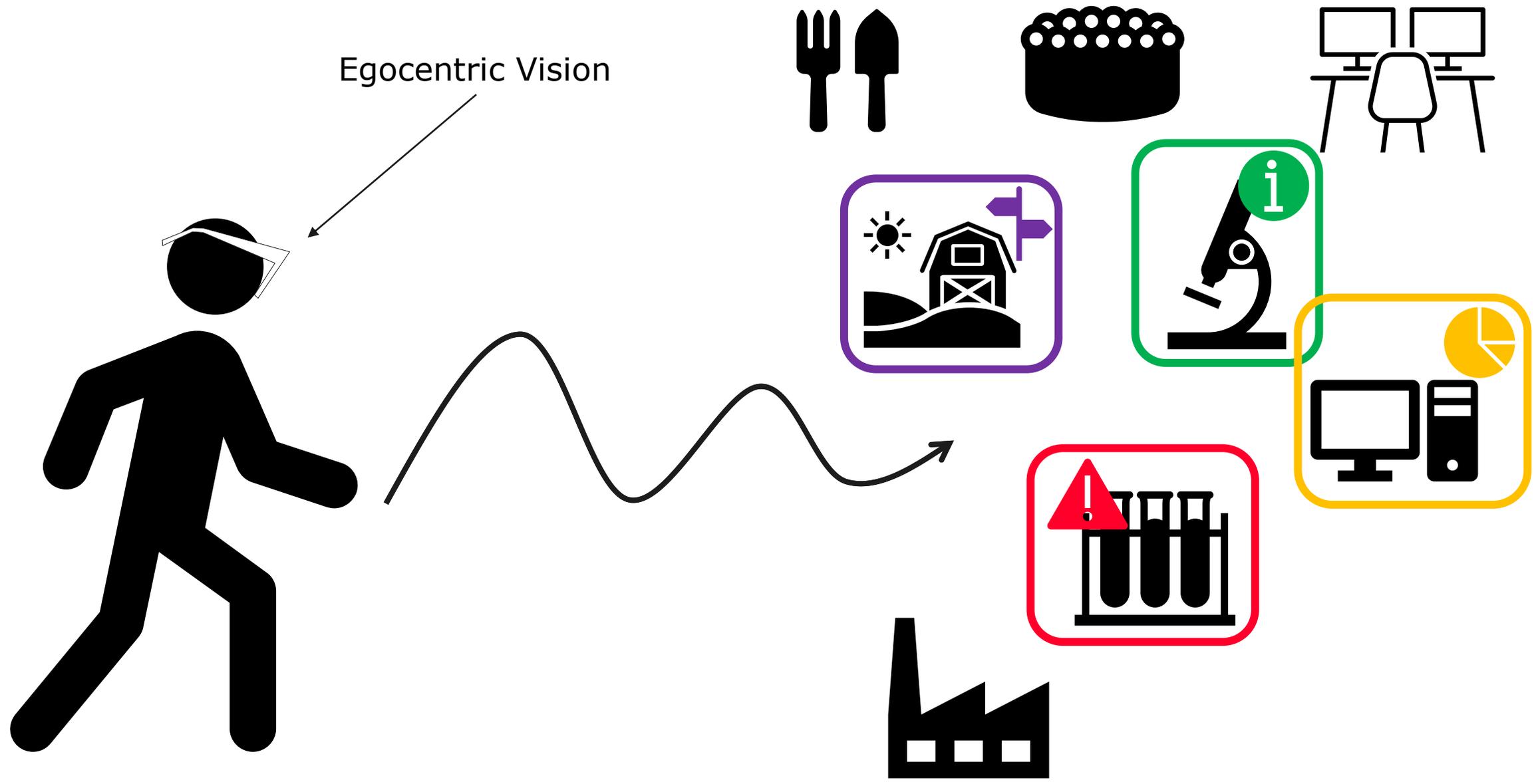
*Hands-free, always-on, egocentric vision.
Computing everywhere with minimal
switch between real and digital worlds*



"her" 2013 movie



A wearable device which perceives the world from our "egocentric" point of view is perfect for implementing a virtual assistant





**(Egocentric) Computer Vision is
Fundamental!**



Exocentric

- ✓ Easy to setup
- ✓ Controlled Field of View
- × Doesn't always see everything
- × Not really portable



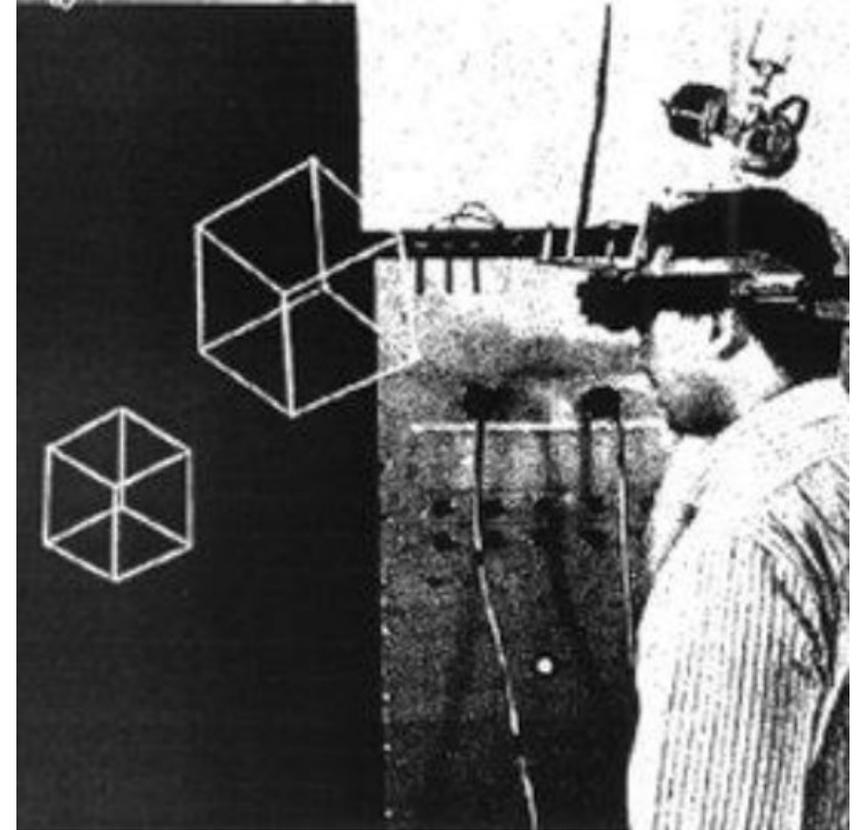
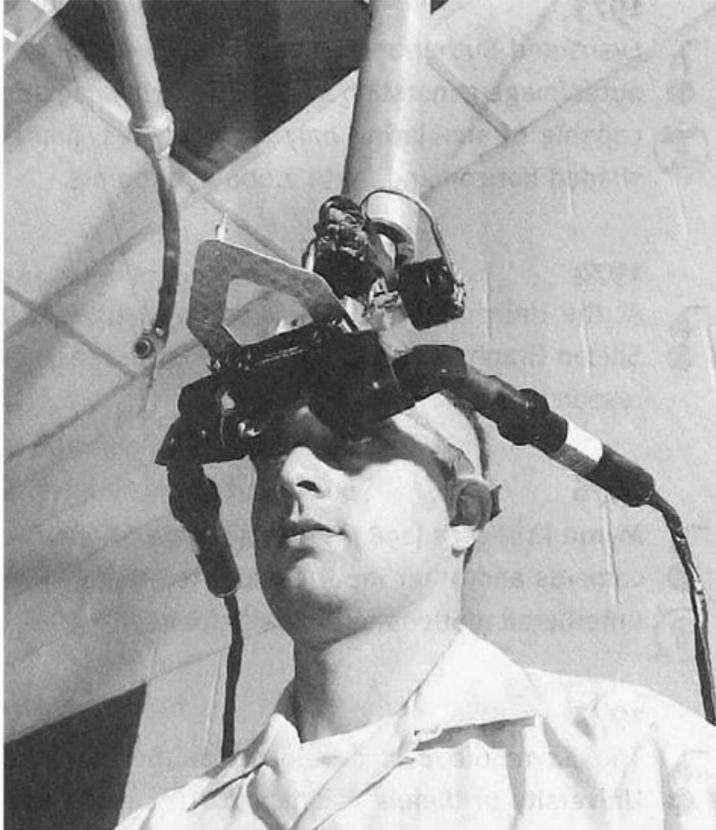
Egocentric

- ✓ Content is always relevant
- ✓ Intrinsically mobile
- × High variability
- × Operational constraints



Receive/Acquire Information

In 1968 Ivan Sutherland invented the first "head mounted display" (HMD), a stereoscopic display mounted on the head of the user which allowed to show wireframe rooms.



Due to its weight, the display was fixed to the ceiling with a pipe, for which it was called «sword of Damocles».

Steve Mann's "wearable computer" and "reality mediator" inventions of the 1970s have evolved into what looks like ordinary eyeglasses.



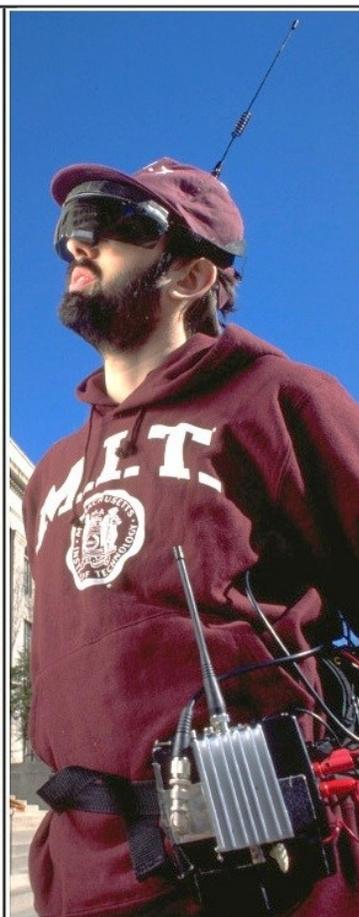
(a)
1980



(b)
Mid 1980s



(c)
Early 1990s



(d)
Mid 1990s



(e)
Late 1990s

In the 80s and 90s Steve Mann (PhD in Media Arts and Sciences at MIT, 1997) invented a number of wearable computers featuring video capabilities, computing capabilities, and a wearable screen for feedback. **Steve Mann is often referred to as «the father of wearable computing»**



- EyeTap Digital Eye Glass
- SWIM (Sequential Wave Imprinting Machine)
- High-dynamic range imaging (HDR)
- Smartwatch
- Visual Orbits



Steve Mann. "Compositing multiple pictures of the same scene." *Proc. IS&T Annual Meeting, 1993.*

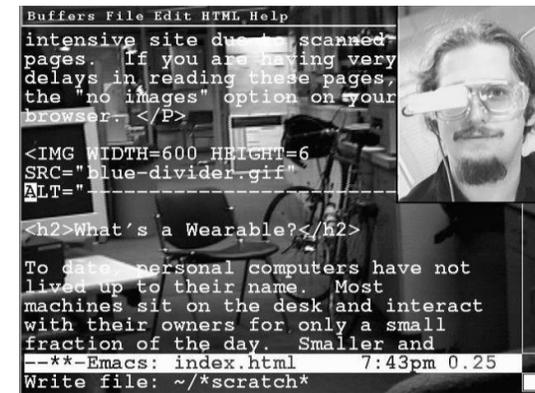
Steve Mann, "Wearable computing: a first step toward personal imaging," in *Computer*, vol. 30, no. 2, pp. 25-32, Feb. 1997.



Augmented Reality Through Wearable Computing

Thad Starner, Steve Mann, Bradley Rhodes, Jeffrey Levine
Jennifer Healey, Dana Kirsch, Roz Picard, and Alex Pentland

The Media Laboratory
Massachusetts Institute of Technology
(augmented reality)



1997

1998



Visual Contextual Awareness in Wearable Computing

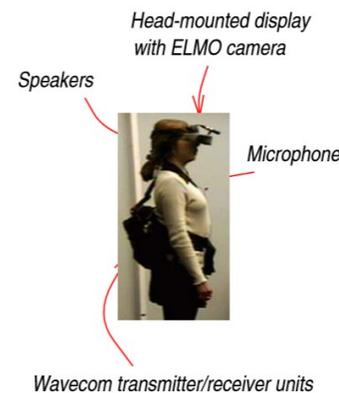
Thad Starner Bernt Schiele Alex Pentland
Media Laboratory, Massachusetts Institute of Technology

(location and task recognition)

An Interactive Computer Vision System DyPERS: Dynamic Personal Enhanced Reality System

Bernt Schiele, Nuria Oliver, Tony Jebara, and Alex Pentland
Vision and Modeling Group
MIT Media Laboratory, Cambridge, MA 02139, USA

(object recognition, media memories)



VISUAL TRIGGER



ASSOCIATED SEQUENCE



GARBAGE NO PLAY-BACK

1999

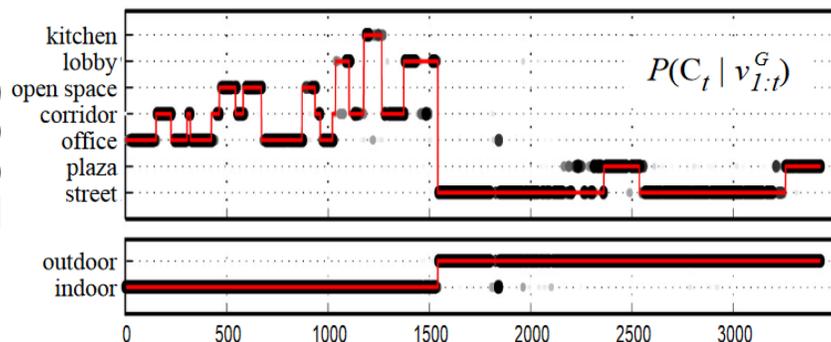
Wearable Visual Robots

W.W. Mayol, B. Tordoff and D.W. Murray
 University of Oxford, Parks Road, Oxford OX1 3PJ, UK
 (active vision)



2000

2003



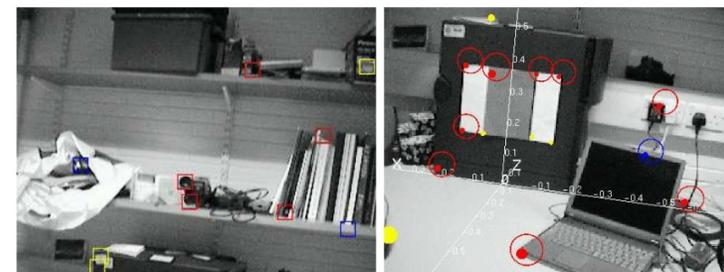
Context-based vision system for place and object recognition

Antonio Torralba MIT AI lab Cambridge, MA 02139	Kevin P. Murphy MIT AI lab Cambridge, MA 02139	William T. Freeman MIT AI lab Cambridge, MA 02139	Mark A. Rubin Lincoln Labs Lexington, MA 02420
---	--	---	--

(location/object recognition)

Real-Time Localisation and Mapping with Wearable Active Vision *

Andrew J. Davison, Walterio W. Mayol and David W. Murray
 Robotics Research Group
 Department of Engineering Science, University of Oxford, Oxford OX1 3PJ, UK
 (active vision, SLAM)



2003

Wearable Hand Activity Recognition for Event Summarization

W.W. Mayol

Department of Computer Science
University of Bristol

D.W. Murray

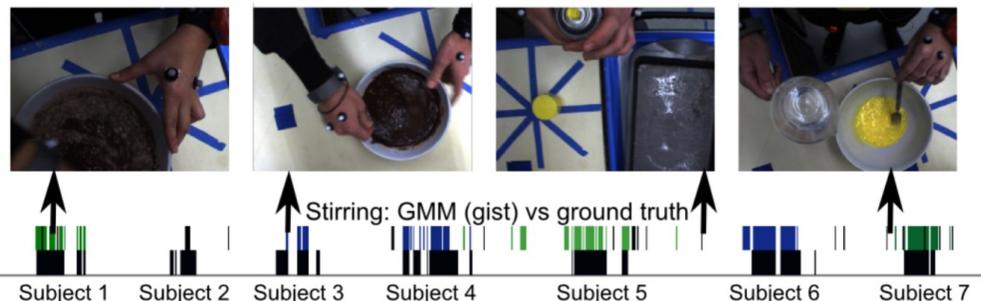
Department of Engineering Science
University of Oxford

(hand activity recognition)



2005

2009



Temporal Segmentation and Activity Classification from First-person Sensing

Ekaterina H. Spriggs, Fernando De La Torre, Martial Hebert
Carnegie Mellon University.

(activity classification)

Figure-Ground Segmentation Improves Handled Object Recognition in Egocentric Video

Xiaofeng Ren

Intel Labs Seattle

1100 NE 45th Street, Seattle, WA 98105

Chunhui Gu

University of California at Berkeley

Berkeley, CA 94720

(handheld object recognition)



2010



"A day in Rome"



- SenseCam is a wearable camera that takes photos automatically;
- Originally conceived as a «personal blackbox» accident recorder;
- Used in the MyLifeBits project, inspired by Bush's Memex;
- Inspired a series of conferences and many research papers.

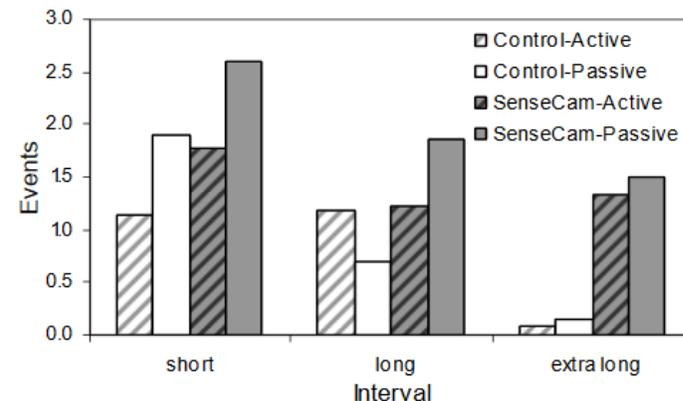
<https://www.microsoft.com/en-us/research/project/sensecam/>

Bell, Gordon, and Jim Gemmell. *Your life, uploaded: The digital way to better memory, health, and productivity*. Penguin, 2010.

Do Life-Logging Technologies Support Memory for the Past? An Experimental Study Using SenseCam

Abigail Sellen, Andrew Fogg, Mike Aitken*, Steve Hodges, Carsten Rother and Ken Wood
 Microsoft Research Cambridge *Behavioural & Clinical Neuroscience Institute
 7 JJ Thomson Ave, Cambridge, UK, CB3 0FB Dept. of Psychology, University of Cambridge

(health, memory augmentation)



2007

2008



(a) Reading in bed



(b) Having dinner

MyPlaces: Detecting Important Settings in a Visual Diary

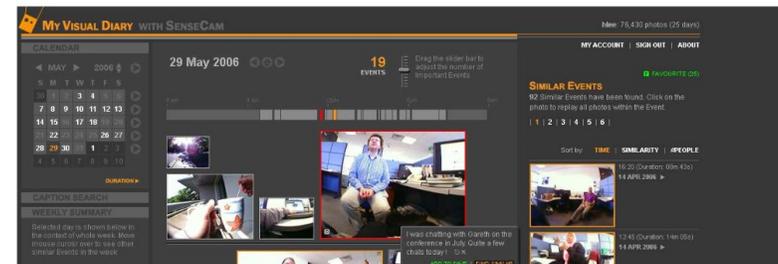
Michael Blighe and Noel E. O'Connor
 Centre for Digital Video Processing, Adaptive Information Cluster
 Dublin City University, Ireland
 {blighem, oconnorn}@eeng.dcu.ie

(lifelogging, place recognition)

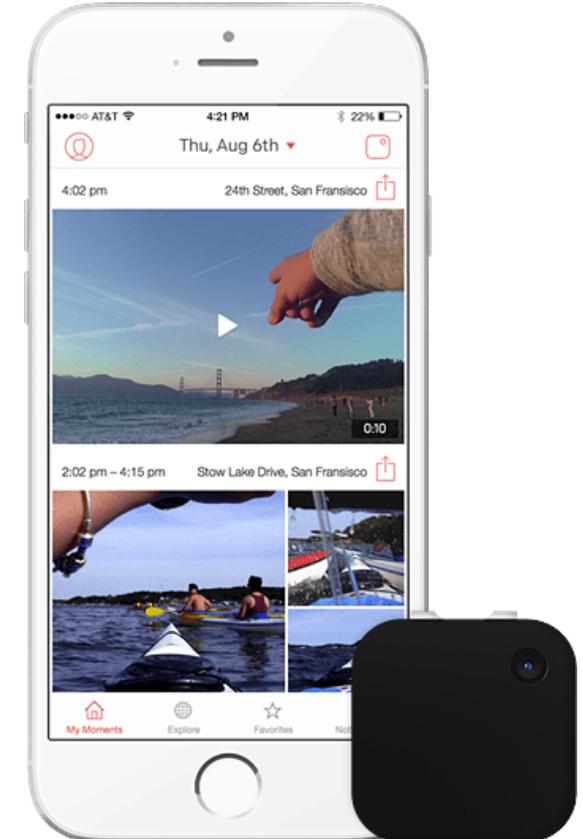
Constructing a SenseCam Visual Diary as a Media Process

Hyowon Lee, Alan F. Smeaton, Noel O'Connor, Gareth Jones, Michael Blighe, Daragh Byrne, Aiden Doherty, and Cathal Gurrin
 Centre for Digital Video Processing & Adaptive Information Cluster,
 Dublin City University

(lifelogging, multimedia retrieval)



2008



<http://getnarrative.com/>

Multi-face tracking by extended bag-of-tracklets in egocentric photo-streams

Maedeh Aghaei^{a,*}, Mariella Dimiccoli^{a,b}, Petia Radeva^{a,b}
(lifelogging, face tracking)



2016

2017

Day's Lifelog:



Event Segmentation

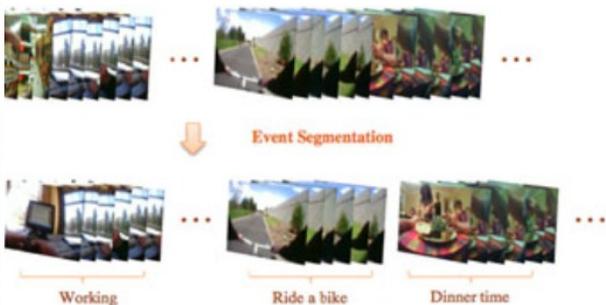
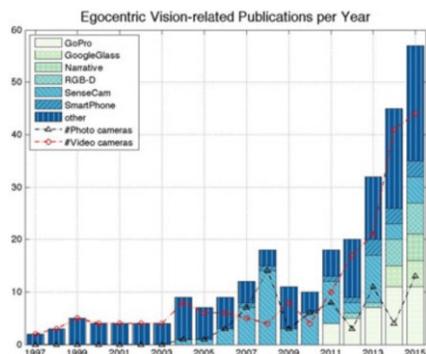
Multiple Events:



SR-clustering: Semantic regularized clustering for egocentric photo streams segmentation

Mariella Dimiccoli^{a,c,1,*}, Marc Bolaños^{a,1,*}, Estefania Talavera^{a,b}, Maedeh Aghaei^a, Stavri G. Nikolov^d, Petia Radeva^{a,c,*}

(lifelogging, event segmentation)



Toward Storytelling From Visual Lifelogging: An Overview

Marc Bolaños, Mariella Dimiccoli, and Petia Radeva

(lifelogging, survey)

2017



different wearing modalities

<https://www.youtube.com/watch?v=D4iU-EOJYK8>



head-mounted



chest-mounted



wrist-mounted

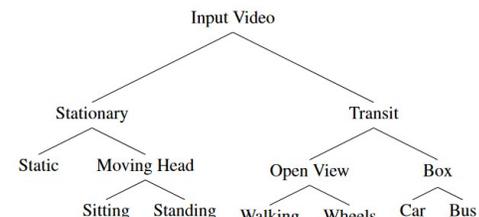


helmet-mounted



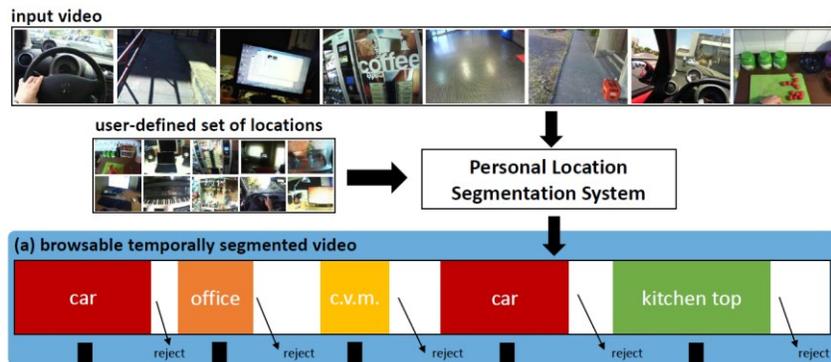
Temporal Segmentation of Egocentric Videos

Yair Poleg Chetan Arora* Shmuel Peleg
 (egocentric video indexing)



2014

2016



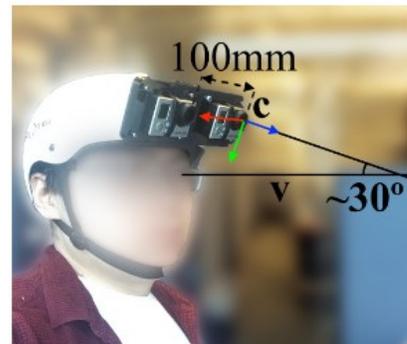
Recognizing Personal Locations from Egocentric Videos

Antonino Furnari, Giovanni Maria Farinella, *Senior Member, IEEE*, and Sebastiano Battiato, *Senior Member, IEEE*

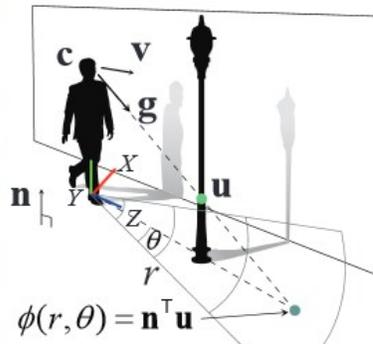
(localization, indexing, context-aware computing)

Egocentric Future Localization

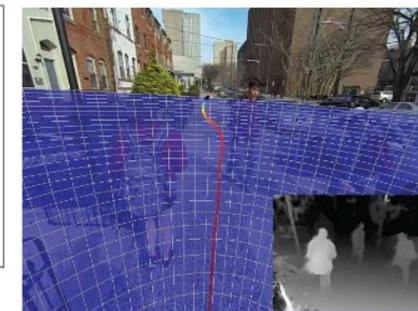
Hyun Soo Park Jyh-Jing Hwang Yedong Niu Jianbo Shi
 (future localization, navigation)



(a) Ego-stereo cameras



(b) Geometry

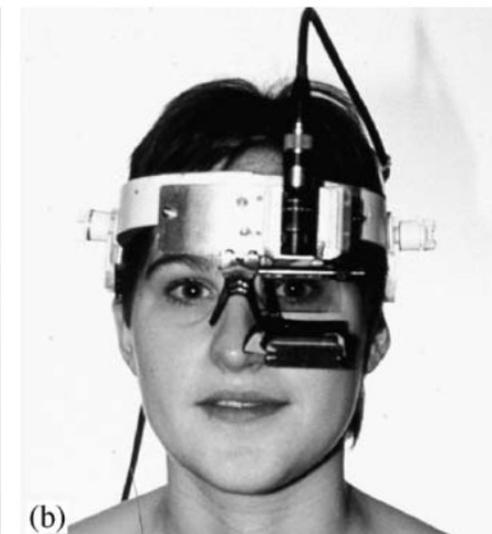
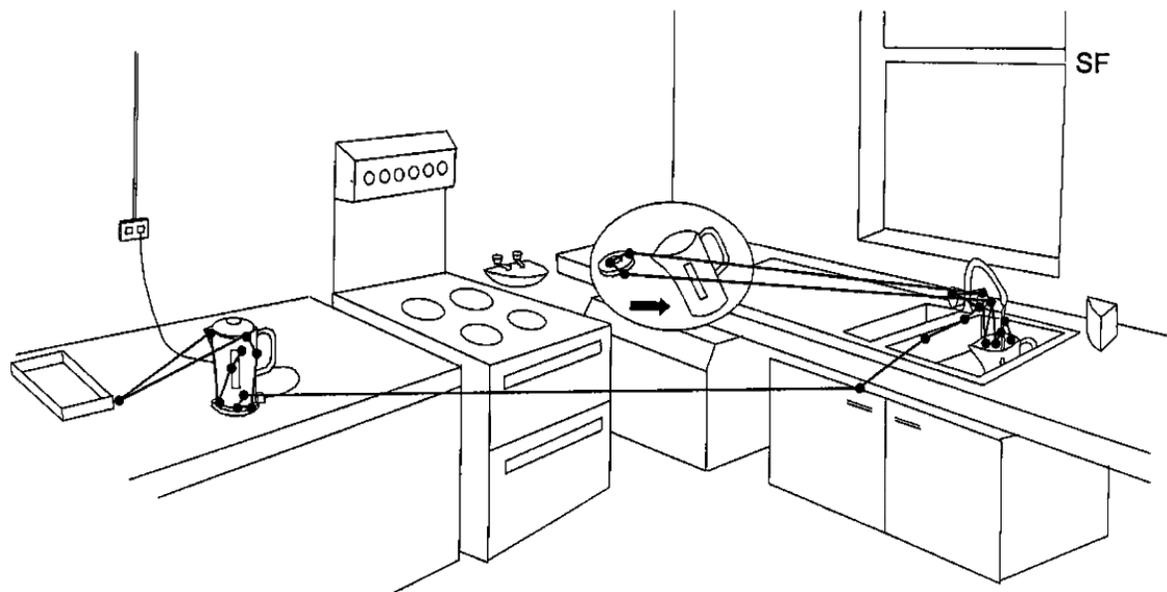


(c) Egocentric RGBD image

2016

Eye movements and the control of actions in everyday life

Michael F. Land



Prototype by Land (1993)

Gaze is important in Egocentric Vision!



Tobii Pro Glasses 2 (2014)



Microsoft HoloLens 2 (2016)



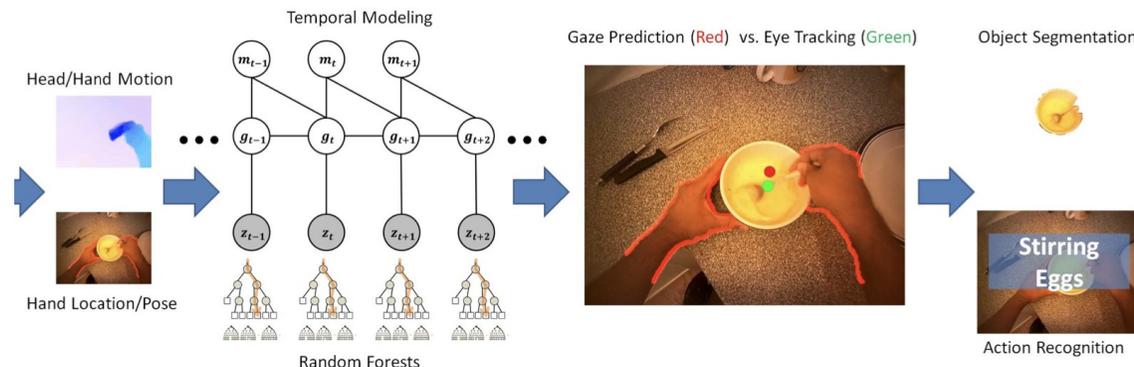
Mobile Eye-XG (2013)



Pupil Eye Trackers (2014 -)

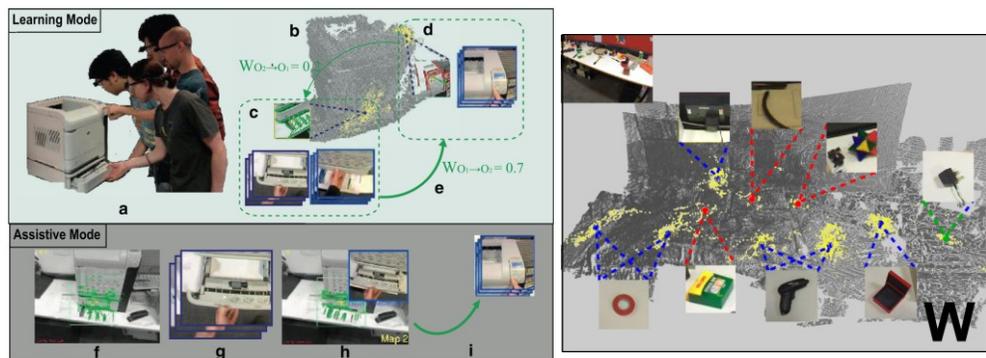
Learning to Predict Gaze in Egocentric Video

Yin Li, Alireza Fathi, James M. Rehg
(gaze prediction, action recognition)



2012

2016



You-Do, I-Learn: Egocentric unsupervised discovery of objects and their modes of interaction towards video-based guidance

Dima Damen*, Teesid Leelasawassuk, Walterio Mayol-Cuevas

(object usage discovery, assistance)

MECCANO: A multimodal egocentric dataset for humans behavior understanding in the industrial-like domain

Francesco Ragusa*, Antonino Furnari, Giovanni Maria Farinella

(gaze prediction, procedural video)



2023



Health, assistive technologies

<https://www.orcam.com/>



<https://www.orcham.com/>

Mixed Reality

<https://www.microsoft.com/hololens>



<https://youtu.be/eqFqtAJMtYE>



HoloLens 2

An ergonomic, untethered self-contained holographic device with enterprise-ready applications to increase user accuracy and output.

\$3,500



HoloLens 2 Industrial Edition

A HoloLens 2 that is designed and tested to support regulated environments such as clean rooms and hazardous locations.

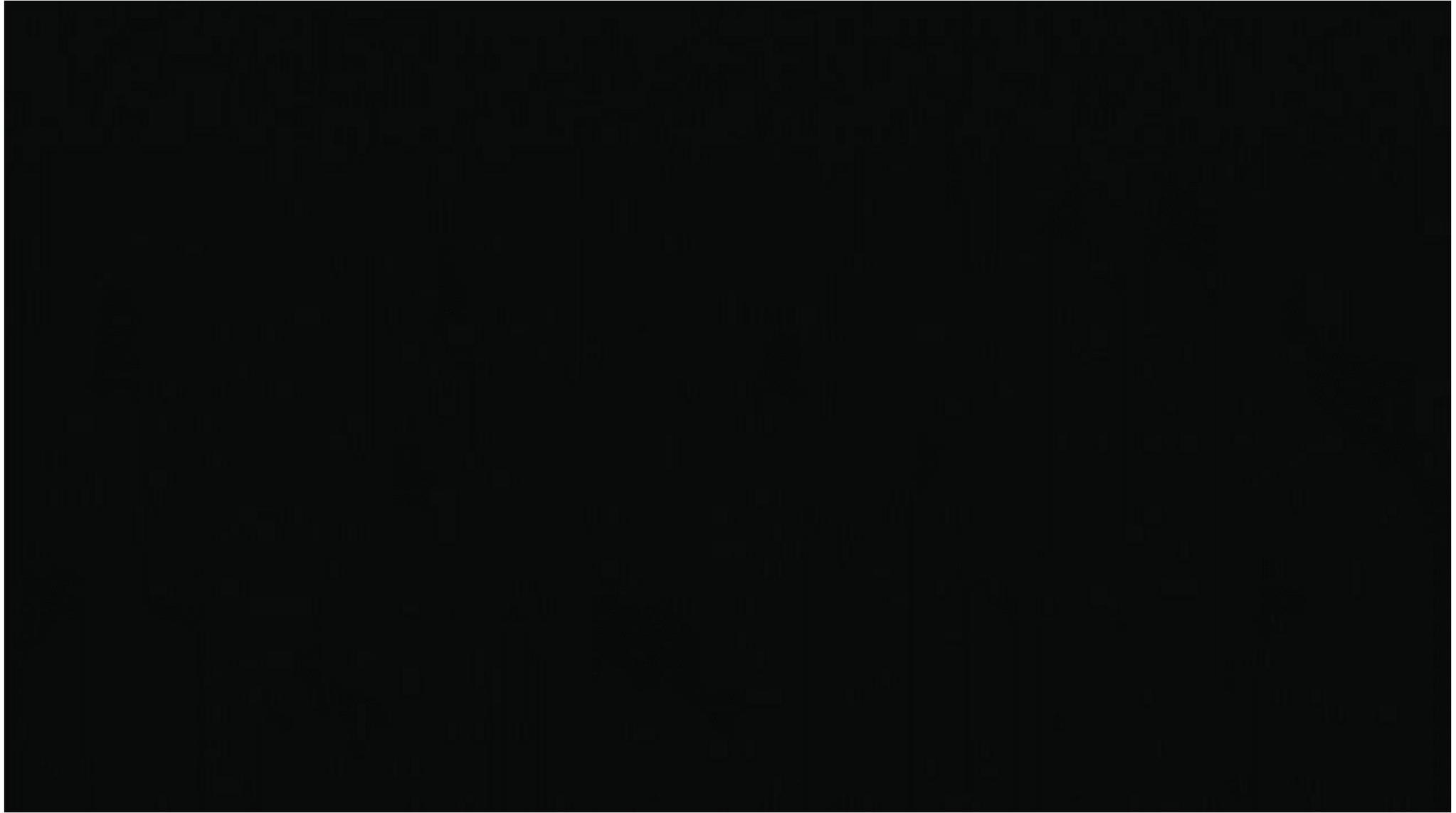
\$4,950



Trimble XR10 with HoloLens 2

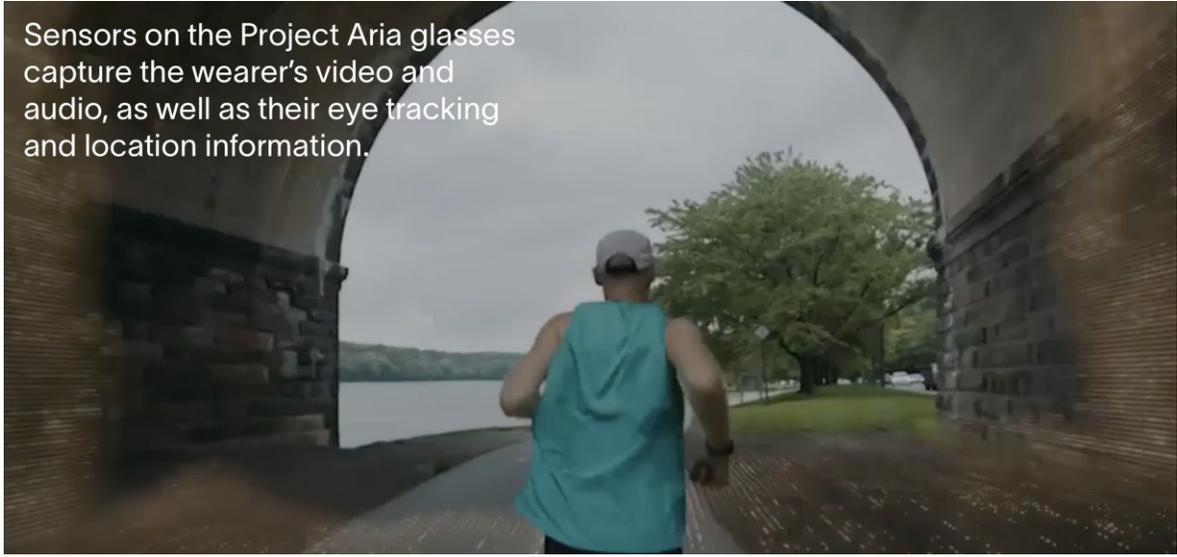
A hardhat-integrated HoloLens 2 that is purpose-built for personnel in dirty, loud, and safety-controlled work site environments.

\$5,199



<https://www.magicleap.com/magic-leap-2>

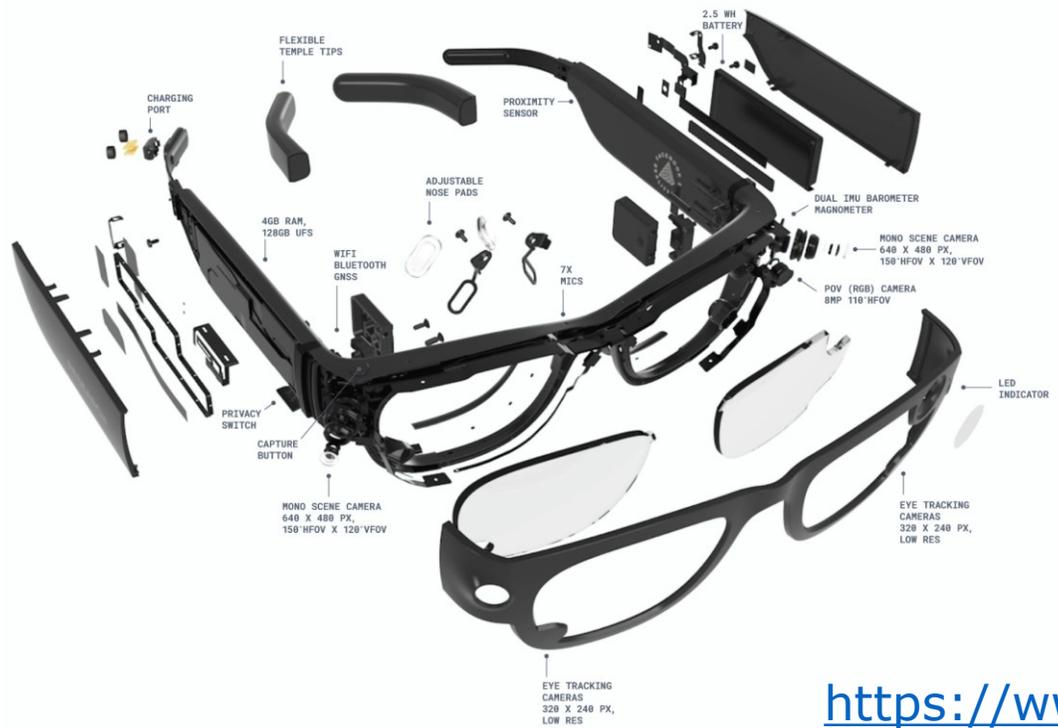
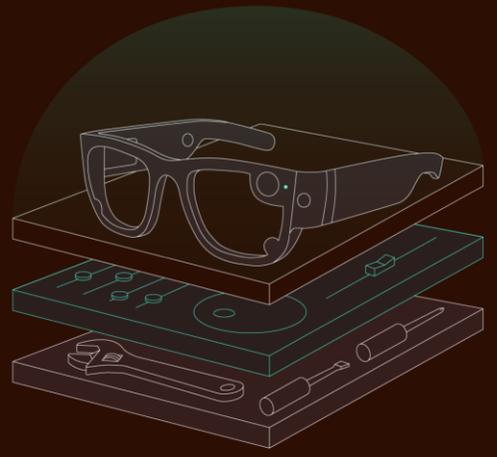
Sensors on the Project Aria glasses capture the wearer's video and audio, as well as their eye tracking and location information.



Aria Research Kit

For approved research partners, Meta offers a kit that includes Project Aria glasses and SDK, so that researchers can conduct independent studies and help shape the future of AR.

[LEARN MORE ABOUT PARTNERING WITH PROJECT ARIA](#)



<https://www.projectaria.com>

52° FOV



Development Kit



6 DoF Positional Tracking

Glasses track real-time position relative to the world, detect planes and images, and obtain environmental depth information.

Image Tracking

Recognizing physical images for AR experiences using multiple reference images in a single session.

Plane Detection

Detection flat surfaces (horizontal/vertical) like tables and walls.

Hand Tracking

Interact with AR content using natural hand gestures, enabling seamless manipulation of virtual objects without additional controllers.

Depth Mesh

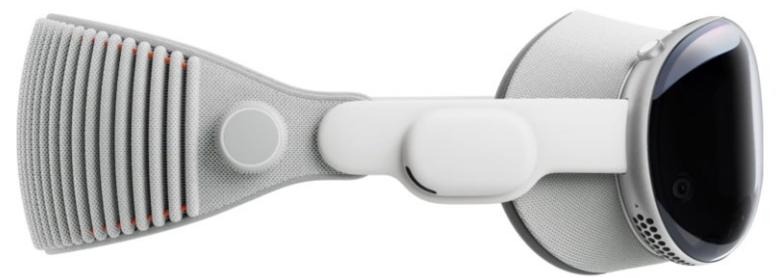
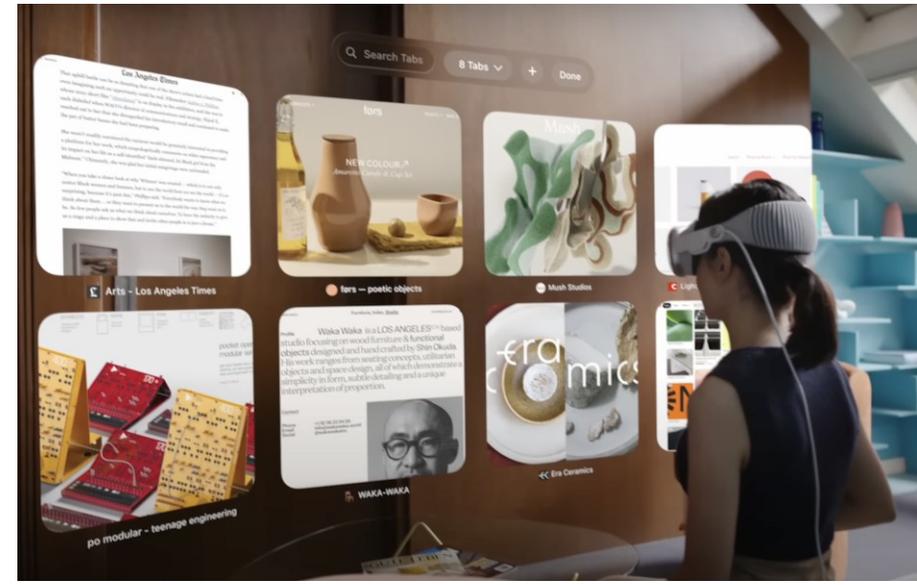
Allowing 3D surface and object detection for realistic AR integration with the real world.

Optimized Rendering

Automatically applied to reduce latency, jitter, and enhance user experience.

Spatial Anchor

Precisely anchor virtual objects to real-world locations, maintaining accurate positioning for collaborative AR experiences and persistent content.









Too Many Devices?

towards standardization...

Unified API supported by many AR and VR devices



XR APPLICATIONS

Head & Hand Pose Information
Controller Input State
Display Configuration



Image(s) to Display
Audio
Haptic Responses

XR PLATFORMS & DEVICES

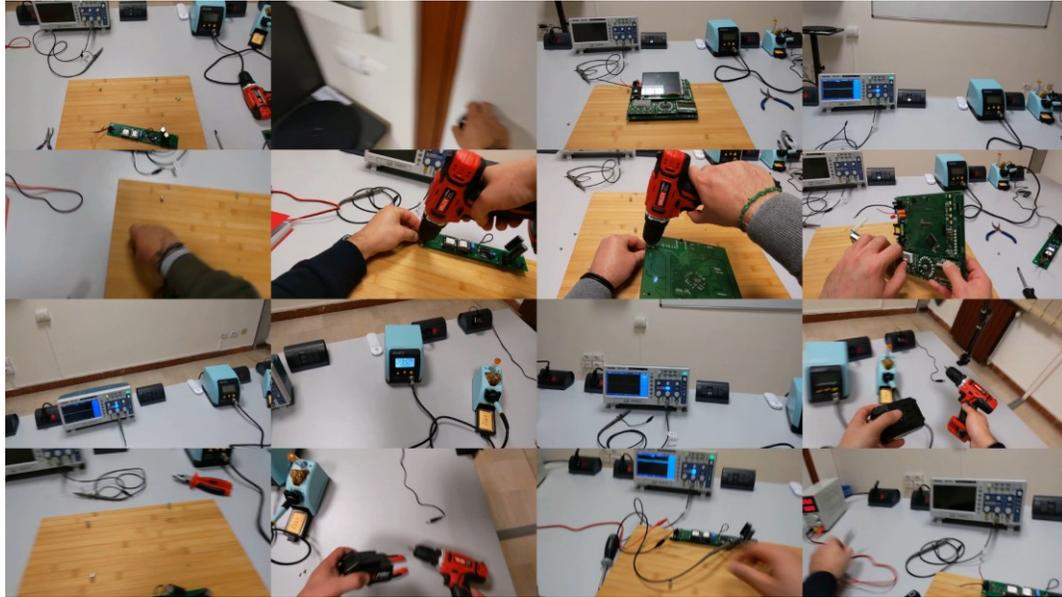




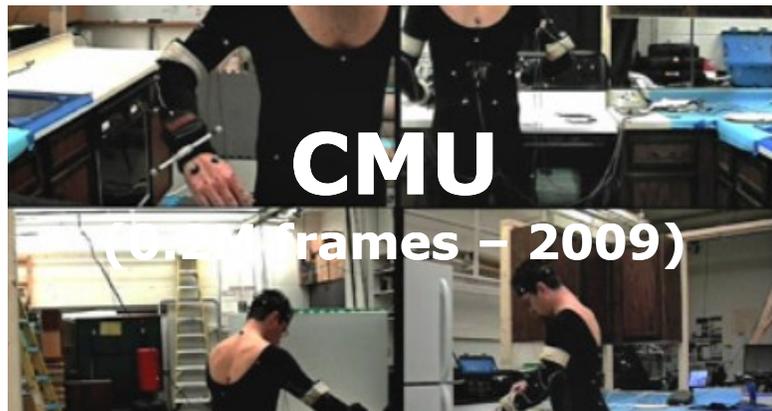
Workshop on Egocentric (First Person) Vision

ACVR





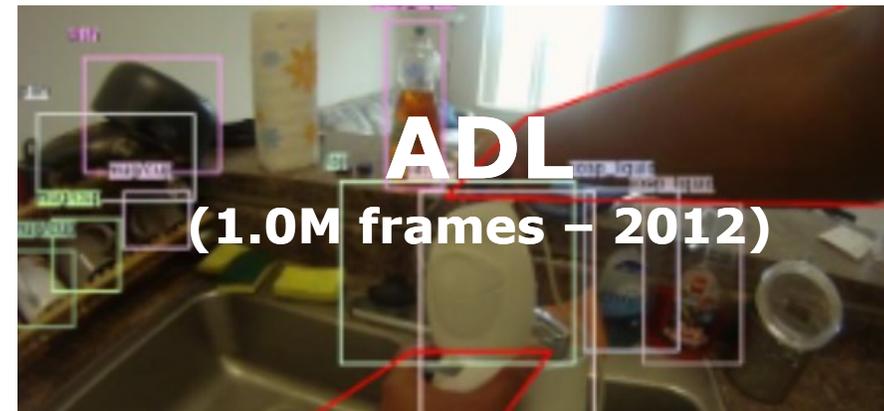
Digital Information



<http://www.cs.cmu.edu/~espriggs/cmu-mmacc/annotations/>



<http://www.cbi.gatech.edu/fpv/>



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



<https://allenai.org/plato/charades/>



<http://www.cbi.gatech.edu/fpv/>

Scaling Egocentric Vision: The EPIC-KITCHENS Dataset

Dima Damen¹, Hazel Doughty¹, Giovanni Maria Farinella¹, Antonino Furnari², Evangelos Kazakos¹, Davide Moltisanti¹, Jonathan Munro³, Toby Perrett¹, Will Price¹, and Michael Wray¹

¹Uni. of Bristol, UK ²Uni. of Catania, Italy, ³Uni. of Toronto, Canada

Abstract. First-person vision is gaining interest as it offers a unique viewpoint on people's interaction with objects, their attention, and even intention. However, progress in this challenging domain has been relatively slow due to the lack of sufficiently large datasets. In this paper, we introduce EPIC-KITCHENS, a large-scale egocentric video benchmark recorded by 32 participants in their native kitchen environments. Our videos depict non-scripted daily activities: we simply asked each participant to start recording every time they entered their kitchen. Recording took place in 4 cities (in North America and Europe) by participants belonging to 10 different nationalities, resulting in highly diverse cooking styles. Our dataset features 55 hours of video consisting of 11.5M frames, which we densely labelled for a total of 39.6K action segments and 454.3K object bounding boxes. Our annotation is unique in that we had the participants narrate their own videos (after recording), thus reflecting true intention, and we crowd-sourced ground-truths based on these. We describe our object, action and anticipation challenges, and evaluate several baselines over two test splits, seen and unseen kitchens.

Keywords: Egocentric Vision, Dataset, Benchmarks, First-Person Vision, Egocentric Object Detection, Action Recognition and Anticipation

EPIC-Kitchens 55

EPIC-SOUNDS: A Large-Scale Dataset of Actions that Sound

Jaesung Huh¹, Jacob Chalk², Evangelos Kazakos³, Dima Damen², Andrew Zisserman¹

¹Visual Geometry Group, Department of Engineering Science, University of Oxford, UK
²Department of Computer Science, University of Bristol, UK
³CIIRC, Czech Technical University in Prague, Czech Republic

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

Fig. 1: Sample video with corresponding audio from EPIC-KITCHENS-100. We compare the already published visual labels with our collected EPIC-SOUNDS audio labels. We demonstrate the differences between the modality annotations, both in temporal extent and class labels, highlighting: **Misaligned intervals:** temporal boundaries are distinct; **Invisible action:** action not seen in the video, but which produces distinct sounds (0-to-1 matching); **Indistinguishable sounds:** sounds from two distinct visual actions, but are audibly inseparable; **Silent action:** visual action that does not have audible sounds (1-to-0); and visual actions containing multiple **repetitive sounds** (1-to-N).

Abstract—We introduce EPIC-SOUNDS, a large-scale dataset of audio annotations capturing temporal extents and class labels visible in the video stream of 2.5 hours of unscripted kitchen recordings.

Index Terms—audio recognition, action recognition, audio event detection, audio dataset, data collection, dataset

EPIC-SOUNDS

Rescaling Egocentric Vision: Collection Pipeline and Challenges for EPIC-KITCHENS-100

Dima Damen¹, Hazel Doughty¹, Giovanni Maria Farinella¹, Antonino Furnari², Evangelos Kazakos¹, Jian Ma¹, Davide Moltisanti¹, Jonathan Munro¹, Toby Perrett¹, Will Price¹, Michael Wray¹

Received: 18 Jan 2021, Revised: 23 Aug 2021, Accepted: 17 Sep 2021

Abstract This paper introduces the pipeline to extend the largest dataset in egocentric vision, EPIC-KITCHENS. The effort culminates in EPIC-KITCHENS-100, a collection of 100 hours, 20M frames, 60K actions in 700 variable-length videos, capturing long-term unscripted activities in 45 environments, using head-mounted cameras. Compared to its previous version [1], EPIC-KITCHENS-100 has been annotated using a novel pipeline that allows denser (54% more actions per minute) and more complete annotations of fine-grained actions (+128% more action segments). This collection enables new challenges such as action detection and evaluating the “test of time” — i.e. whether models trained on data collected in 2018 can generalise to new footage collected two years later.

1 Introduction and Related Datasets

Since the dawn of machine learning for computer vision, datasets have been curated to train models, for single tasks from classification [2,3] to detection [4,5], captioning [6,7] and segmentation [8,9]. Increasingly, datasets have been used for novel tasks, through pre-training [10], self-supervision [12,13] or additional annotations [14,15]. However, task adaptation demonstrates that models

recognition. For each challenge, we define the task, provide baselines and evaluation metrics [5].

Keywords Video Dataset, Egocentric Vision, First-Person Vision, Action Understanding, Multi-Benchmark Large-Scale Dataset, Annotation Quality

EPIC-Kitchens 100

EPIC Fields Marrying 3D Geometry and Video Understanding

Vadim Tschernecki¹, Ahmad Darkhalil², Zhifan Zhu², David Fouhey¹, Iro Laina¹, Diane Larlus¹, Dima Damen¹, Andrea Vedaldi¹

¹VGG, University of Oxford ²University of Bristol
³New York University ⁴NAVER LABS Europe ^{*}: Equal Contribution

Abstract

Neural rendering is fuelling a unification of learning, 3D geometry and video understanding that has been waiting for more than two decades. Progress, however, is still hampered by a lack of suitable datasets and benchmarks. To address this gap, we introduce EPIC Fields, an augmentation of EPIC-KITCHENS with 3D camera information. Like other datasets for neural rendering, EPIC Fields removes the complex and expensive step of reconstructing cameras using photogrammetry, and allows researchers to focus on modelling problems. We illustrate the challenge of photogrammetry in egocentric videos of dynamic actions and propose innovations to address them. Compared to other neural rendering datasets, EPIC Fields is better tailored to video understanding because it is paired with labelled action segments and the recent VISOR segment annotations. To further motivate the community, we also evaluate three benchmark tasks in neural rendering and segmenting dynamic objects, with strong baselines that showcase what is not possible today. We also highlight the advantage of geometry in semi-supervised video object segmentations on the VISOR annotations. EPIC Fields reconstructs 96% of videos in EPIC-KITCHENS, registering 19M frames in 99 hours recorded in 45 kitchens, and is available from: <http://epic-kitchens.github.io/epic-fields>

EPIC-FIELDS

EPIC-KITCHENS VISOR Benchmark Video Segmentations and Object Relations

Ahmad Darkhalil¹, Dandan Shan², Bin Zhu³, Jian Ma³, Amlan Kar¹, Richard Higgins¹, Sanja Fidler¹, David Fouhey¹, Dima Damen¹

¹Uni. of Bristol, UK ²Uni. of Michigan, US ³Uni. of Toronto, CA ^{*}: Co-First Authors

Abstract

We introduce VISOR, a new dataset of pixel annotations and a benchmark suite for segmenting hands and active objects in egocentric video. VISOR annotates videos from EPIC-KITCHENS, which comes with a new set of challenges not encountered in current video segmentation datasets. Specifically, we need to ensure both short- and long-term consistency of pixel-level annotations as objects undergo transformative interactions, e.g. an onion is peeled, diced and cooked - where we aim to obtain accurate pixel-level annotations of the peel, onion pieces, chopping board, knife, pan, as well as the acting hands. VISOR introduces an annotation pipeline, AI-powered in parts, for scalability and quality. In total, we publicly release 272K manual semantic masks of 257 object classes, 9.9M interpolated dense masks, 67K hand-object relations, covering 36 hours of 179 untrimmed videos. Along with the annotations, we introduce three challenges in video object segmentation, interaction understanding and long-term reasoning.

For data, code and leaderboards: <http://epic-kitchens.github.io/VISOR>

EPIC-Kitchens VISOR

HD-EPIC: A Highly-Detailed Egocentric Video Dataset

Toby Perrett¹, Ahmad Darkhalil², Saptarshi Sinha³, Omar Emarat⁴, Sam Pollard⁵, Kranti Parida⁶, Kaiting Liu⁷, Prajwal Gatti⁸, Siddhant Bansal⁹, Kevin Flanagan¹⁰, Jacob Chalk¹¹, Zhifan Zhu¹², Rhodri Guerrier¹³, Fahd Abdelazim¹⁴, Bin Zhu¹⁵, Davide Moltisanti¹⁶, Michael Wray¹⁷, Hazel Doughty¹⁸, Dima Damen¹⁹

¹Uni. of Bristol ²Leiden Uni. ³Singapore Management Uni. ⁴Uni. of Bath ^{*}: Equal Contribution

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

Abstract

We show the potential of our highly-detailed annotations through a challenging VQA benchmark of 26K questions assessing the capability to recognise recipes, ingredients, nutrition, fine-grained actions, 3D perception, object motion and gaze direction. The powerful long-context Gemini Pro only achieves 38.5% on this benchmark, showcasing its difficulty and highlighting shortcomings in current VLMs. We additionally assess action recognition, sound recognition and long-term video-object segmentation on HD-EPIC.

HD-EPIC is 41 hours of video in 9 kitchens with digital twins of 413 kitchen fixtures, capturing 69 recipes, 59K fine-grained actions, 51K audio events, 20K object movements and 37K object masks lifted to 3D. On average, we have 263 annotations per minute of our unscripted videos.

HD-EPIC

EPIC-KITCHENS TEAM



Dima Damen
Principal Investigator
University of Bristol
United Kingdom



Sanja Fidler
Co-Investigator
University of Toronto
Canada



Giovanni Maria Farinella
Co-Investigator
University of Catania
Italy



Davide Moltisanti
(Apr 2017 -)
University of Bristol



Michael Wray
(Apr 2017 -)
University of Bristol



Hazel Doughty
(Apr 2017 -)
University of Bristol



Toby Perrett
(Apr 2017 -)
University of Bristol



Antonino Furnari
(Jul 2017 -)
University of Catania



Jonathan Munro
(Sep 2017 -)
University of Bristol



Evangelos Kazakos
(Sep 2017 -)
University of Bristol



Will Price
(Oct 2017 -)
University of Bristol

Dima Damen, Hazel Doughty, Giovanni Maria Farinella, Sanja Fidler, Antonino Furnari, Evangelos Kazakos, Davide Moltisanti, Jonathan Munro and Toby Perrett, Will Price, Michael Wray (2021). The EPIC-KITCHENS Dataset: Collection, Challenges and Baselines. PAMI, 43(11), pp. 4125-4141.



32 KITCHENS

EPIC-KITCHENS-100



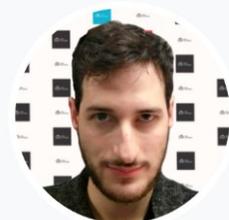
Dima Damen
University of Bristol



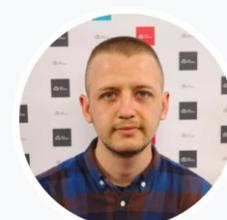
Hazel Doughty
University of Bristol



Giovanni M. Farinella
University of Catania



Antonino Furnari
University of Catania



Evangelos Kazakos
University of Bristol



Jian Ma
University of Bristol



Davide Moltisanti
University of Bristol



Jonathan Munro
University of Bristol



Toby Perrett
University of Bristol

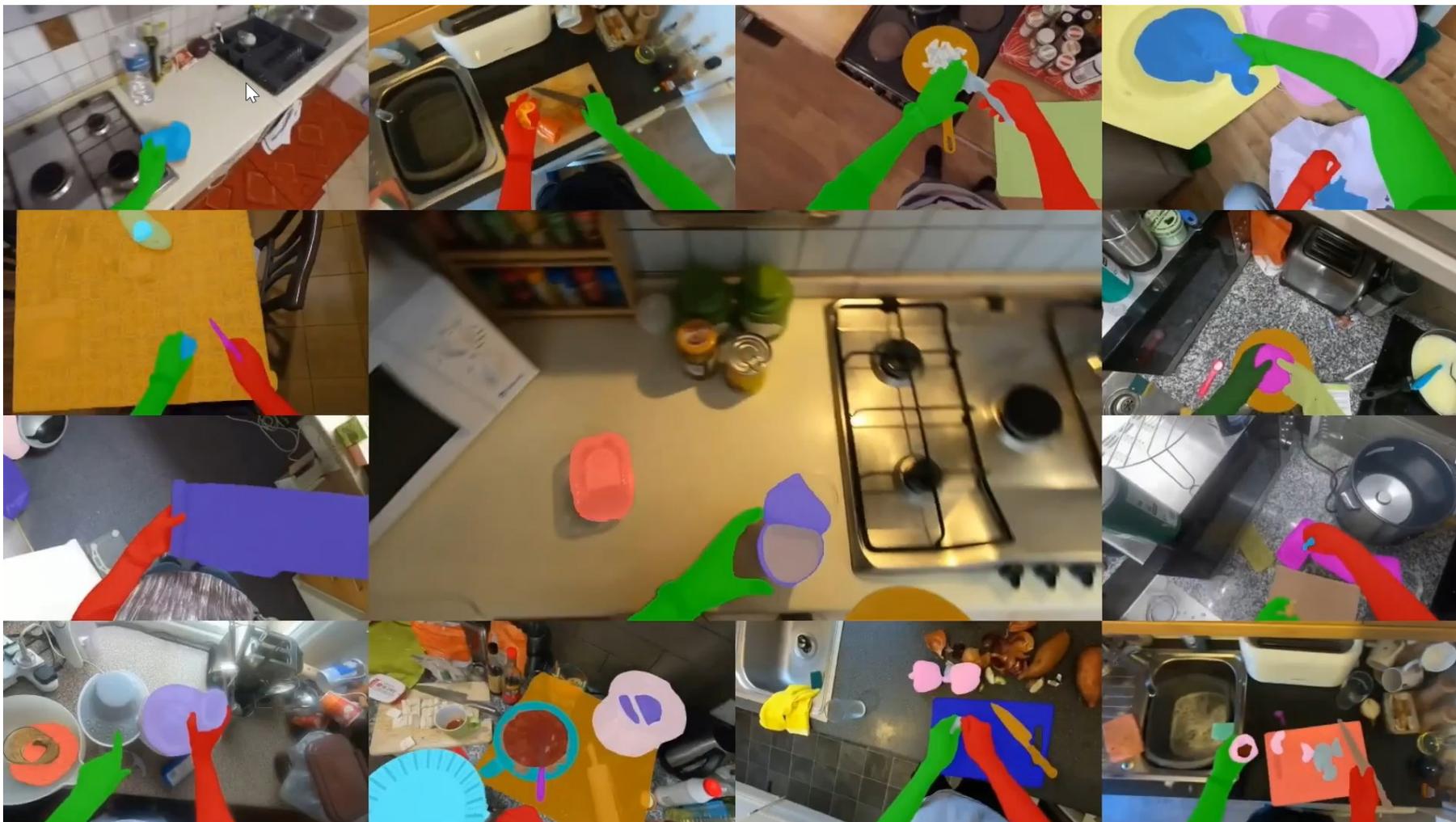


Will Price
University of Bristol



Michael Wray
University of Bristol

	EPIC-KITCHENS-55	EPIC-KITCHENS-100
No. of Hours	55	100
No. of Kitchens	32	45
No. of Videos	432	700
No. of Action Segments	39,432	89,979
Action Classes	2,747	4,025
Verb Classes	125	97
Noun Classes	331	300
Splits	Train/Test	Train/Val/Test
No. of Challenges	3	6 (4 new challenges)

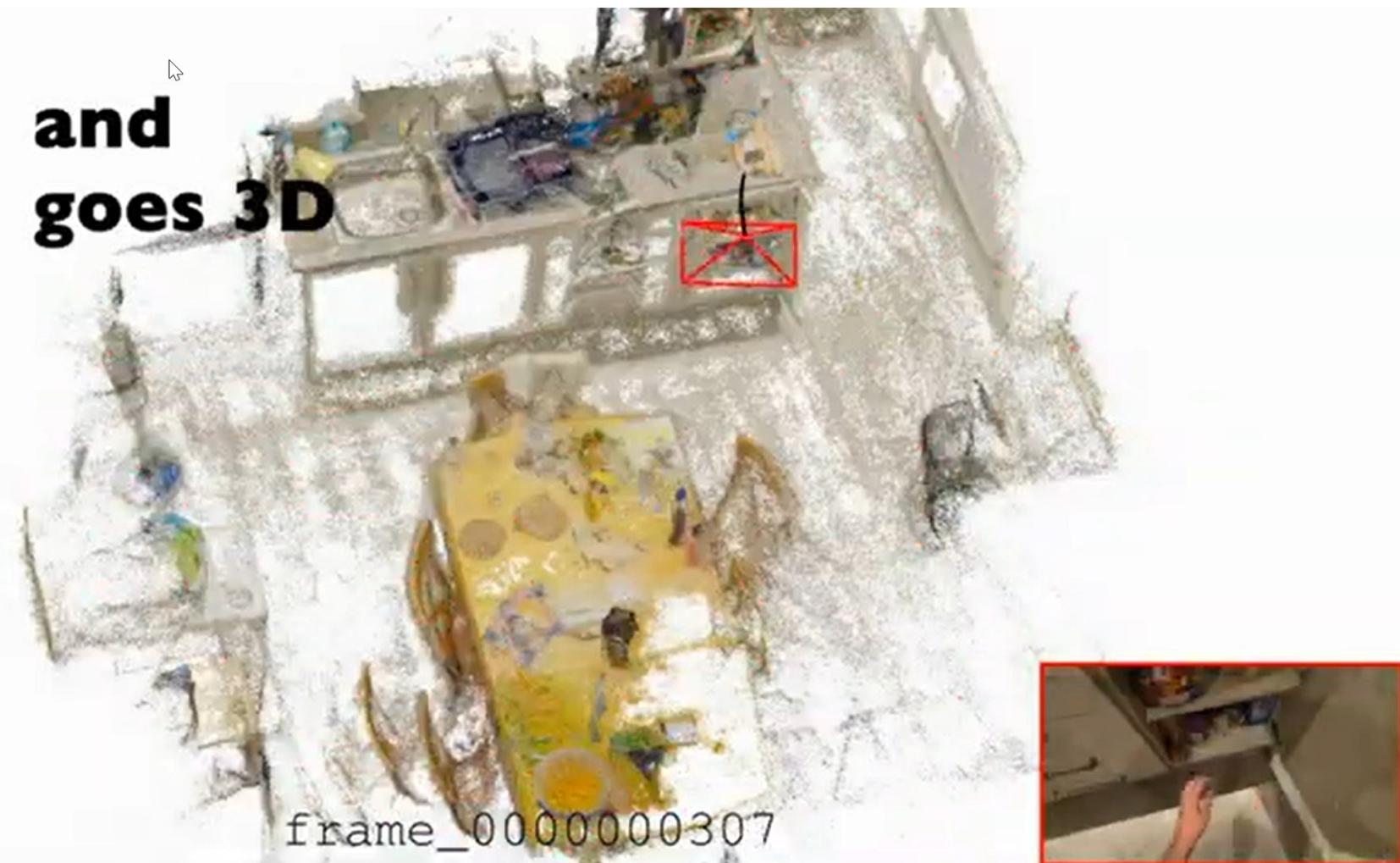


- 272K manual sparse masks for hands and active objects;
- Hand-object contact relations;
- 1477 unique entities;
- 22 categories.



- 74.8K categorised audio segments;
- Material-based collision sounds;
- Repetitive sounds;
- 44 classes.

and
goes 3D



- 19M registered frames;
- Camera poses;
- 3D reconstruction;
- Paired with VISOR annotations.

Preps and Steps

- Recipe and Nutrition;
- Preparation and Step;
- Narrations;
- Audio Annotations;
- Digital Twin;
- Gaze Priming;

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

EPIC-KITCHENS-100- 2022 Challenges Report

RESULTS - 2024 CHALLENGES (JUNE 2024)

EPIC-Kitchens Challenges @CVPR2024, Seattle, US

2024 CHALLENGE WINNERS

	Team	Member	Affiliations
Action Recognition	① KAUST-4Paradigm -MoonshotAI-Nvidia	Shuming Liu	King Abdullah University of Science and Technology
		Lin Sui	4Paradigm Inc
		Chen-Lin Zhang	Moonshot AI
		Fangzhou Mu	NVIDIA
		Chen Zhao	King Abdullah University of Science and Technology
	② DeepGlint (dg_team)	Bernard Ghanem	King Abdullah University of Science and Technology
		Yingxin Xia	DeepGlint and Harbin Institute of Technology
		Ninghua Yang	DeepGlint
		Kaicheng Yang	DeepGlint
		Xiang An	DeepGlint
	③ Shanghai AI Laboratory (Aiyiyai)	Xiangzi Dai	DeepGlint
		Weimo Deng	DeepGlint
		Ziyong Feng	DeepGlint
		Baoqi Pei	Shanghai AI Laboratory and Zhejiang University
		Yifei Huang	Shanghai AI Laboratory
	Guo Chen	Shanghai AI Laboratory and Nanjing University	
	Jilan Xu	Shanghai AI Laboratory and Fudan University	
	Yicheng Liu	Nanjing University	
	Yuping He	Nanjing University	
	Kanghua Pan	Nanjing University	
	Tong Lu	Nanjing University	
	Limin Wang	Shanghai AI Laboratory	
	Yali Wang	Shanghai AI Laboratory	
	Yu Qiao	Shanghai AI Laboratory	

EPIC@CVPR19

The fourth international workshop on Egocentric Perception, Interaction and Computing

EPIC@CVPR2020

The Sixth International Workshop on Egocentric Perception, Interaction and Computing

EPIC@CVPR2021

The Eighth International Workshop on Egocentric Perception, Interaction and Computing

EPIC@CVPR22

Tenth International Workshop on Egocentric Perception, Interaction and Computing
held in conjunction with the 1st Ego4D Workshop

EPIC@CVPR23

Monday 20th June 2022

First Joint Egocentric Vision (EgoVis) Workshop

Second Joint Egocentric Vision (EgoVis) Workshop

Held in Conjunction with CVPR 2025

12 June 2025 - Nashville, USA





Consortium



Ego4D: Around the World in 3,000 Hours of Egocentric Video 84 authors

Kristen Grauman^{1,2}, Andrew Westbury¹, Eugene Byrne^{*1}, Zachary Chavis^{*3}, Antonino Furnari^{*4}, Rohit Girdhar^{*1}, Jackson Hamburger^{*1}, Hao Jiang^{*5}, Miao Liu^{*6}, Xingyu Liu^{*7}, Miguel Martin^{*1}, Tushar Nagarajan^{*1,2}, Ilija Radosavovic^{*8}, Santhosh Kumar Ramakrishnan^{*1,2}, Fiona Ryan^{*6}, Jayant Sharma^{*3}, Michael Wray^{*9}, Mengmeng Xu^{*10}, Eric Zhongcong Xu^{*11}, Chen Zhao^{*10}, Siddhant Bansal¹⁷, Dhruv Batra¹, Vincent Cartillier^{1,6}, Sean Crane⁷, Tien Do³, Morrie Doulaty¹³, Akshay Erapalli¹³, Christoph Feichtenhofer¹, Adriano Fragomeni⁹, Qichen Fu⁷, Christian Fuegen¹³, Abraham Gebreselasie¹², Cristina González¹⁴, James Hillis⁵, Xuhua Huang⁷, Yifei Huang¹⁵, Wenqi Jia⁶, Weslie Khoo¹⁶, Jachym Kolar¹³, Satwik Kottur¹³, Anurag Kumar⁵, Federico Landini¹³, Chao Li⁵, Zhenqiang Li¹⁵, Karttikeya Mangalam^{1,8}, Raghava Modhugu¹⁷, Jonathan Munro⁹, Tullie Murrell¹, Takumi Nishiyasu¹⁵, Will Price⁹, Paola Ruiz Puentes¹⁴, Mery Ramazanova¹⁰, Leda Sari⁵, Kiran Somasundaram⁵, Audrey Southerland⁶, Yusuke Sugano¹⁵, Ruijie Tao¹¹, Minh Vo⁵, Yuchen Wang¹⁶, Xindi Wu⁷, Takuma Yagi¹⁵, Yunyi Zhu¹¹, Pablo Arbeláez¹⁴, David Crandall¹⁶, Dima Damen¹⁹, Giovanni Maria Farinella¹⁴, Bernard Ghanem¹⁰, Vamsi Krishna Ithapu¹⁵, C. V. Jawahar¹⁷, Hanbyul Joo¹¹, Kris Kitani¹⁷, Haizhou Li¹¹, Richard Newcombe¹⁵, Aude Oliva¹⁸, Hyun Soo Park¹³, James M. Rehg¹⁶, Yoichi Sato¹⁵, Jianbo Shi¹⁹, Mike Zheng Shou¹¹, Antonio Torralba¹⁸, Lorenzo Torresani^{11,20}, Mingfei Yan¹⁵, Jitendra Malik^{1,8}

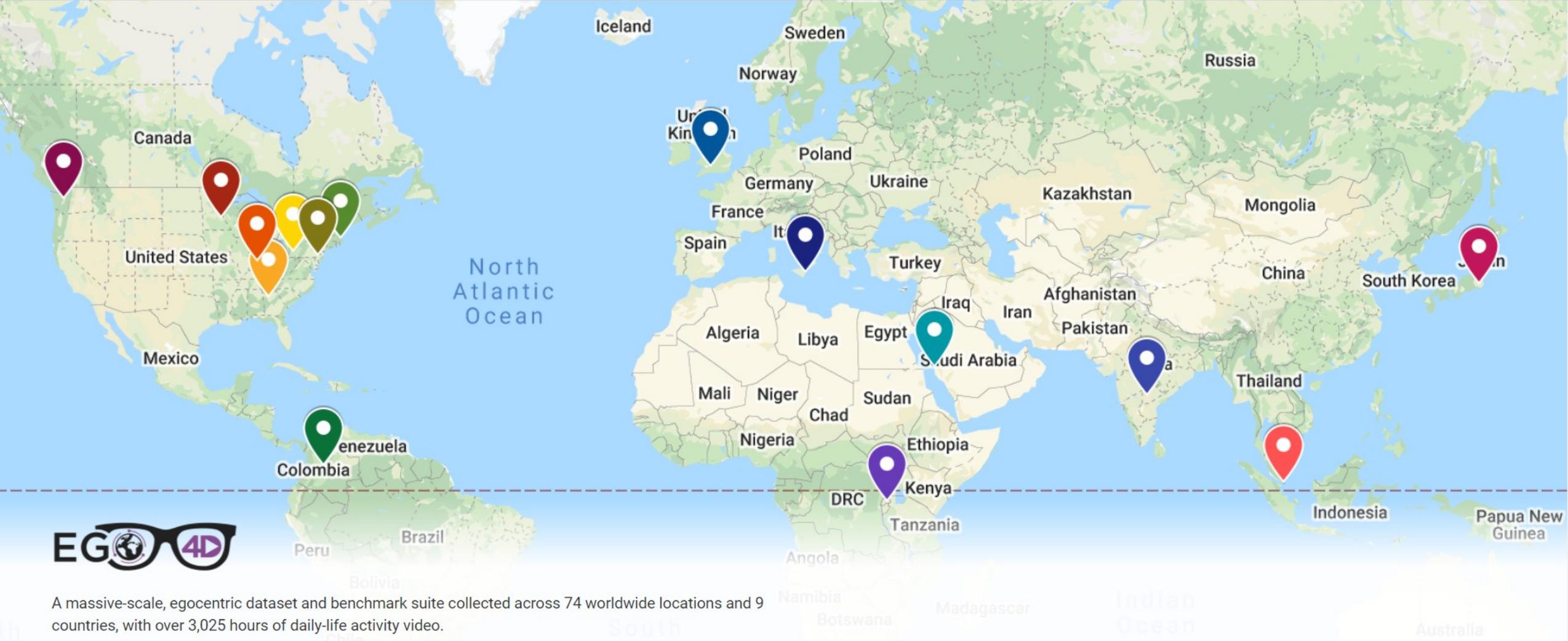
¹Facebook AI Research (FAIR), ²University of Texas at Austin, ³University of Minnesota, ⁴University of Catania,

⁵Facebook Reality Labs, ⁶Georgia Tech, ⁷Carnegie Mellon University, ⁸UC Berkeley, ⁹University of Bristol,

¹⁰King Abdullah University of Science and Technology, ¹¹National University of Singapore,

¹²Carnegie Mellon University Africa, ¹³Facebook, ¹⁴Universidad de los Andes, ¹⁵University of Tokyo, ¹⁶Indiana University,

¹⁷International Institute of Information Technology, Hyderabad, ¹⁸MIT, ¹⁹University of Pennsylvania, ²⁰Dartmouth



855 Subjects



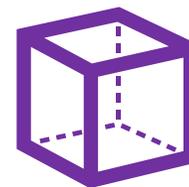
74 Locations



9 Countries



3025 Hours



3D Scans



Audio



Gaze

 120 Parts.
120 hours

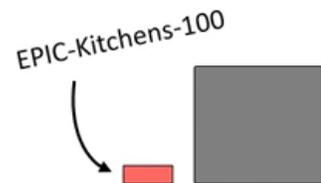
Ego4D – A Massive-Scale Egocentric Dataset

3,025 Hours

855 Participants

5 Benchmark Tasks

Find out more: <https://ego4d-data.org/>



Animation by Michael Wray – <https://mwrap.github.io>

Animation by Michael Wray - <https://www.youtube.com/watch?v=p78-V2RiKo>



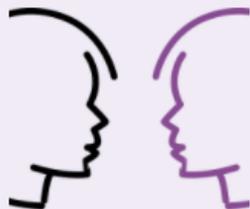
Episodic Memory



**Hand-Object
Interactions**



AV Diarization



Social



Forecasting

1st Ego4D Workshop @ CVPR 2022

2nd International Ego4D Workshop @ ECCV 2022

3rd International Ego4D Workshop @ CVPR 2023

First Joint Egocentric Vision (EgoVis) Workshop

Second Joint Egocentric Vision (EgoVis) Workshop

Held in Conjunction with CVPR 2025

12 June 2025 - Nashville, USA

Room: Grand B1



EGO-EXO4D





Universidad de los Andes

Colombia



SIMON FRASER UNIVERSITY

Carnegie Mellon University



Università di Catania

Carnegie Mellon University Africa



NUS

National University of Singapore

Meta



INTERNATIONAL INSTITUTE OF INFORMATION TECHNOLOGY

HYDERABAD



東京大学

THE UNIVERSITY OF TOKYO



جامعة الملك عبد الله للعلوم والتقنية

King Abdullah University of Science and Technology



UNIVERSITY OF MINNESOTA



Penn

UNIVERSITY of PENNSYLVANIA



University of BRISTOL



INDIANA UNIVERSITY BLOOMINGTON



THE UNIVERSITY of NORTH CAROLINA at CHAPEL HILL

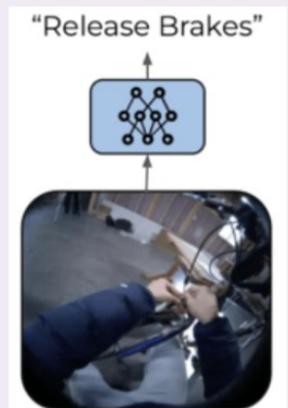


UNIVERSITY OF ILLINOIS URBANA-CHAMPAIGN

Ego-Exo4D: Understanding Skilled Human Activity from First- and Third-Person Perspectives

Kristen Grauman^{1,2}, Andrew Westbury¹, Lorenzo Torresani¹, Kris Kitani^{1,3}, Jitendra Malik^{1,4},
 Triantafyllos Afouras^{*1}, Kumar Ashutosh^{*1,2}, Vijay Baiyya^{*5}, Siddhant Bansal^{*6,7}, Bikram Boote^{*8},
 Eugene Byrne^{*1,9}, Zach Chavis^{*10}, Joya Chen^{*11}, Feng Cheng^{*1}, Fu-Jen Chu^{*1}, Sean Crane^{*9}, Avijit
 Dasgupta^{*7}, Jing Dong^{*5}, Maria Escobar^{*12}, Cristhian Forigua^{*12}, Abrahm Gebreselasie^{*9}, Sanjay
 Haresh^{*13}, Jing Huang^{*1}, Md Mohaiminul Islam^{*14}, Suyog Jain^{*1}, Rawal Khirodkar^{*9}, Devansh
 Kukreja^{*1}, Kevin J Liang^{*1}, Jia-Wei Liu^{*11}, Sagnik Majumder^{*1,2}, Yongsen Mao^{*13}, Miguel Martin^{*1},
 Effrosyni Mavroudi^{*1}, Tushar Nagarajan^{*1}, Francesco Ragusa^{*15}, Santhosh Kumar Ramakrishnan^{*2},
 Luigi Seminara^{*15}, Arjun Somayazulu^{*2}, Yale Song^{*1}, Shan Su^{*16}, Zihui Xue^{*1,2}, Edward Zhang^{*16},
 Jinxu Zhang^{*16}, Angela Castillo¹², Changan Chen², Xinzhu Fu¹¹, Ryosuke Furuta¹⁷, Cristina
 González¹², Prince Gupta⁵, Jiabo Hu¹⁸, Yifei Huang¹⁷, Yiming Huang¹⁶, Weslie Khoo¹⁹, Anush
 Kumar¹⁰, Robert Kuo¹⁸, Sach Lakhavani⁵, Miao Liu¹⁸, Mi Luo², Zhengyi Luo³, Brigid Meredith¹⁸,
 Austin Miller¹⁸, Oluwatumininu Oguntola¹⁴, Xiaqing Pan⁵, Penny Peng¹⁸, Shraman Pramanick²⁰,
 Merey Ramazanova²¹, Fiona Ryan²², Wei Shan¹⁴, Kiran Somasundaram⁵, Chenan Song¹¹, Audrey
 Southerland²², Masatoshi Tateno¹⁷, Huiyu Wang¹, Yuchen Wang¹⁹, Takuma Yagi¹⁷, Mingfei Yan⁵,
 Xitong Yang¹, Zecheng Yu¹⁷, Shengxin Cindy Zha¹⁸, Chen Zhao²¹, Ziwei Zhao¹⁹, Zhifan Zhu⁶, Jeff
 Zhuo¹⁴, Pablo Arbeláez^{†12}, Gedas Bertasius^{†14}, David Crandall^{†19}, Dima Damen^{†6}, Jakob Engel^{†5},
 Giovanni Maria Farinella^{†15}, Antonino Furnari^{†15}, Bernard Ghanem^{†21}, Judy Hoffman^{†22}, C. V.
 Jawahar^{†7}, Richard Newcombe^{†5}, Hyun Soo Park^{†10}, James M. Rehg^{†8}, Yoichi Sato^{†17}, Manolis
 Savva^{†13}, Jianbo Shi^{†16}, Mike Zheng Shou^{†11}, and Michael Wray^{†6}

<https://ego-exo4d-data.org/>



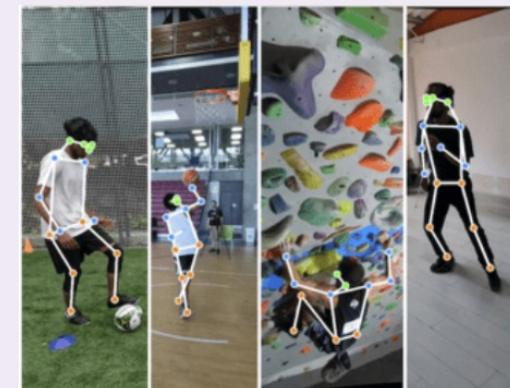
Keystep Recognition



Proficiency Estimation



Relation



Pose Estimation

Third Joint Egocentric Vision (EgoVis) Workshop

Held in Conjunction with CVPR 2026

3/4 June 2026 - Denver, CO, USA



Ego-Exo4D



Ego4D



EPIC-Kitchens

The MECCANO Dataset: Understanding Human-Object Interactions from Egocentric Videos in an Industrial-like Domain

F. Ragusa^{1,3}, A. Furnari¹, S. Livatino², G. M. Farinella¹

¹IPLab, Department of Mathematics and Computer Science - University of Catania, IT

²University of Hertfordshire, Hatfield, Hertfordshire, U.K.

³Xenia Gestione Documentale s.r.l. - Xenia Progetti s.r.l., Acicastello, Catania, IT

The new version of MECCANO is available here!

Assembly101: A Large-Scale Multi-View Video Dataset for Understanding Procedural Activities

Fadime Sener¹

Dibyadip Chatterjee²

Daniel Shelepov¹

Kun He¹

Dipika Singhania²

Robert Wang¹

Angela Yao²

¹Reality Labs at Meta

²National University of Singapore

CVPR 2022

[Paper](#)
[Dataset](#)
[Code](#)
[Sample](#)
[Codalab Challenge](#)



IndustReal: A Dataset for Procedure Step Recognition Handling Execution Errors in Egocentric Videos in an Industrial-Like Setting

Tim J. Schoonbeek¹, Tim Houben¹, Hans Onvlee², Peter H.N. de With¹, Fons van der Sommen¹,

¹Eindhoven University of Technology, ²ASML Research

Published in: WACV 2024

[Paper](#)
[arXiv](#)
[Video](#)
[Code](#)
[Data](#)
[Poster](#)



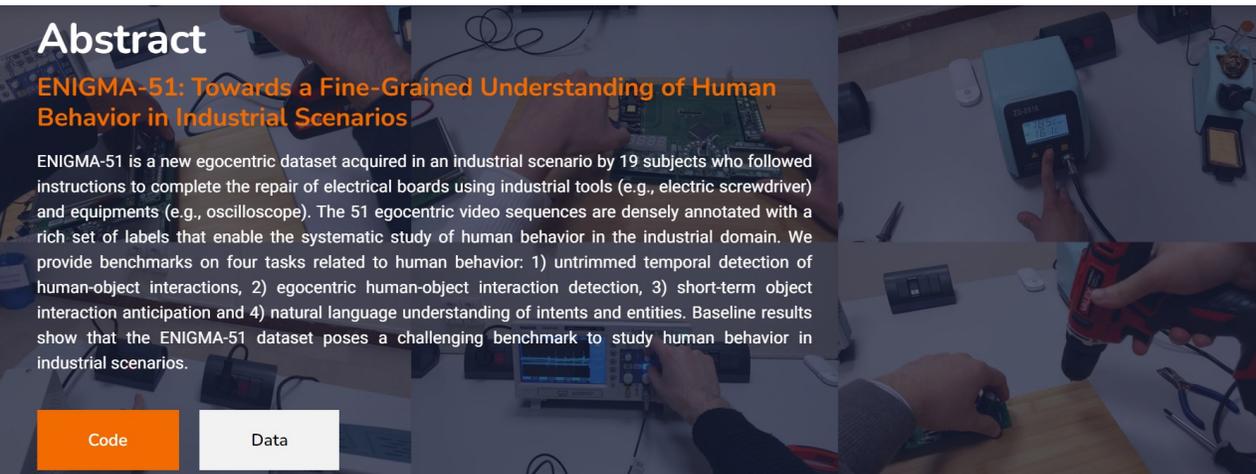
Abstract

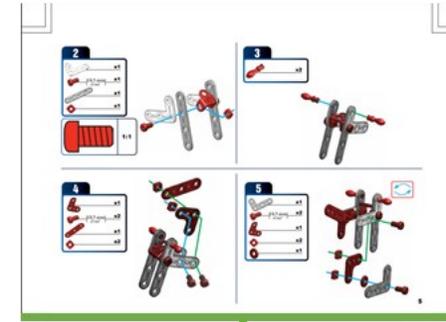
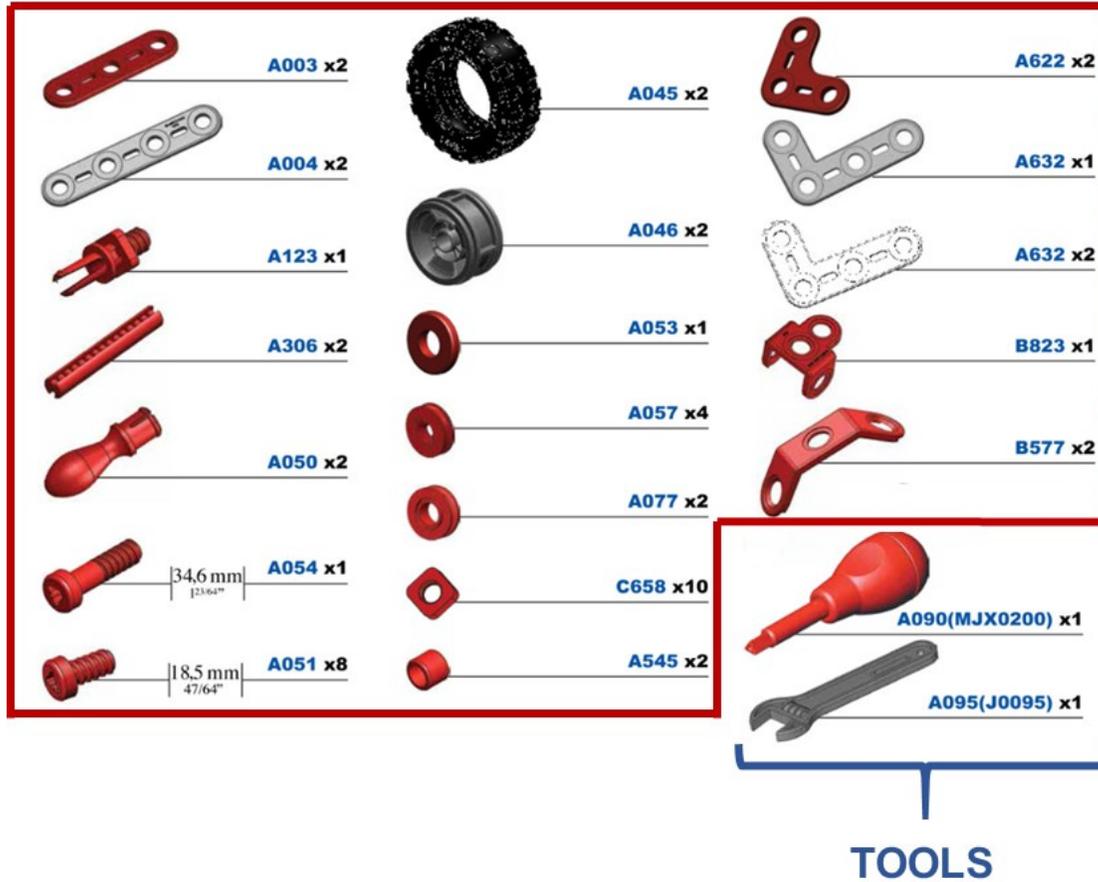
ENIGMA-51: Towards a Fine-Grained Understanding of Human Behavior in Industrial Scenarios

ENIGMA-51 is a new egocentric dataset acquired in an industrial scenario by 19 subjects who followed instructions to complete the repair of electrical boards using industrial tools (e.g., electric screwdriver) and equipments (e.g., oscilloscope). The 51 egocentric video sequences are densely annotated with a rich set of labels that enable the systematic study of human behavior in the industrial domain. We provide benchmarks on four tasks related to human behavior: 1) untrimmed temporal detection of human-object interactions, 2) egocentric human-object interaction detection, 3) short-term object interaction anticipation and 4) natural language understanding of intents and entities. Baseline results show that the ENIGMA-51 dataset poses a challenging benchmark to study human behavior in industrial scenarios.

Code

Data



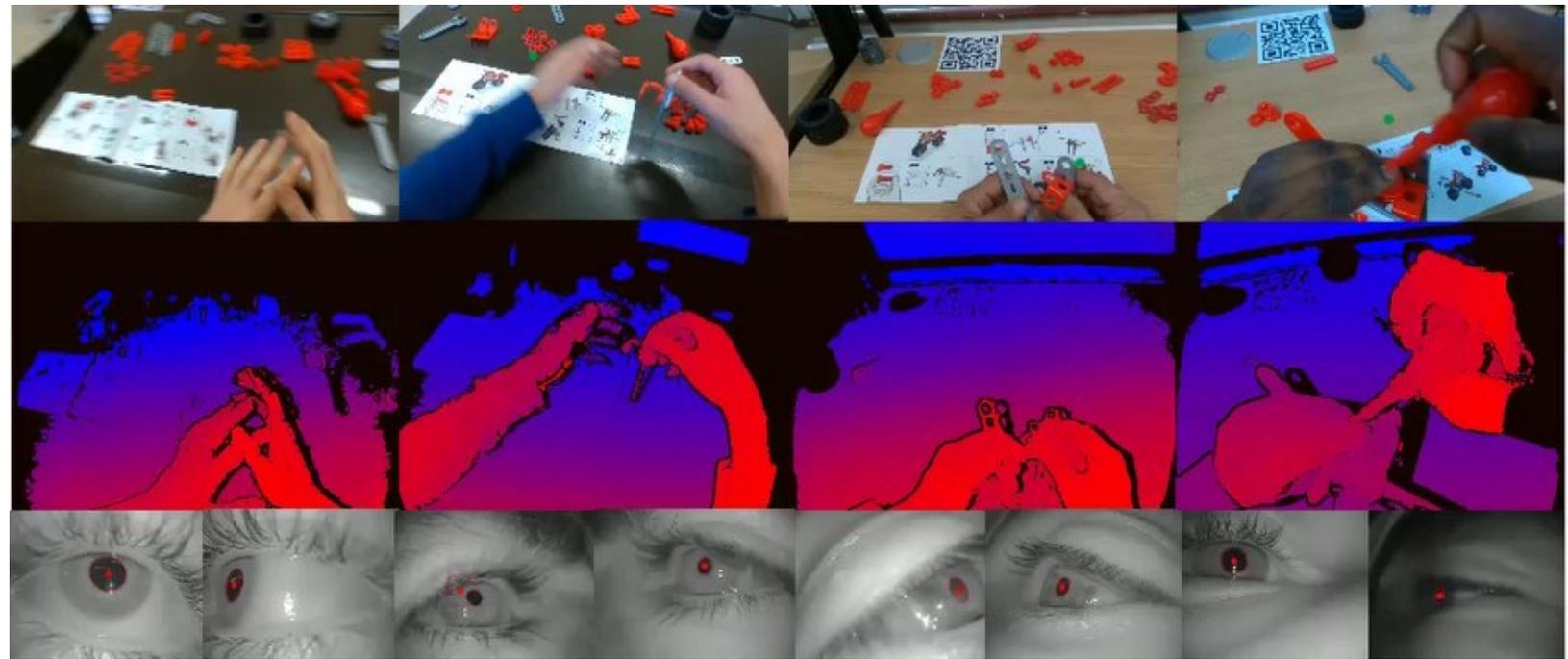
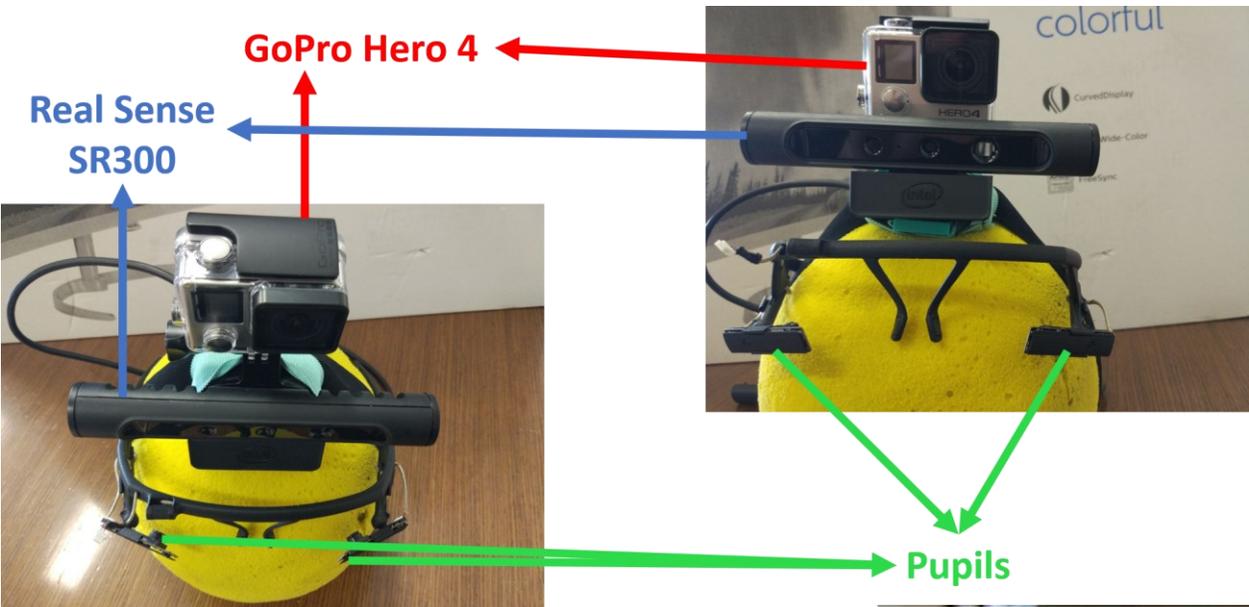


BOOKLET

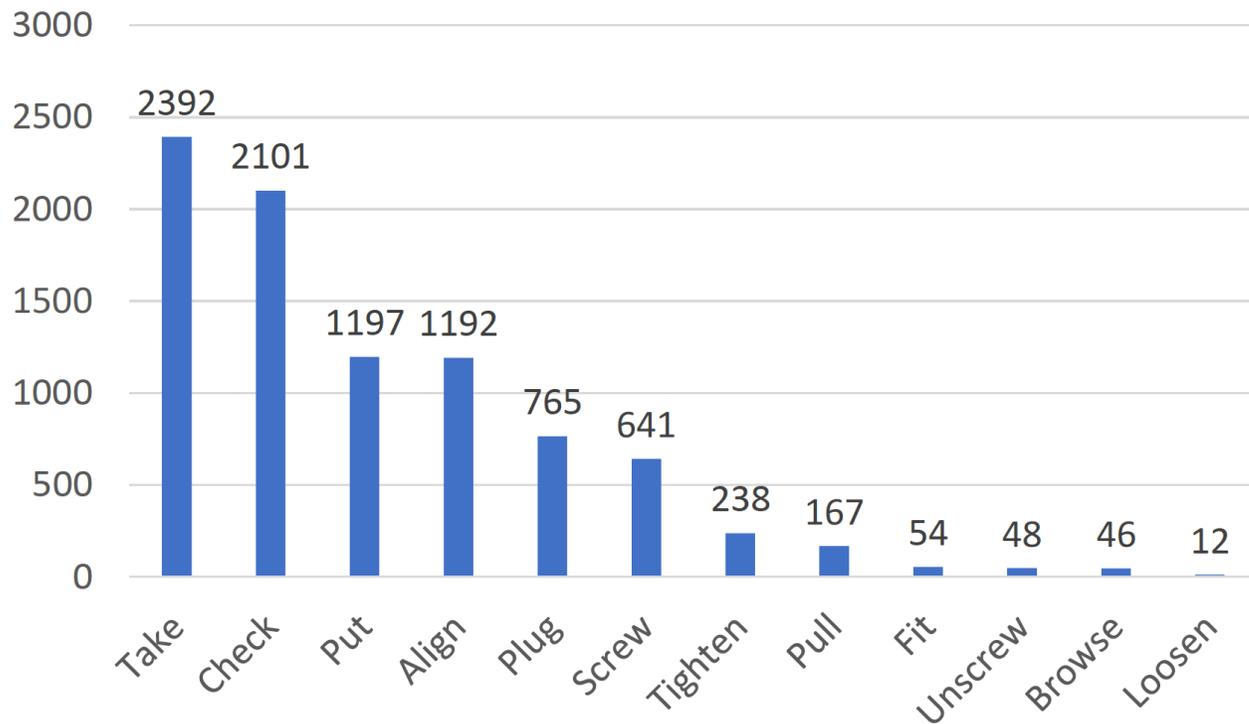
COMPONENTS



[Project page:](https://iplab.dmi.unict.it/MECCANO/)
<https://iplab.dmi.unict.it/MECCANO/>

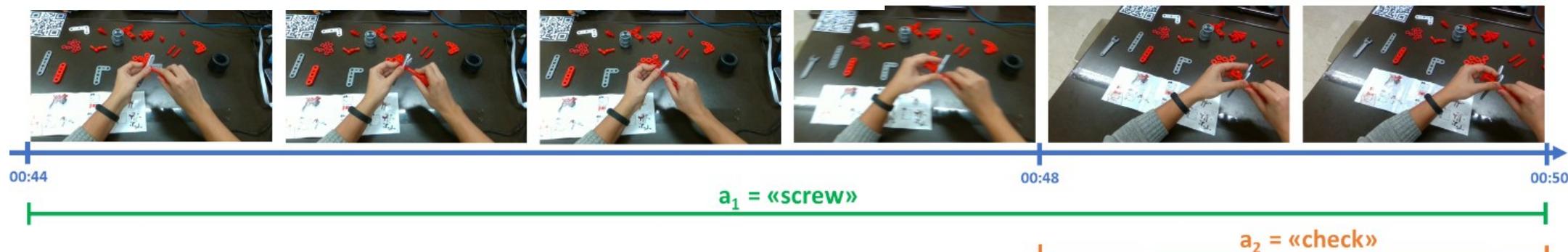


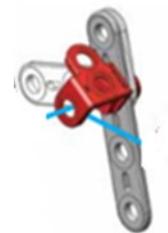
Verbs Classes



8857 video segments

1401 overlap segments (15.82%)





red_perforated_bar

gray_bar

wheels_axle

bar

handlebar

partial_model

gray_angled_bar

bolt

red_3_junction_bar

wrench



tire

rim

washer

white_bar

instruction_booklet

cylinder

red_angled_bar

screw

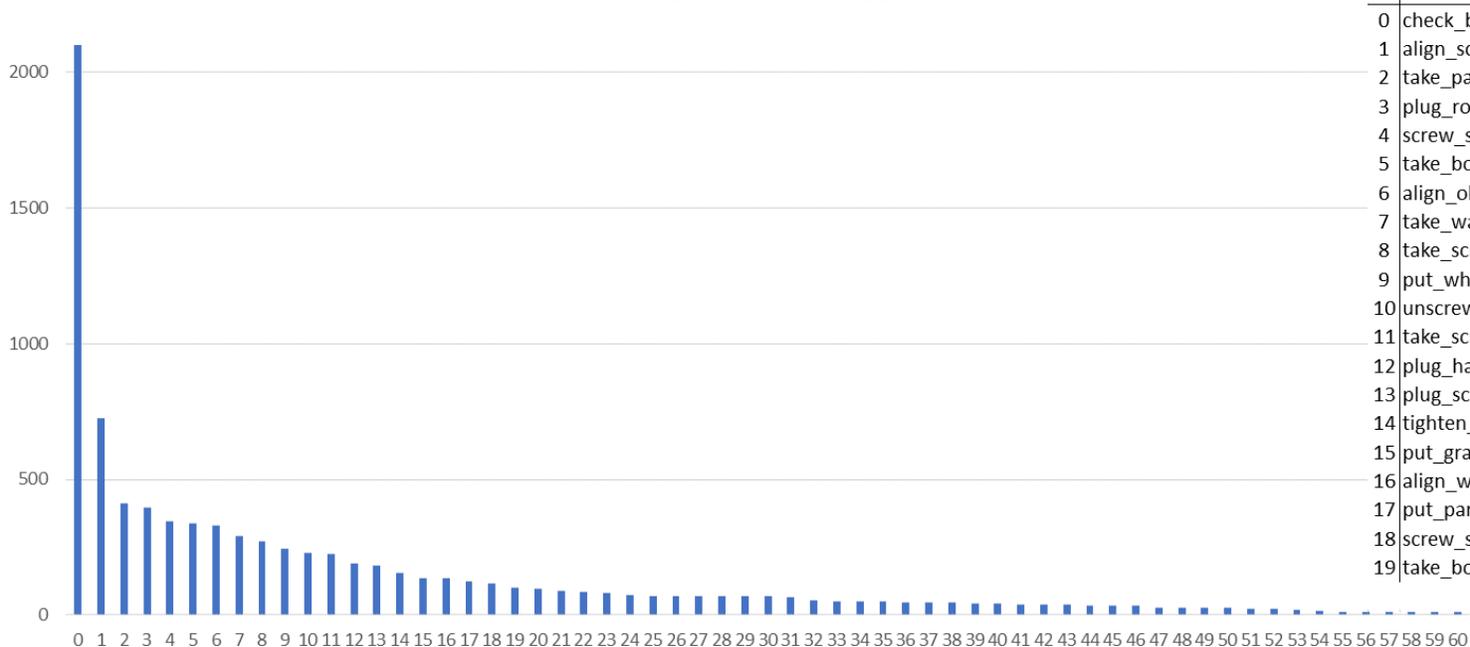
red_4_junction_bar

screwdriver



**64439
frames**

Action instances



ID	Action
0	check_booklet
1	align_screwdriver_to_screw
2	take_partial_model
3	plug_rod
4	screw_screw_with_screwdriver
5	take_bolt
6	align_objects
7	take_washer
8	take_screw
9	put_white_angled_perforated_bar
10	unscrew_screw_with_hands
11	take_screwdriver
12	plug_handlebar
13	plug_screw
14	tighten_nut_with_wrench
15	put_gray_perforated_bar
16	align_wrench_to_bolt
17	put_partial_model
18	screw_screw_with_hands
19	take_booklet

ID	Action
20	put_screwdriver
21	put_red_perforated_junction_bar
22	put_gray_angled_perforated_bar
23	take_red_perforated_bar
24	take_gray_perforated_bar
25	take_red_angled_perforated_bar
26	tighten_nut_with_hands
27	take_white_angled_perforated_bar
28	take_rod
29	put_tire
30	put_roller
31	pull_partial_model
32	pull_screw
33	take_gray_angled_perforated_bar
34	take_tire
35	pull_rod
36	take_wrench
37	browse_booklet
38	take_roller
39	take_handlebar

ID	Action
40	take_red_perforated_junction_bar
41	fit_rim_tire
42	take_rim
43	take_red_4_perforated_junction_bar
44	put_screw
45	put_rod
46	put_washer
47	unscrew_screw_with_screwdriver
48	put_red_perforated_bar
49	put_wrench
50	put_bolt
51	take_wheels_axle
52	put_wheels_axle
53	put_red_angled_perforated_bar
54	put_red_4_perforated_junction_bar
55	take_objects
56	put_objects
57	loosen_bolt_with_hands
58	put_booklet
59	put_rim
60	put_handlebar

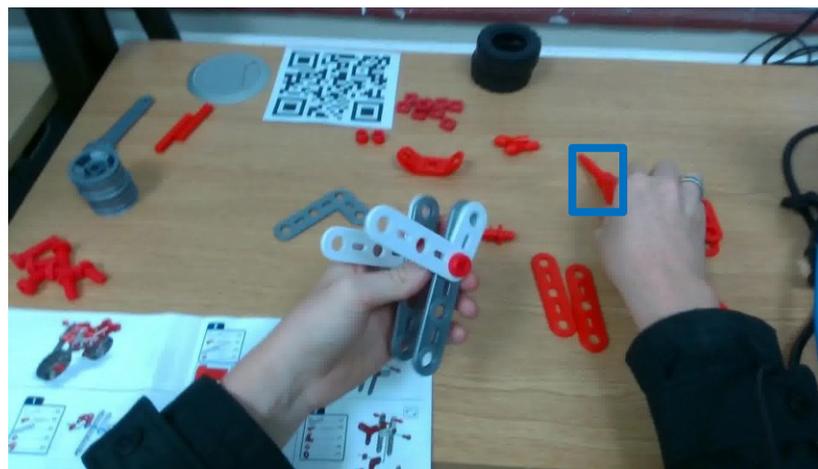
align screadriver **to** screw

Egocentric Human-Object Interaction

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

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

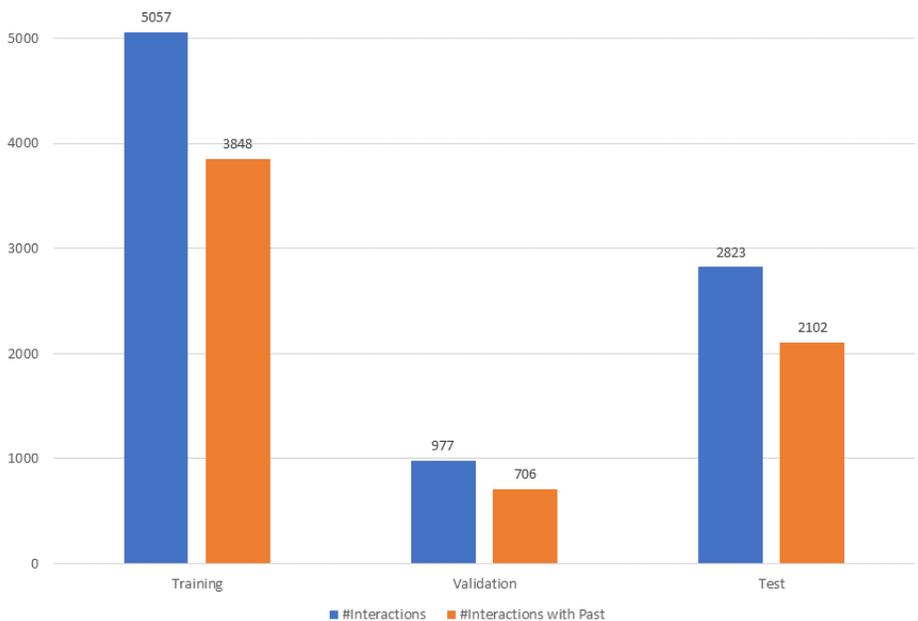
$$e = (v_h, \{o_1, o_2, \dots, o_i\})$$



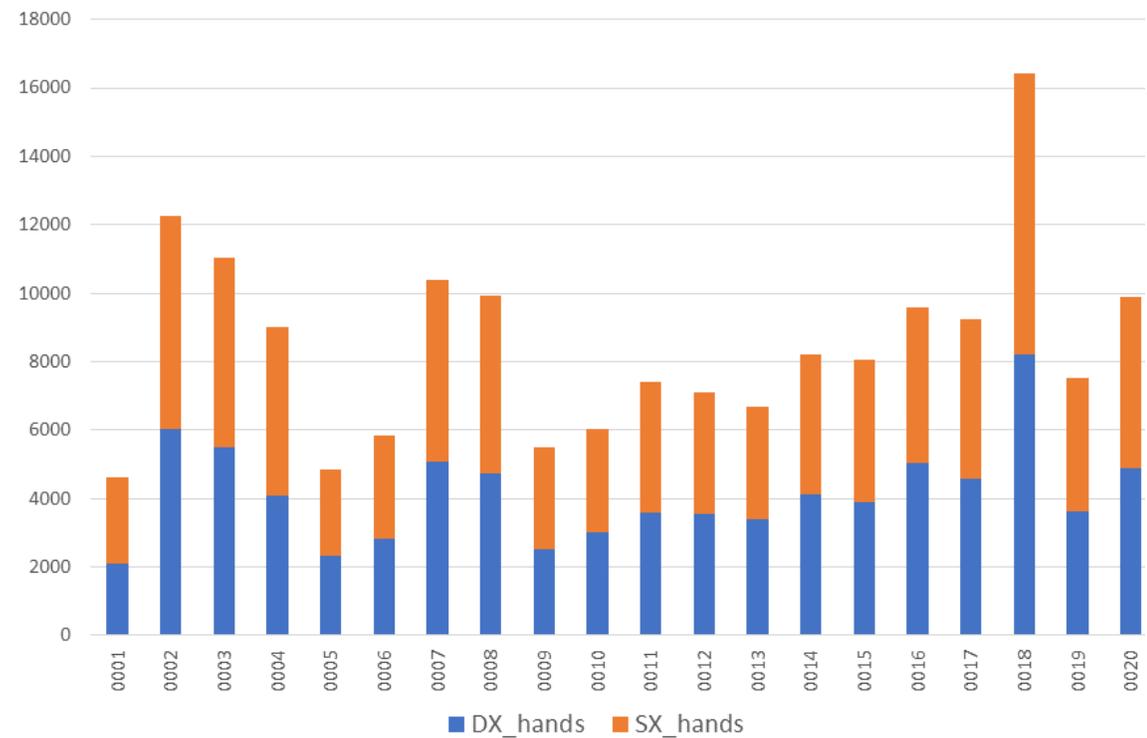
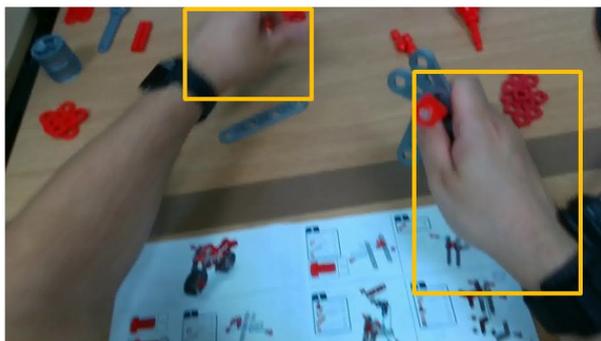
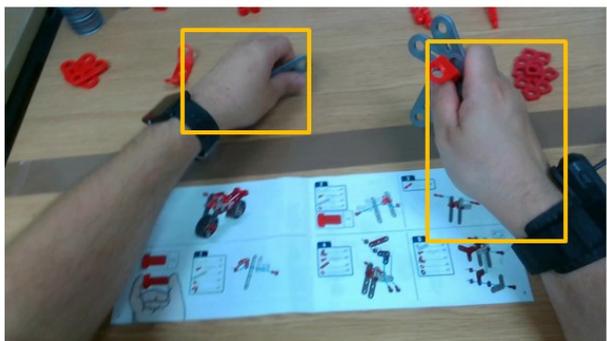
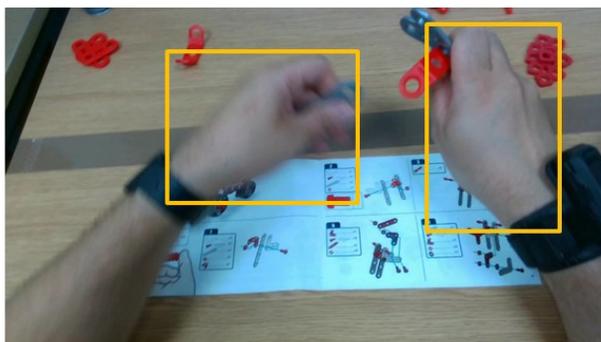
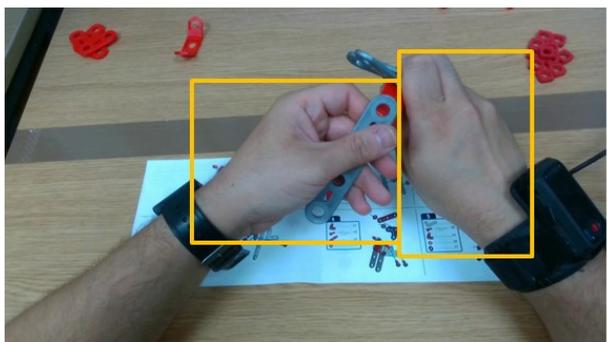
<take, screwdriver>



<screw, {screwdriver, screw, partial_model}>



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



1) Action Recognition

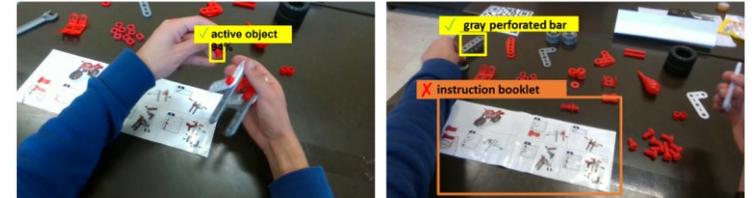
start frame

end frame

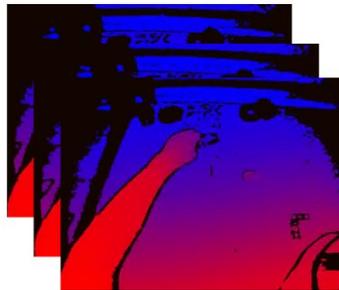
GT
RGB+Gaze
Depth+Gaze
All

take screwdriver	align objects	take screwdriver	take screwdriver
---------------------	------------------	---------------------	---------------------

2) Active Object Detection and Recognition



3) EHOI Detection



<take>



<gray perforated bar>

4) Action Anticipation

Ground Truth action: **take bolt**

$\tau_a = 2.00$



take bolt, align objects, tighten **bolt**, plug screw, check booklet

$\tau_a = 1.50$



take bolt, align objects, plug screw, tighten **bolt**, check booklet

$\tau_a = 1.00$



take bolt, align objects, plug screw, check booklet, tighten **bolt**

$\tau_a = 0.25$

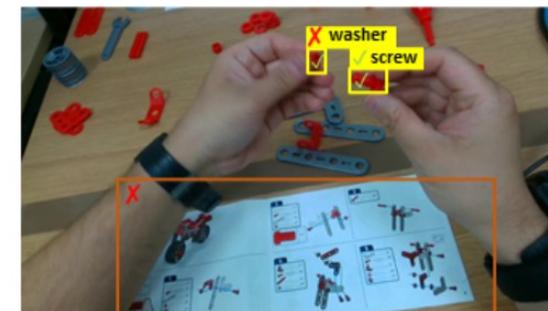


take bolt, align objects, plug screw, check booklet, take screwdriver

5) Next-Active Object (NAO) Detection

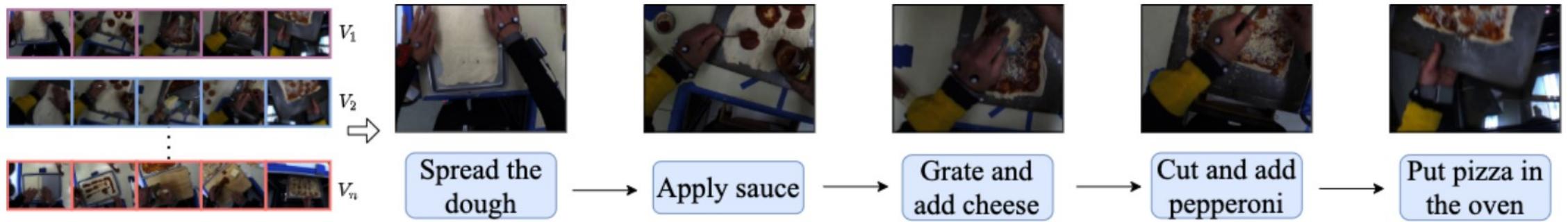


Time to start = 1.6s



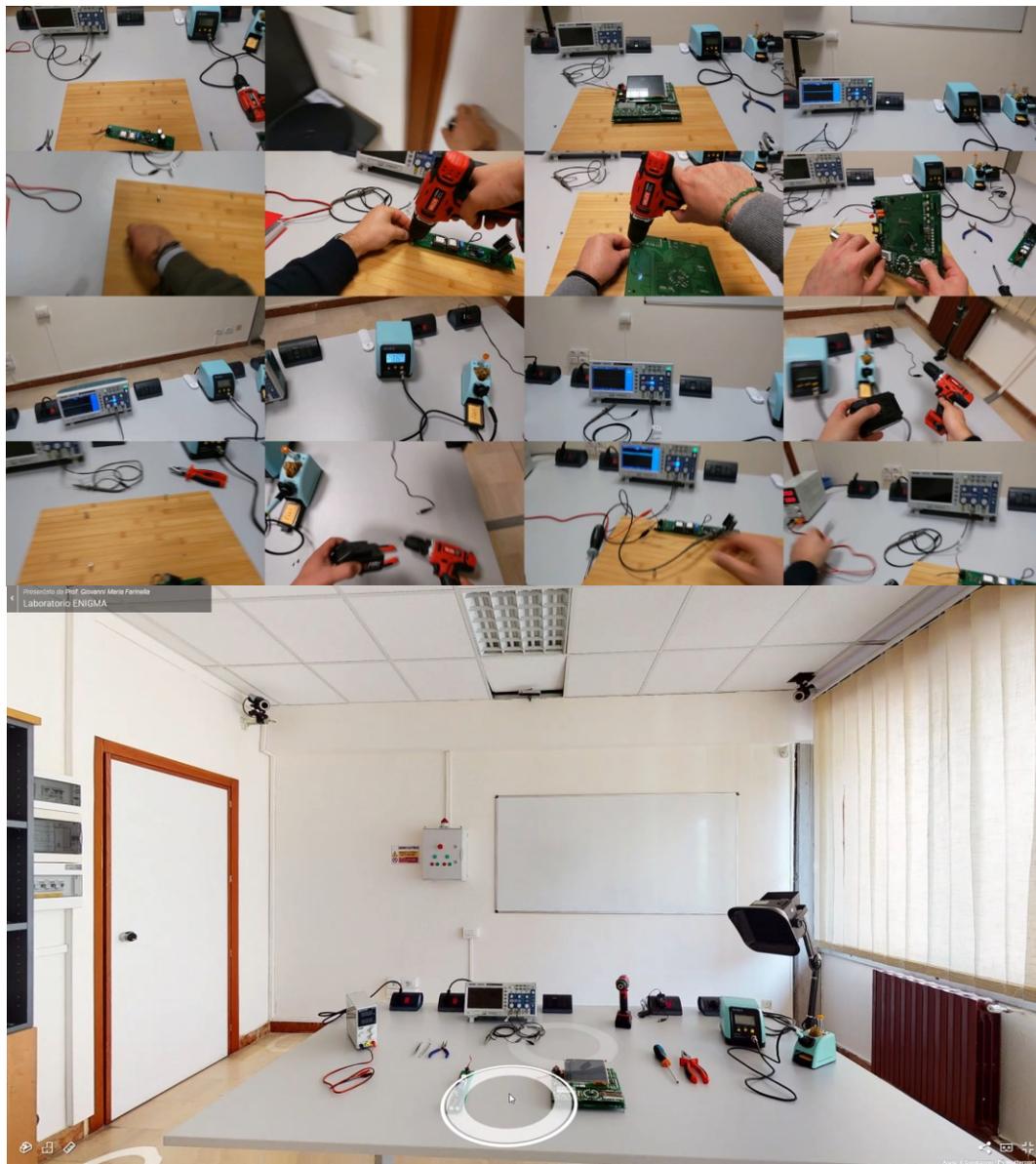
Time to start = 0.8s

Given multiple videos of a task, the goal is to identify the key-steps and their order to perform the task.



- 1) EgoProceL (proposed)
- 2) CMU-MMAC
- 3) EGTEA Gaze+

- 4) MECCANO
- 5) EPIC-Tent

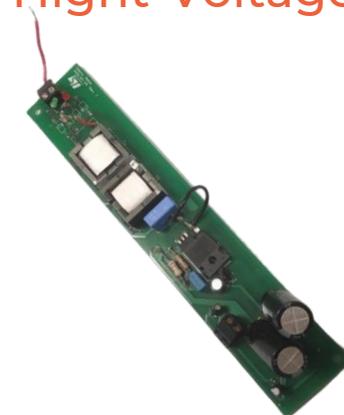


We designed two procedures consisting of instructions that involve humans interacting with the objects present in the laboratory to achieve the goal of repairing two electrical boards

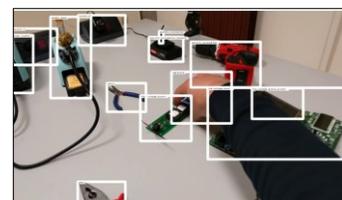
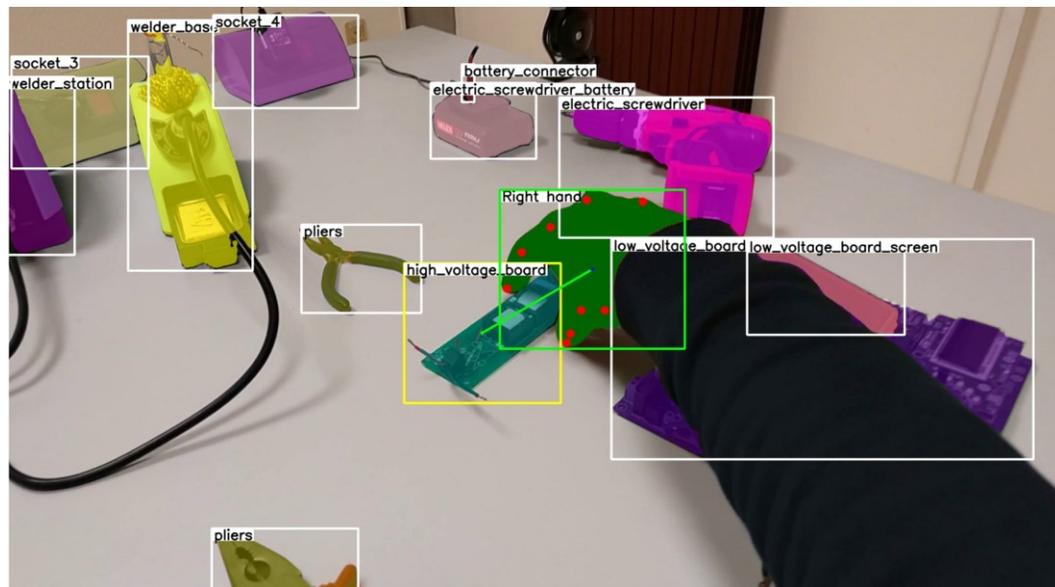
Low-Voltage



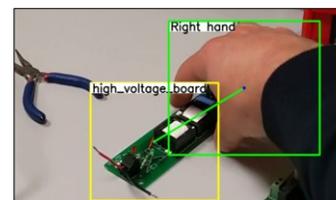
High-Voltage



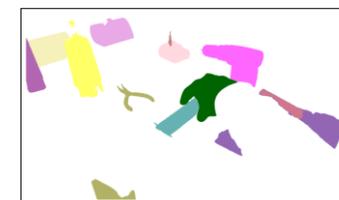
ENIGMA-51: Towards a Fine-Grained Understanding of Human Behavior in Industrial Scenarios. F. Ragusa R. Leonardi, M. Mazzamuto, C. Bonanno, R. Scavo, A. Furnari, G. M. Farinella. WACV (2024).



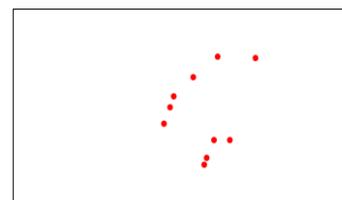
Hand-Object boxes



Human-Object Interactions



Hand-Object Masks



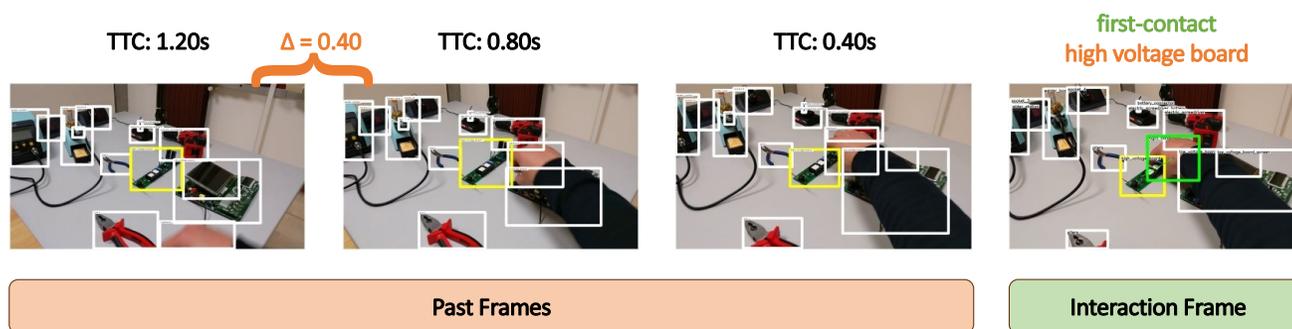
Hand Keypoints



Environment 3D Model



Object 3D Models



Procedure :

.....

4. Take the high voltage board and put it on the working area

5. Take the screwdriver

.....

22. Turn on the welder using the switch on the corresponding socket (second from right)

23. Set the temperature of the welder to 480 °C using the yellow "UP" button

.....

Untrimmed temporal detection of human-object interactions

Egocentric human-object interaction detection

Short-term object interaction anticipation

Natural language understanding of intents and entities



35
subjects



7
tasks

Exocentric



Egocentric



- **Temporal Action Segmentation**
- **Keystep Recognition**
- **Mistake Detection**



F. Ragusa, M. Mazzamuto, R. Forte, I. D'Ambra, J. Fort, J. Engel, A. Furnari, G. M. Farinella (2026). Ego-EXTRA: video-language Egocentric Dataset for EXpert-TRAinee assistance. In IEEE Winter Conference on Application of Computer Vision (WACV).

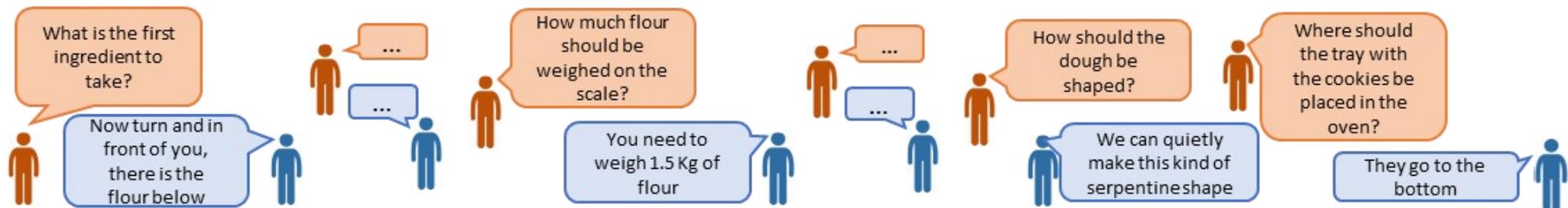
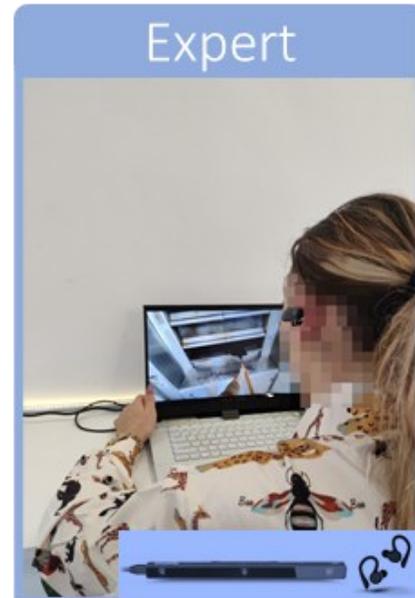


The Ideal Personal Assistant





The Ego-Extra Dataset

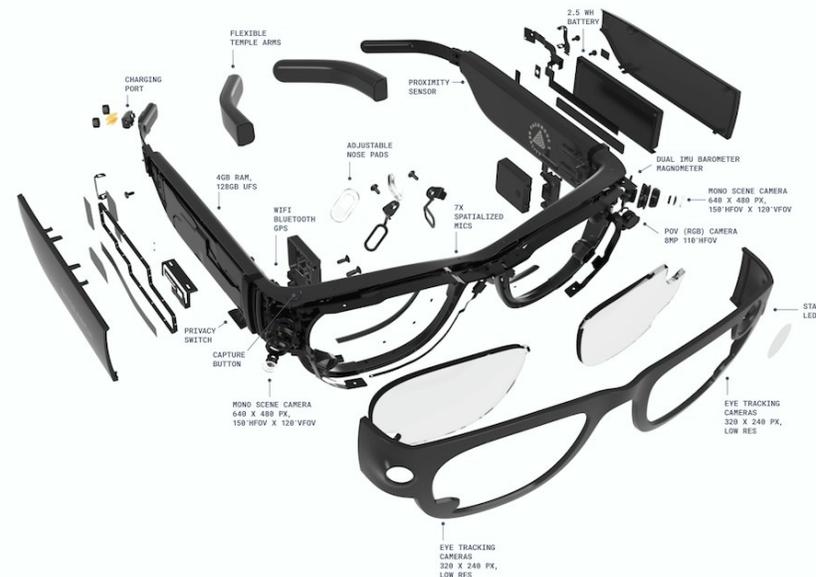


Time →

Trainee: RGB videos, Gaze, SLAM, Hand poses

Expert: Gaze

Trainee-Expert: Text (transcriptions)



Expert: Now you need to fix the electric board to the working area

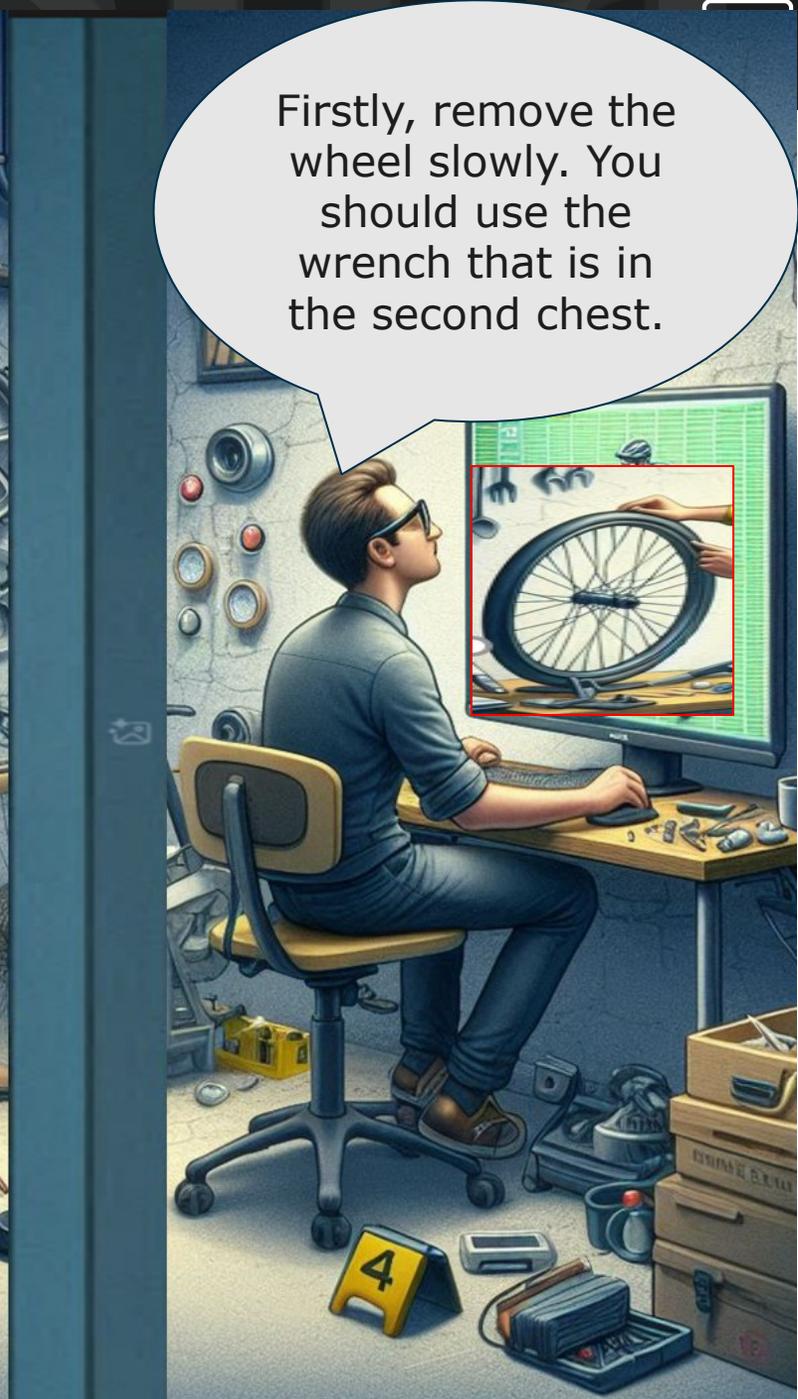
Trainee: With the screwdriver?

Expert: Yes

Pro-active:

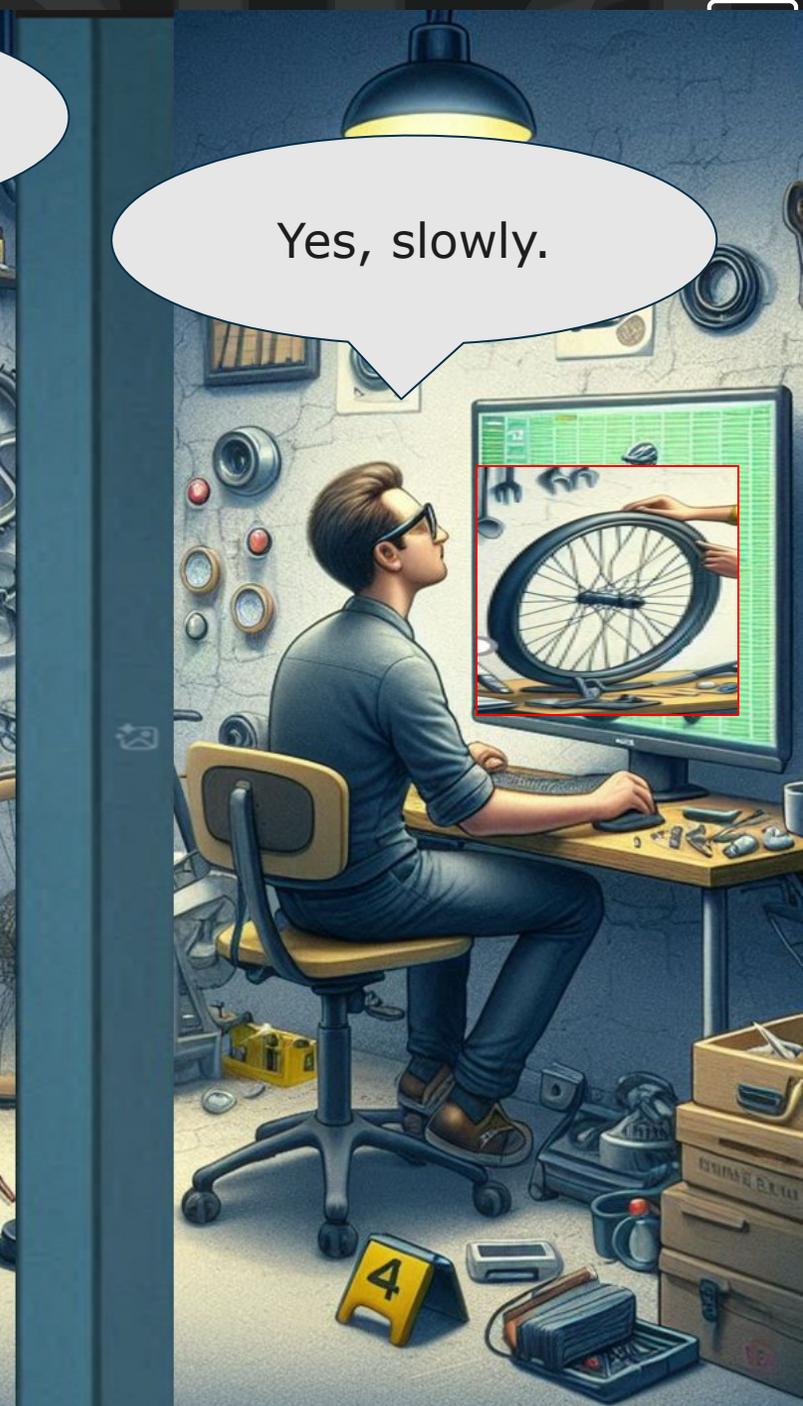
At the beginning, the trainee has not a knowledge about the environment, the objects and the procedure to perform.

The expert speaks freely with the trainee, suggesting next steps, instructions and anything that may be useful.

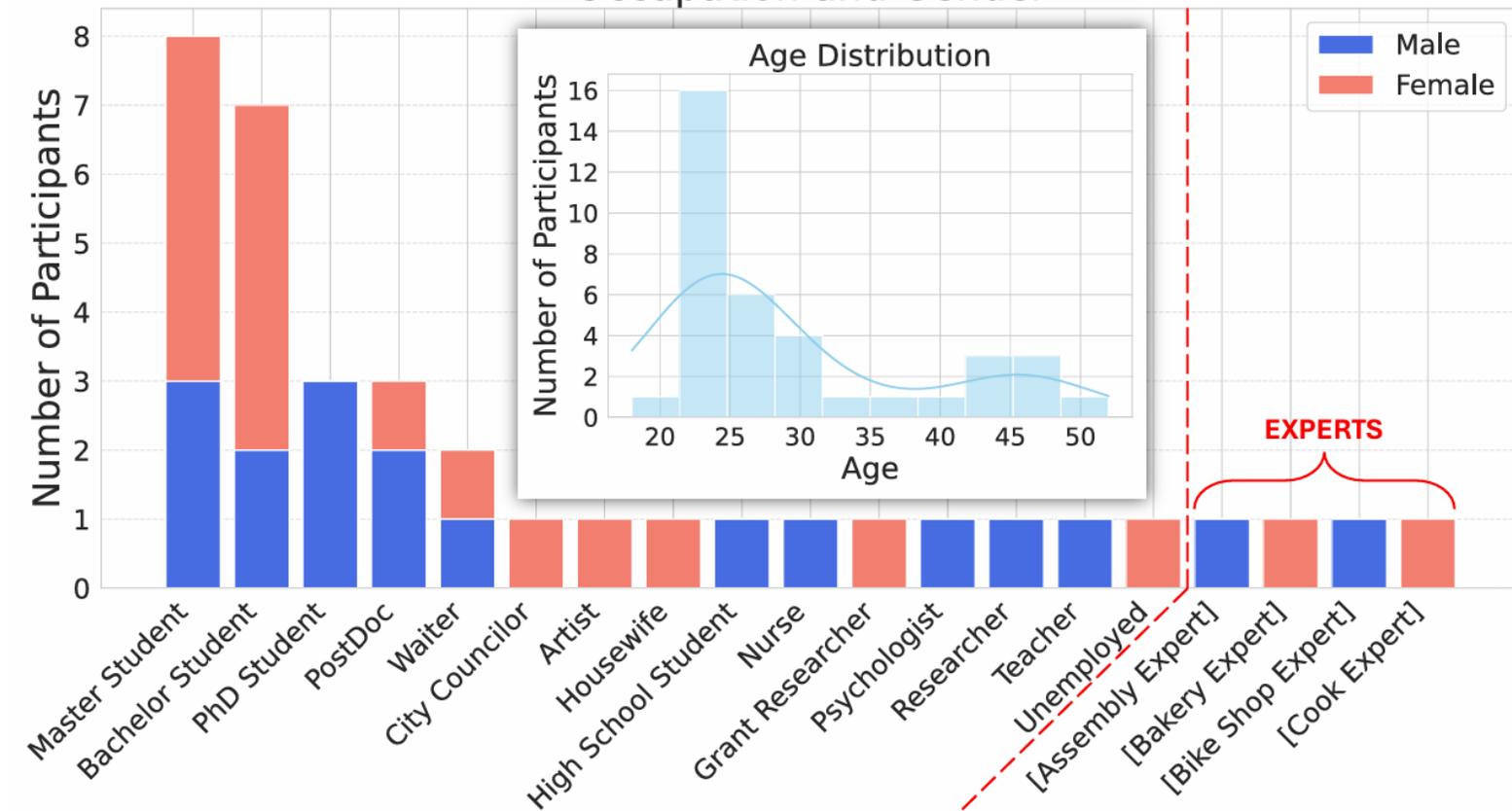


On-Demand:

The expert may answer only the trainee's questions or alert him if a mistake is happening.



Occupation and Gender



50 Hours



4 Scenarios



1 expert for each scenario



8-10 real trainees for each scenario

Bike Workshop

Procedures:

1. Lubrification brakes
2. Air Chamber change



Bakery

Procedures:

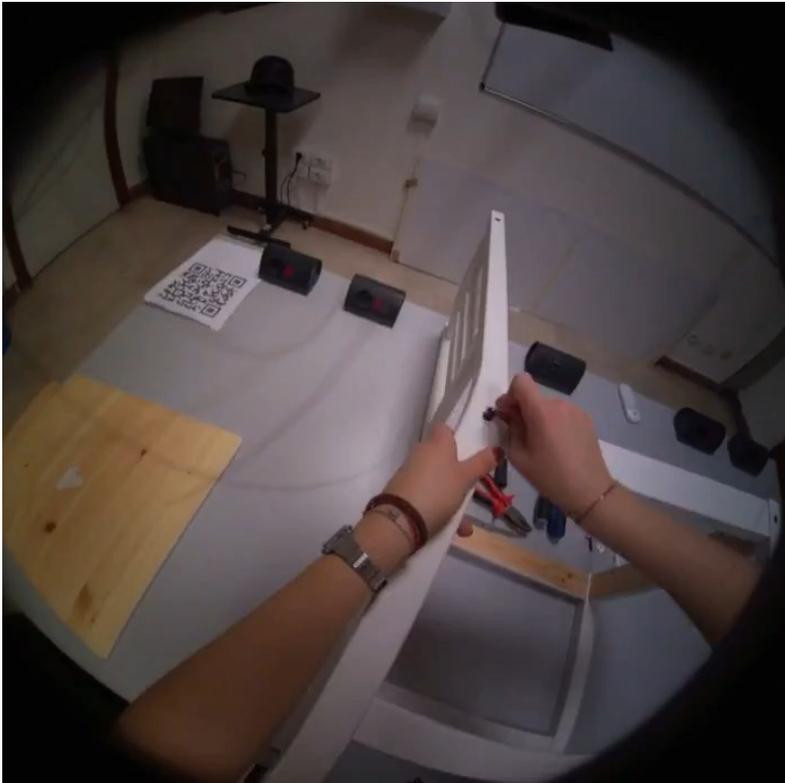
1. Make tart
2. Make biscuits



Assembly

Procedures:

1. Assemble a chair
2. Assemble chest of drawers



Kitchens

Procedures:

1. Make lasagna
2. Make vegetarian meat balls



Step1: Initial QA set Extraction

E: There are some wooden pegs.

T: Yes, what should I do with the wooden pegs?

E: You can insert them into the large holes.



Q: What should I do with the wooden pegs?

A: Insert the wooden pegs into the big holes.

B: Use the wooden pegs as reference.

C: Give the wooden pegs to someone else.

D: Put the wooden pegs in a corner.

E: Use the hammer to break the wooden pegs.

Step1: Initial QA set Extraction

E: There are some wooden pegs.
 T: Yes, what should I do with the wooden pegs?
 E: You can insert them into the large holes.



Q: What should I do with the wooden pegs?

A: Insert the wooden pegs into the big holes.

- B: Use the wooden pegs as reference.
- C: Give the wooden pegs to someone else.
- D: Put the wooden pegs in a corner.
- E: Use the hammer to break the wooden pegs.

Step2: Human Validation



Q: What tool should I use to tighten the black bolt?

Accept ✓

Q: What is the expert suggesting will help with the current task?

Discard ✗

Q: What should I use to pry open the package?

- A: A screwdriver
- B: A hammer
- C: Pliers
- D: A wrench
- E: A genevile (a type of lever) ←

Transcription Error ✗

Q: Am I correctly holding the pieces in place?

- A: No, I need to rotate them first
- B: I'm not sure, I need more guidance
- C: I'm holding them too loosely
- D: Yes, I'm holding them tightly
- E: I'm holding the wrong pieces

To Revise ✗

Step1: Initial QA set Extraction

E: There are some wooden pegs.
 T: Yes, what should I do with the wooden pegs?
 E: You can insert them into the large holes.



Q: What should I do with the wooden pegs?

A: Insert the wooden pegs into the big holes.
 B: Use the wooden pegs as reference.
 C: Give the wooden pegs to someone else.
 D: Put the wooden pegs in a corner.
 E: Use the hammer to break the wooden pegs.

Step2: Human Validation



Q: What tool should I use to tighten the black bolt?

Accept ✓

Q: What is the expert suggesting will help with the current task?

Discard ✗

Q: What should I use to pry open the package?

A: A screwdriver
 B: A hammer
 C: Pliers
 D: A wrench
 E: A genevile (a type of lever) ←

Transcription Error ✗

Q: Am I correctly holding the pieces in place?

A: No, I need to rotate them first
 B: I'm not sure, I need more guidance
 C: I'm holding them too loosely
 D: Yes, I'm holding them tightly
 E: I'm holding the wrong pieces

To Revise ✗

Step3: Video Grounding Validation



Q: What should I do with the part that is a bit open?

A: Loosen the bolts
 B: Tighten it more with the bolts
 C: Leave it as it is
 D: Use a different tool
 E: Remove the part

Video Grounded ✓

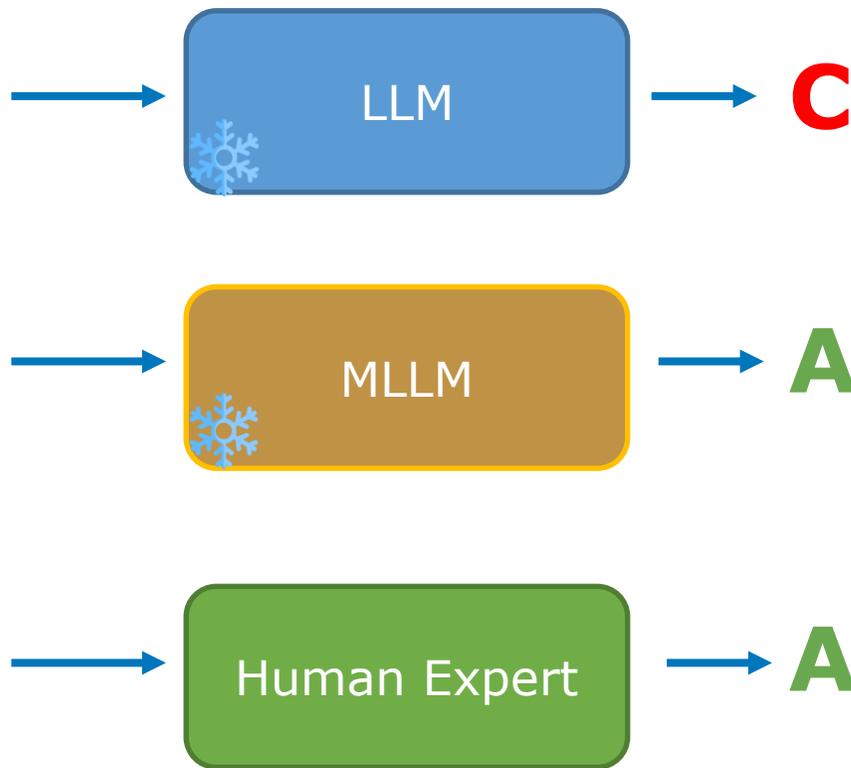
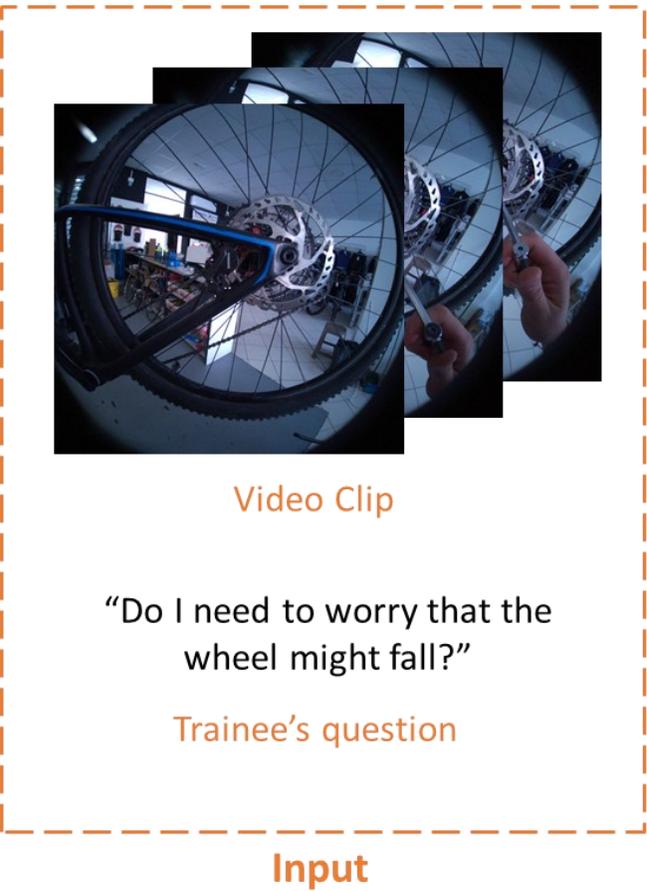


Q: Why shouldn't I apply too much pressure when unscrewing?

A: To avoid stripping the screw
 B: To avoid breaking the furniture
 C: To avoid using the hammer
 D: To avoid touching the camera
 E: To avoid making a mess

Not Grounded ✗





- A) "No, not at this moment. Now, hold it like that. "
- B) "Maybe we should stop and secure everything again to be absolutely sure."
- c) "No, but it’s better to use additional supports or have someone assist you just in case."
- D) "No, just let go and see if it stays in place."

5 What is the purpose of the wooden pins? *

- To attach the seat to the chair
- To hold the sticks together
- To tighten the screws
- To loosen the bolts



	Model	Bike Workshop	Bakery	Assembly	Kitchen	Avg.
Language Only	Llama 3.1 Instruct 8B	07.63	08.62	07.45	10.96	08.67
	Llama 3.1 Instruct 70B	27.57	22.54	25.19	31.30	26.65
	Llama 3.3 Instruct Turbo	27.14	18.61	24.67	30.42	25.21
	Qwen 2.5 Instruct 72B	20.27	15.28	19.01	21.54	19.02
	DeepSeek-R1 Turbo	24.22	21.94	21.73	26.15	23.51
Video-Language	MiniGPT4-video	06.62	07.09	08.26	15.74	10.68
	LLaVa Video	27.01	27.16	26.12	32.09	28.55
	LLaVa-OneVision	32.03	33.13	30.88	<u>35.77</u>	33.06
	Qwen 2.5-VL	<u>29.99</u>	<u>28.59</u>	<u>27.47</u>	35.87	<u>31.11</u>
	Sample Human Baseline	87.50	90.91	100	81.82	89.65

Language Only

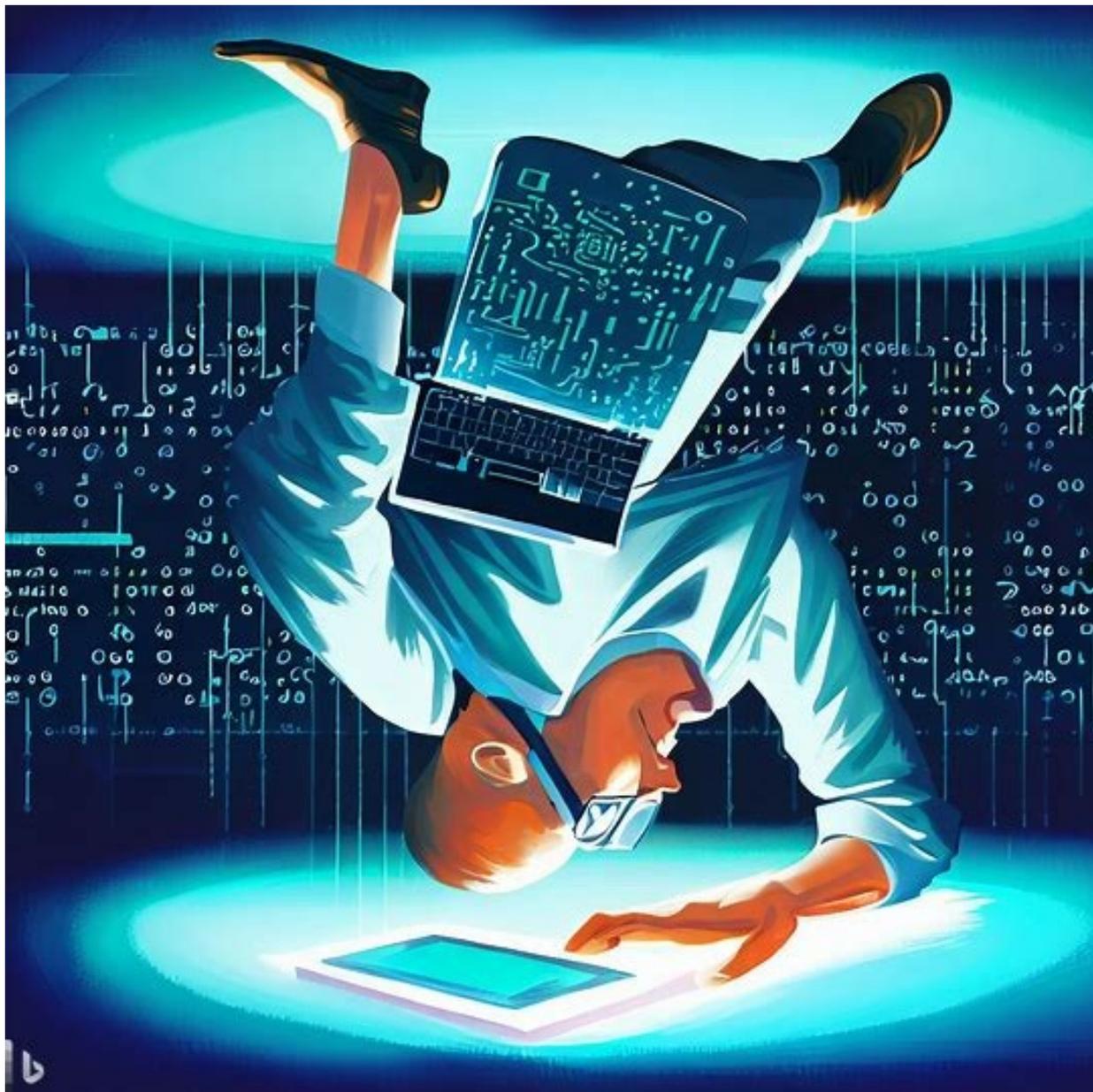
Video Language

Model	Bike Workshop	Bakery	Assembly	Kitchen	Avg.
Llama 3.1 Instruct 8B	07.63	08.62	07.45	10.96	08.67
Llama 3.1 Instruct 70B	27.57	22.54	25.19	31.30	26.65
Llama 3.3 Instruct Turbo	27.14	18.61	24.67	30.42	25.21
Qwen 2.5 Instruct 72B	20.27	15.28	19.01	21.54	19.02
DeepSeek-R1 Turbo	24.22	21.94	21.73	26.15	23.51
MiniGPT4-video	06.62	07.09	08.26	15.74	10.68
LLaVa Video	27.01	27.16	26.12	32.09	28.55
LLaVa-OneVision	32.03	33.13	30.88	<u>35.77</u>	33.06
Qwen 2.5-VL	<u>29.99</u>	<u>28.59</u>	<u>27.47</u>	35.87	<u>31.11</u>
Sample Human Baseline	87.50	90.91	100	81.82	89.65

What's Next?



An Outlook into the Future



Imagine the Future

Write Stories in Different Scenarios

Extract Important Tasks from the Stories

Go in-depth with Tasks and Datasets

A lot of data!



Rather than being extensive, we considered **seminal** and **state-of-the-art** works

An Outlook into the Future of Egocentric Vision

Chiara Plizzari* · Gabriele Goletto* · Antonino Furnari* · Siddhant Bansal* · Francesco Ragusa* · Giovanni Maria Farinella† · Dima Damen† · Tatiana Tommasi†



Politecnico di Torino



University of BRISTOL



Università di Catania

Received: date / Accepted: date

Abstract *What will the future be? We wonder!*

In this survey, we explore the gap between current research in egocentric vision and the ever-anticipated future, where wearable computing, with outward facing cameras and digital overlays, is expected to be integrated in our every day lives. To understand this gap, the article starts by envisaging the future through character-based stories, showcasing through examples the limitations of current technology. We then provide a mapping between this future and previously defined research tasks. For each task, we survey its seminal works, current state-of-the-art methodologies and available datasets, then reflect on shortcomings that limit its applicability to future research. Note that this survey focuses on software models for egocentric vision, independent of any specific hardware. The paper concludes with recommendations for areas of immediate explorations so as to unlock our path to the future always-on, personalised and life-enhancing egocentric vision.

Keywords Egocentric Vision, Future, Survey, Localisation, Scene Understanding, Recognition, Anticipation, Gaze Prediction, Social Understanding, Body Pose Estimation, Hand and Hand-Object Interaction, Person Identification, Summarisation, Dialogue, Privacy

Contents

1	Introduction	1
2	Imagining the Future	2

*: Equal Contribution/First Author

†: Equal Senior Author

C. Plizzari, G. Goletto and T. Tommasi, Politecnico di Torino, Italy · A. Furnari, F. Ragusa and G. M. Farinella, University of Catania, Italy · S. Bansal and D. Damen, University of Bristol, UK. E-mail: Tatiana.Tommasi@polito.it

2.1	EGO-Home	2
2.2	EGO-Worker	4
2.3	EGO-Tourist	5
2.4	EGO-Police	6
2.5	EGO-Designer	7
3	From Narratives to Research Tasks	8
4	Research Tasks and Capabilities	10
4.1	Localisation	10
4.2	3D Scene Understanding	14
4.3	Recognition	16
4.4	Anticipation	21
4.5	Gaze Understanding and Prediction	23
4.6	Social Behaviour Understanding	24
4.7	Full-body Pose Estimation	28
4.8	Hand and Hand-Object Interactions	30
4.9	Person Identification	36
4.10	Summarisation	38
4.11	Dialogue	40
4.12	Privacy	43
4.13	Beyond individual tasks	45
5	General Datasets	45
6	Conclusion	49

1 Introduction

Designing and building tools able to support human activities, improve quality of life, and enhance individuals' abilities to achieve their goals is the ever-lasting aspiration of our species. Among all inventions, digital computing has already had a revolutionary effect on human history. Of particular note is mobile technology, currently integrated in our lives through hand-held devices, i.e. *mobile smart phones*. These are nowadays the de facto for outdoor navigation, capturing static and moving footage of our everyday and connecting us to both familiar and novel connections and experiences.

However, humans have been dreaming about the next-version of such mobile technology — wearable computing, for a considerable amount of time. Imaginations

OpenReview.net

An Outlook into the Future of Egocentric Vision



Chiara Plizzari, Gabriele Goletto, Antonino Furnari, Siddhant Bansal, Francesco Ragusa, Giovanni Maria Farinella, Dima Damen, Tatiana Tommasi

14 Aug 2023 OpenReview Archive Direct Upload Readers: Everyone Show Revisions

Abstract: What will the future be? We wonder!

In this survey, we explore the gap between current research in egocentric vision and the ever-anticipated future, where wearable computing, with outward facing cameras and digital overlays, is expected to be integrated in our every day lives. To understand this gap, the article starts by envisaging the future through character-based stories, showcasing through examples the limitations of current technology. We then provide a mapping between this future and previously defined research tasks. For each task, we survey its seminal works, current state-of-the-art methodologies and available datasets, then reflect on shortcomings that limit its applicability to future research. Note that this survey focuses on software models for egocentric vision, independent of any specific hardware. The paper concludes with recommendations for areas of immediate explorations so as to unlock our path to the future always-on, personalised and life-enhancing egocentric vision.

Add Comment

Reply Type: Author: Visible To: Hidden From:

6 Replies

[+] Related work on modeling social interactions, especially multimodal dialogue agents



Jaewoo Ahn

18 Aug 2023 OpenReview Archive Paper22166 Comment Readers: Everyone Show Revisions

Comment:

I've been reading your fascinating work and wanted to contribute a suggestion based on my recent research in multimodal dialogue agents.

In our recent paper [1], we explored the benefits of a multimodal approach to dialogue personalization. Our study showed that incorporating both text and images in defining a persona greatly enriched the dialogue agent's understanding and personalization capabilities. Specifically, the image modality (i.e., egocentric vision) allowed the dialogue agents to access and better understand their personal characteristics and experiences based on their "episodic memory".

Drawing from this, I propose that there is a strong case to be made for the integration of egocentric vision into the domain of personalized dialogue agent responses. Egocentric vision, being intrinsically tied to personal perspective and experience, can serve as a valuable addition to a persona's episodic memory. This integration can enable chatbots to generate more contextually aware, and personalized responses based on the visual experiences of a user. The fusion of such vision-based episodic memory with textual modalities can be also a promising avenue for future research in personalized dialogue agents.

[1] Ahn et al. MPMCHAT: Towards Multimodal Persona-Grounded Conversation, ACL 2023 (<https://aclanthology.org/2023.acl-long.189/>)

Add Comment

[+] Related work on egocentric full-body pose estimation



Jiayi Jiang

17 Aug 2023 (modified: 17 Aug 2023) OpenReview Archive Paper22166 Comment Readers: Everyone Show Revisions

Comment:

Thanks for the nice paper, that's awesome!

I would really appreciate if our work (AvatarPoser [1] and EgoPoser [2]) on the topic of egocentric full-body pose estimation can also be presented in this review paper.

EGO-HOME

Sam is finally home after a long day. EgoAI kept track of Sam's food intake and a tomato soup sounds like the best complementary nutrition

- EgoAI localises Marco and provides route instructions to reach his workstation for the day
- This way the tomato will cook evenly
- A 3D projection of Remy helping with cooking
- Audible 3D projection
- Toaster reminder
- Sam is impressed by how fun it is to cook with his 3D friend
- EgoAI recommends some more spice
- Waves hitting the shore look and sound natural
- Transferred to a beach he visited last summer
- After dinner, Sam enjoys a group card game with his friends, who are connected through their own EgoAI
- EgoAI proposes a short clip from his day, but Sam decides not to share it

While getting ready for bed, Sam feels an itch on the wrist that has annoyed him the whole day. EgoAI stores a picture of the injury and sends it to Sam's doctor for advice

EGO-WORKER

EgoAI verifies if Marco is properly wearing the Personal Protection Equipment (PPE)

- EgoAI localises Marco and provides route instructions to reach his workstation for the day
- Where should I go today in the factory?
- In the past, EgoAI guided Marco to the closest fire extinguisher during a fire
- EgoAI passes a message from the manager about today's goal: testing a set of electric boards
- Since the measuring device is a new brand, EgoAI guides Marco through the basic functionality and tools
- EgoAI detects a risk and turns off the IoT electrical socket while promptly alerting Marco
- For the rest of the day, EgoAI validates Marco's work making sure all the procedures are properly and safely completed
- By the end of the day, EgoAI checks Marco's feedback for improving future sessions

EGO-TOURIST

EgoAI prepares Claire a personalised and exciting one-day itinerary in Turin

- EgoAI suggests an half-day visit to the Egyptian museum
- Claire feels transported to ancient Egypt
- Claire asks Cleopatra for a good place for a pizza
- Claire observes virtual elements being added to the scene, which bring the artwork to life
- Cleopatra leads Claire through the artworks and proposes her the most suited path
- Cleopatra discovers a fantastic pizza place for lunch while also enlightening Claire about the history behind various Italian monuments
- EgoAI has reserved an afternoon at the thermal baths. The next bus is scheduled to arrive in 20 minutes
- EgoAI offers a egocentric view from the chef who prepared her that delicacy
- EgoAI suggests Claire a proper Italian coffee at a nearby cafe, sided by a slice of bunet, Turin-based dessert
- EgoAI actively saved snapshots and videos of the day
- EgoAI retrieves the closest souvenir shop based on Claire's taste and budget

EGO-POLICE

EgoAI is constantly pinpointing Judy's position and would send an alert to the headquarters if she encounters unusual events or dangerous situations

- EgoAI helps Judy navigate the shortest safe path to target places
- One of the fellow officers shared via EgoAI a clip from a surveillance camera one block east: the suspect was moving in Judy's direction
- EgoAI detected and re-identified the man before he passed Judy
- EgoAI accesses the lost-and-found database of the airport
- EgoAI has both thermal and multi-spectral sensors
- Thanks to its sensors, EgoAI calculates a low risk for explosive content
- Judy was able to swiftly arrest him
- Judy also appreciated the help of EgoAI when she had to manage an abandoned backpack
- EgoAI connects Judy with the bomb squad and live-shares the observed scene
- EgoAI projects a clear red circle around the backpack with the minimal stand-off distance
- Thanks to EgoAI, all the relevant events are saved and transformed into a document with related images and video recordings
- EgoAI guides Judy with exact instructions to grasp the backpack and open it
- The sensitive information is properly identified and secured under admin rights to protect citizens' privacy

EGO-DESIGNER

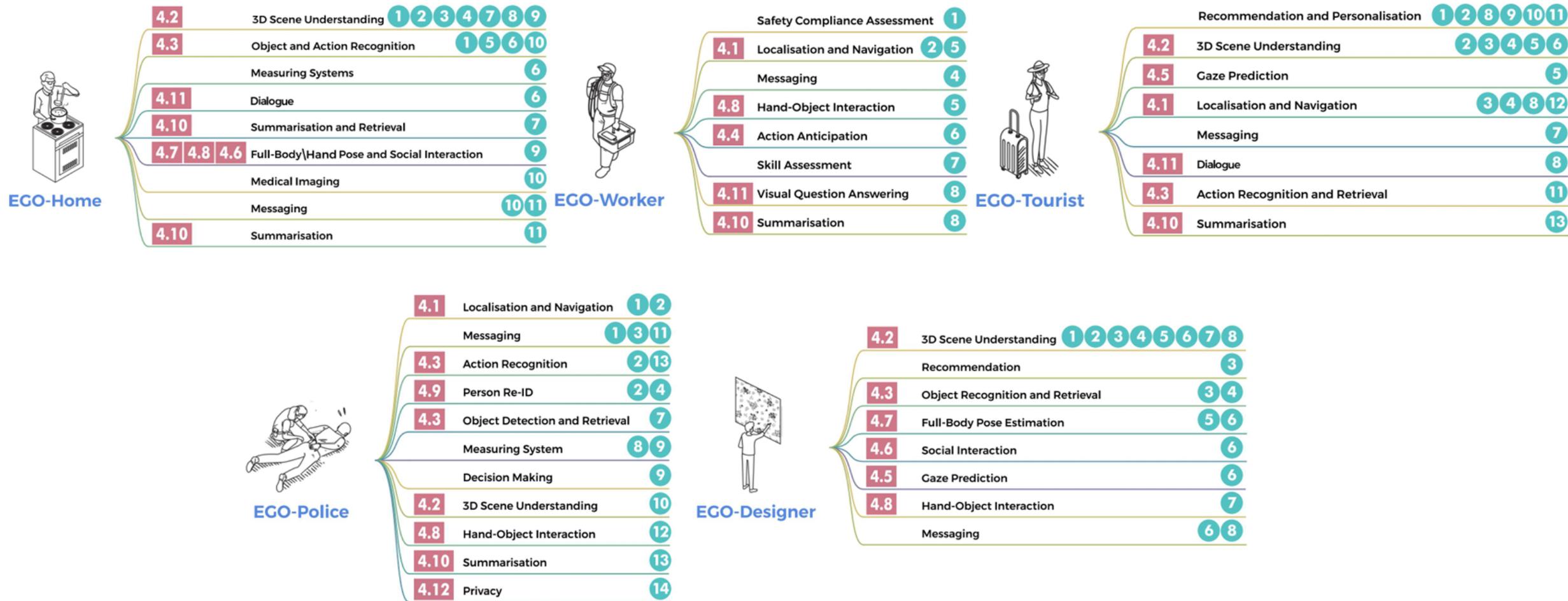
EgoAI helps Stanley (the scenographer) re-design the surrounding environment. The real scene represents the hall of a villa in New York, but it is almost empty

- EgoAI adds a luxurious wallpaper with floral patterns
- EgoAI also suggests adding velvet couches on the right and a carved wooden table on the left
- EgoAI has access to the database of the equipment warehouse; Stanley can search for the available pieces of furniture
- EgoAI also allows Stanley to visualise how the actors should move in the space considering that there will be musicians in the middle of the room
- EgoAI shares the scene with the actors. Through their own EgoAI, they are immersed inside the changing and moving 3D computer-generated environment
- EgoAI assists make-up artists with advanced 3D modelling techniques to project guidelines on the actor's face while applying make-up
- EgoAI also assists the director. He is able to preview the planned scene and light effects in real-time while shooting the scene



12 Egocentric Vision Research Tasks

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



*perspective and provides ego-based assistance. We associate story **P**arts with research tasks (marked by **section number**) and later revisit the link between these*

Table 1 General Egocentric Datasets - Collection Characteristics. †: For EGTEA, Audio was collected but not made public.
*: For Ego4D, apart from RGB, the other modalities are present for subsets of the data.

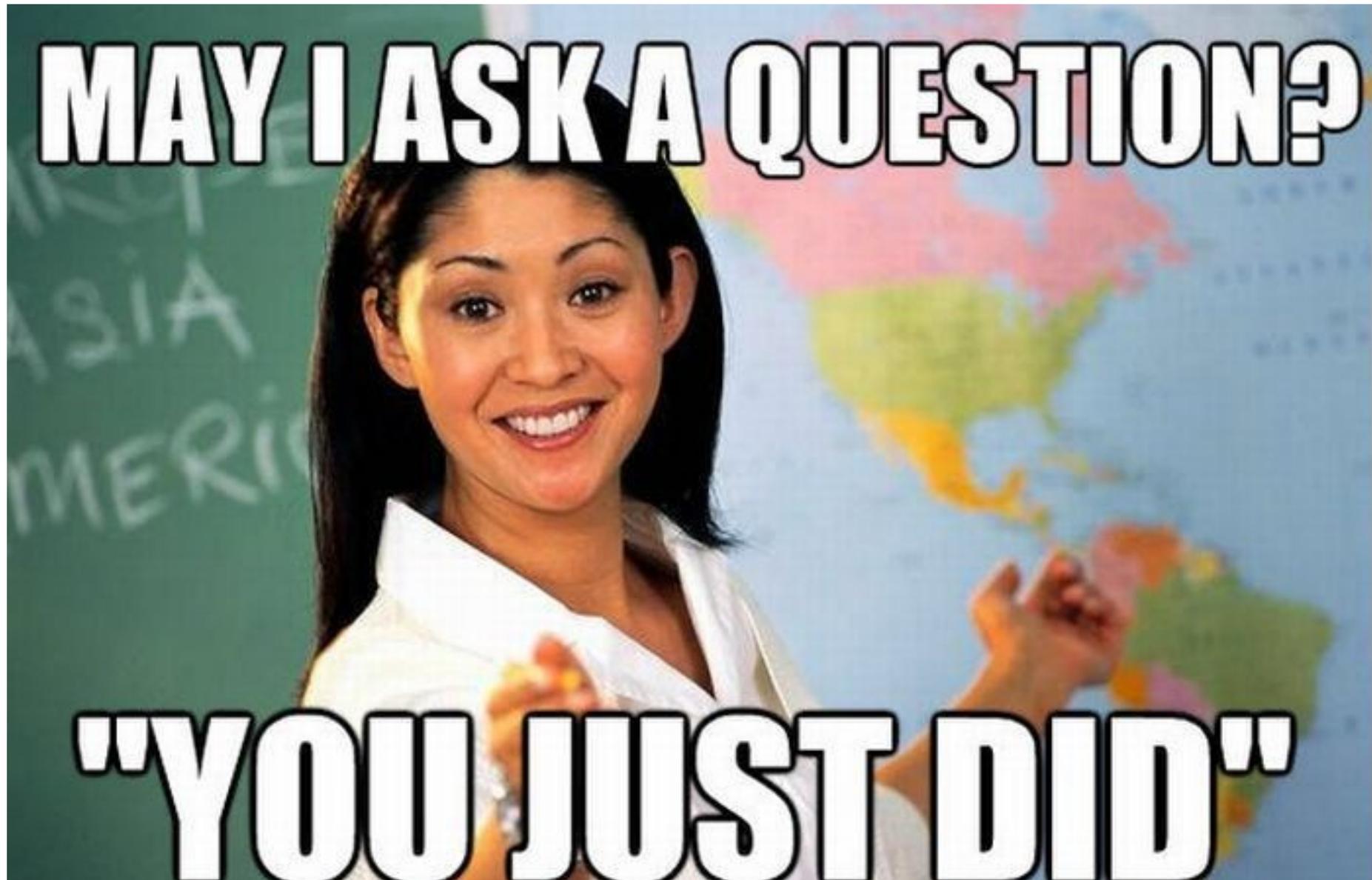
Dataset	Settings	Signals	Hours	Sequences	AVG. video duration	Participants
MECCANO (Ragusa et al 2023b)	Industrial	RGB, depth, gaze	6.9	20	20.79 min	20
ADL (Pirsiavash and Ramanan 2012)	Daily activities	RGB	10.0	20	30.00 min	20
HOI4D (Liu et al 2022c)	Table-Top	RGB, depth	22.2	4000	0.33 min	9
EGTEA Gaze+† (Li et al 2021a)	Kitchen	RGB, gaze	27.9	86	19.53 min	32
UTE (Lee et al 2012)	Daily Activities	RGB	37.0	10	222.00 min	4
EGO-CH (Ragusa et al 2020a)	Cultural Sites	RGB	37.1	180	12.37 min	70
FPSI (Fathi et al 2012a)	Recreational Site	RGB	42.0	8	315.00 min	8
KrishnaCam (Singh et al 2016a)	Daily Routine	RGB, GPS, acc	69.9	460	9.13 min	1
EPIC-KITCHENS-100 (Damen et al 2022)	Kitchens	RGB, audio	100.0	700	8.57 min	37
Assembly101 (Sener et al 2022)	Industrial	RGB, multi-view	167.0	1425	7.10 min	53
Ego4D* (Grauman et al 2022)	Multi Domain	RGB, Audio, 3D, gaze, IMU, multi	3670.0	9650	24.11 min	931

Table 2 General Egocentric Datasets - Current set of annotations. *: For Ego4D, apart from narrations, the remaining annotations are only available for subsets of the dataset depending on the benchmark

Dataset	Annotations
MECCANO (Ragusa et al 2023b)	Temporal action segments, hand & object bounding boxes, hand-object interactions, next-active object
ADL (Pirsiavash and Ramanan 2012)	Temporal action segments, objects bounding boxes, hand-object interactions
HOI4D (Liu et al 2022c)	Temporal action segments, 3D hand poses and object poses, panoptic and motion segmentation, object meshes, scene point clouds
EGTEA Gaze+ (Li et al 2021a)	Temporal action segments, hand masks, gaze
UTE (Lee et al 2012)	Text descriptions, object segmentations
EGO-CH (Ragusa et al 2020a)	Temporal locations, object bounding boxes, surveys, object masks
FPSI (Fathi et al 2012a)	Temporal social interaction segments
KrishnaCam (Singh et al 2016a)	Motion classes, virtual webcams, popular locations
EPIC-KITCHENS-100 (Damen et al 2022)	Temporal action video segments, Temporal audio segments, narrations, hand and objects masks, hand-object interactions, camera poses
Assembly101 (Sener et al 2022)	Temporal action segments, 3D hand poses
Ego4D* (Grauman et al 2022)	Narrations, Temporal action segments, moment queries, speaker labels, diarisation, hand bounding boxes, time to contact, active objects bounding boxes, trajectories, next-active objects bounding boxes

Table 3 General Egocentric Datasets - Current set of tasks: **4.1** Localisation, **4.2** 3D Scene Understanding, **4.3** Recognition, **4.4** Anticipation, **4.5** Gaze Understanding and Prediction, **4.6** Social Behaviour Understanding, **4.7** Full-body Pose Estimation, **4.8** Hand and Hand-Object Interactions, **4.9** Person Identification, **4.10** Summarisation, **4.11** Dialogue, **4.12** Privacy.

Dataset	Task												
	4.1	4.2	4.3	4.4	4.5	4.6	4.7	4.8	4.9	4.10	4.11	4.12	
MECCANO (Ragusa et al 2023b)			✓	✓	✓			✓					
ADL (Pirsiavash and Ramanan 2012)			✓	✓						✓			
HOI4D (Liu et al 2022c)								✓					
EGTEA Gaze+ (Li et al 2021a)			✓	✓	✓			✓					
UTE (Lee et al 2012)								✓		✓			
EGO-CH (Ragusa et al 2020a)	✓												
FPSI (Fathi et al 2012a)							✓			✓		✓	
KrishnaCam (Singh et al 2016a)				✓									
EPIC-KITCHENS-100 (Damen et al 2022)		✓	✓	✓				✓			✓	✓	
Assembly101 (Sener et al 2022)			✓					✓					
Ego4D (Grauman et al 2022)			✓	✓	✓	✓		✓		✓	✓		





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

- 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