



**VISIGRAPP 2024**

19<sup>th</sup> International Joint Conference on Computer Vision, Imaging  
and Computer Graphics Theory and Applications

Rome, Italy 27 - 29 February, 2024

GRAPP HUCAPP IVAPP VISAPP



Università  
di Catania

**NEXT VISION**

Spin-off of the University of Catania



# First Person (Egocentric) Vision: History and Applications

## Francesco Ragusa

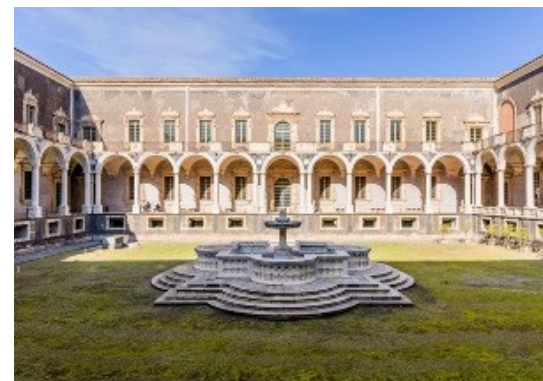
First Person Vision@Image Processing Laboratory - <http://iplab.dmi.unict.it/fpv>

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

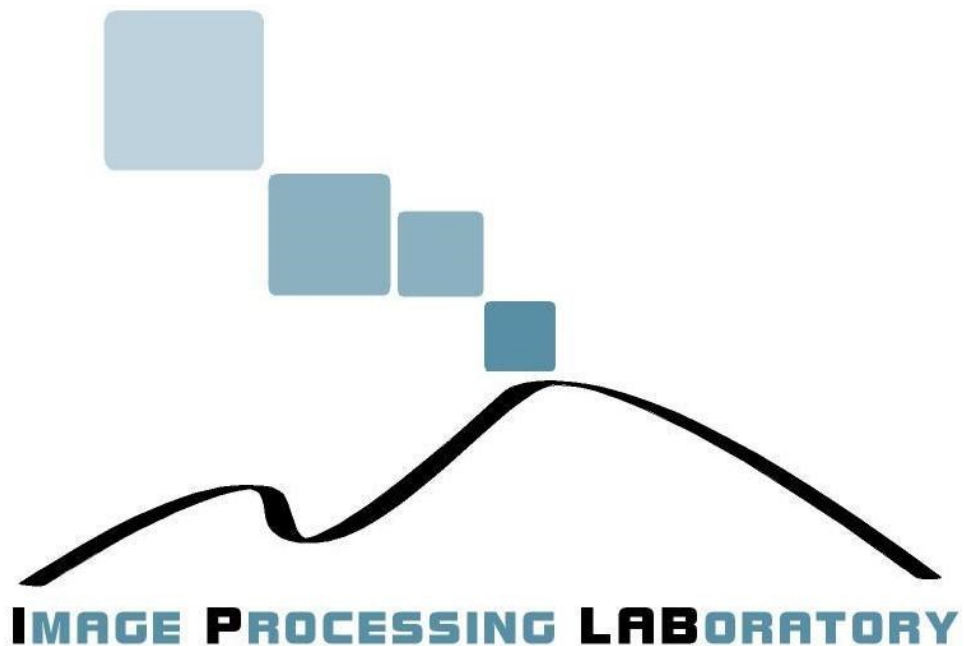
Department of Mathematics and Computer Science - University of Catania

[francesco.ragusa@unict.it](mailto:francesco.ragusa@unict.it) - <https://francescoragusa.github.io/>

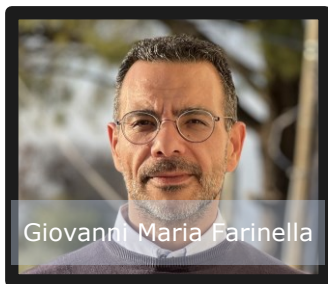




# Università di Catania



**FPV @ IPLAB Group**



Giovanni Maria Farinella



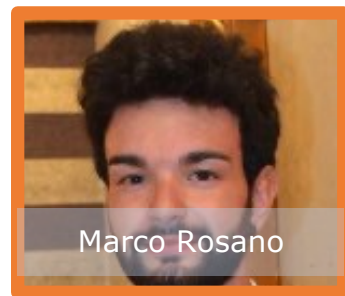
Antonino Furnari



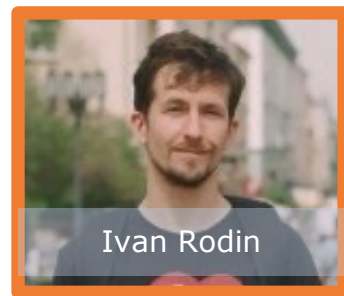
Francesco Ragusa



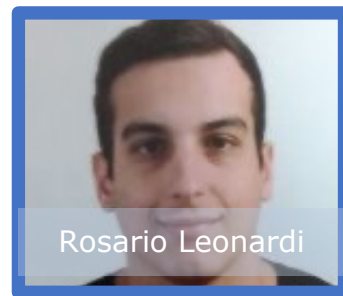
Daniele Di Mauro



Marco Rosano



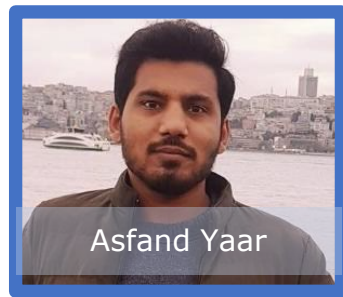
Ivan Rodin



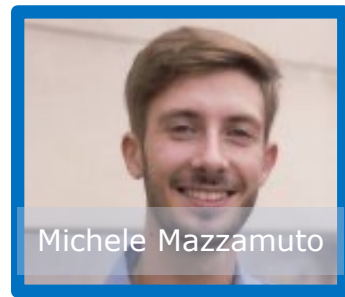
Rosario Leonardi



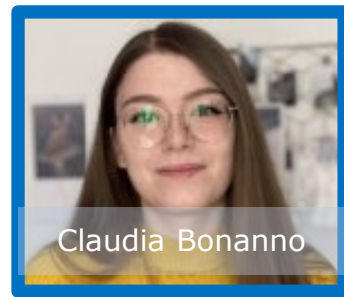
Camillo Quattrocchi



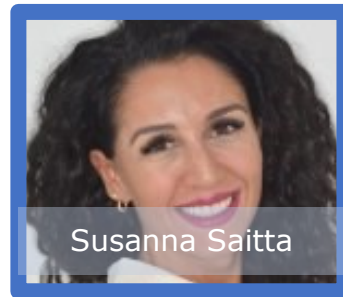
Asfand Yaar



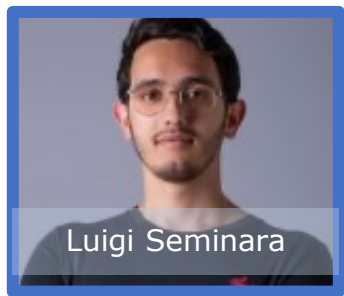
Michele Mazzamuto



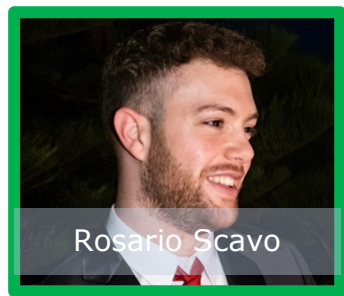
Claudia Bonanno



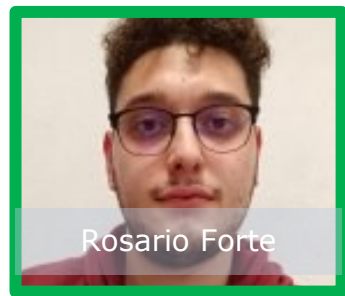
Susanna Saitta



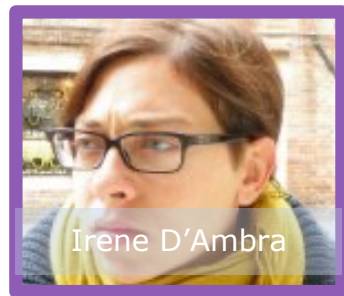
Luigi Seminara



Rosario Scavo



Rosario Forte



Irene D'Ambra

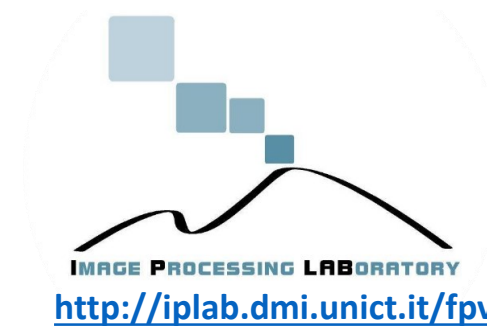


IMAGE PROCESSING LABORATORY

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

NEXT VISION

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

**16 Members**

1 Full Professor

1 Assistant Professor

1 Researcher

3 Post Docs

7 PhD Students

2 Master Students

1 Lab Assistant

The slides of this tutorial are available online at:  
<https://francescoragusa.github.io/visigrapp2024>



**1) Part I: History and motivations [09.00 - 10.30]**

- a) **Agenda of the tutorial;**
- b) **Definitions, motivations, history and research trends of First Person (egocentric) Vision;**
- c) **Seminal works in First Person (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;**

Coffee Break [10.30 – 10.45]

Keynote presentation: Gerhard Rigoll [10.45 – 12.00]

**1) Part II: Fundamental tasks for First Person Vision systems [12.00 – 13.00]**

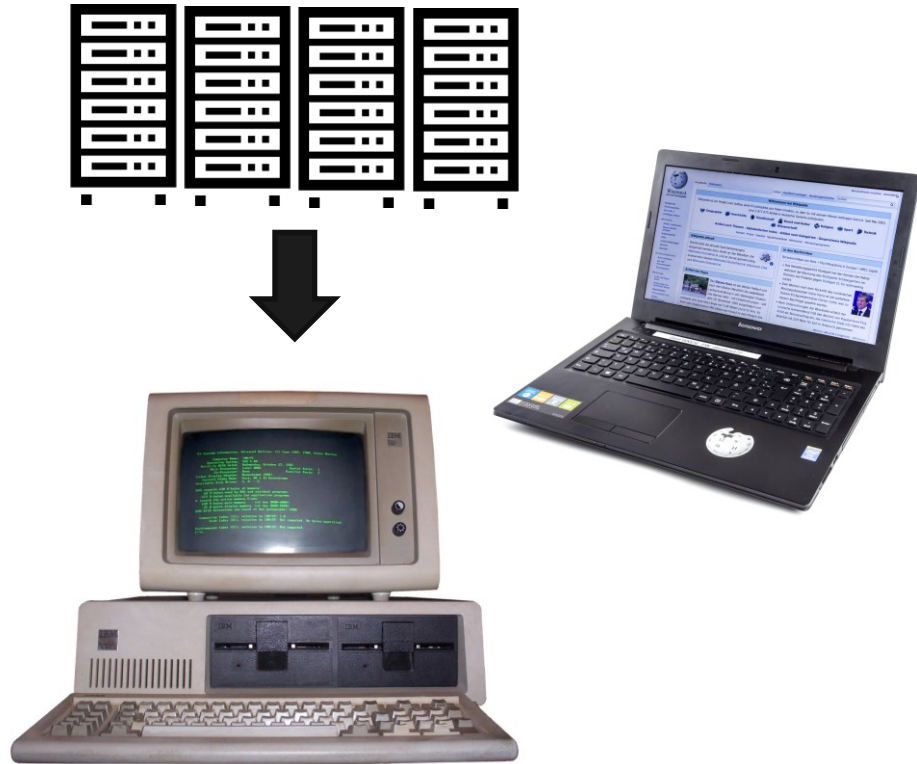
- a) **Localization;**
- b) **Hand/Object Detection;**
- c) **Action/Activity Recognition;**
- d) **Egocentric Human-Object Interaction;**
- e) **Anticipation;**
- f) **Industrial Applications;**
- g) **Conclusion.**

# Part I

## History and Motivations

After personal computers and smartphones, wearable devices are the third wave of computing

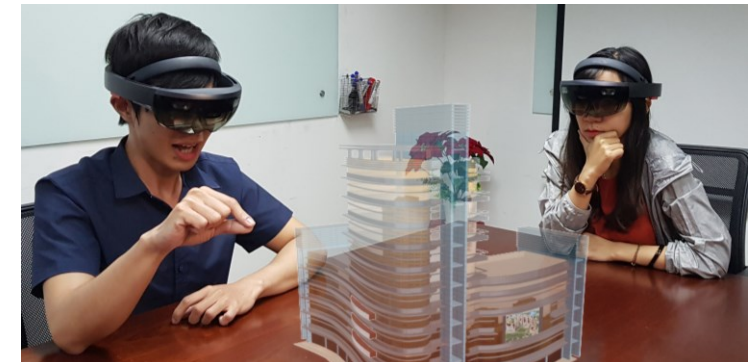
– [Marc Pollefeys](#), Lab Director, Microsoft Mixed Reality and AI Zurich



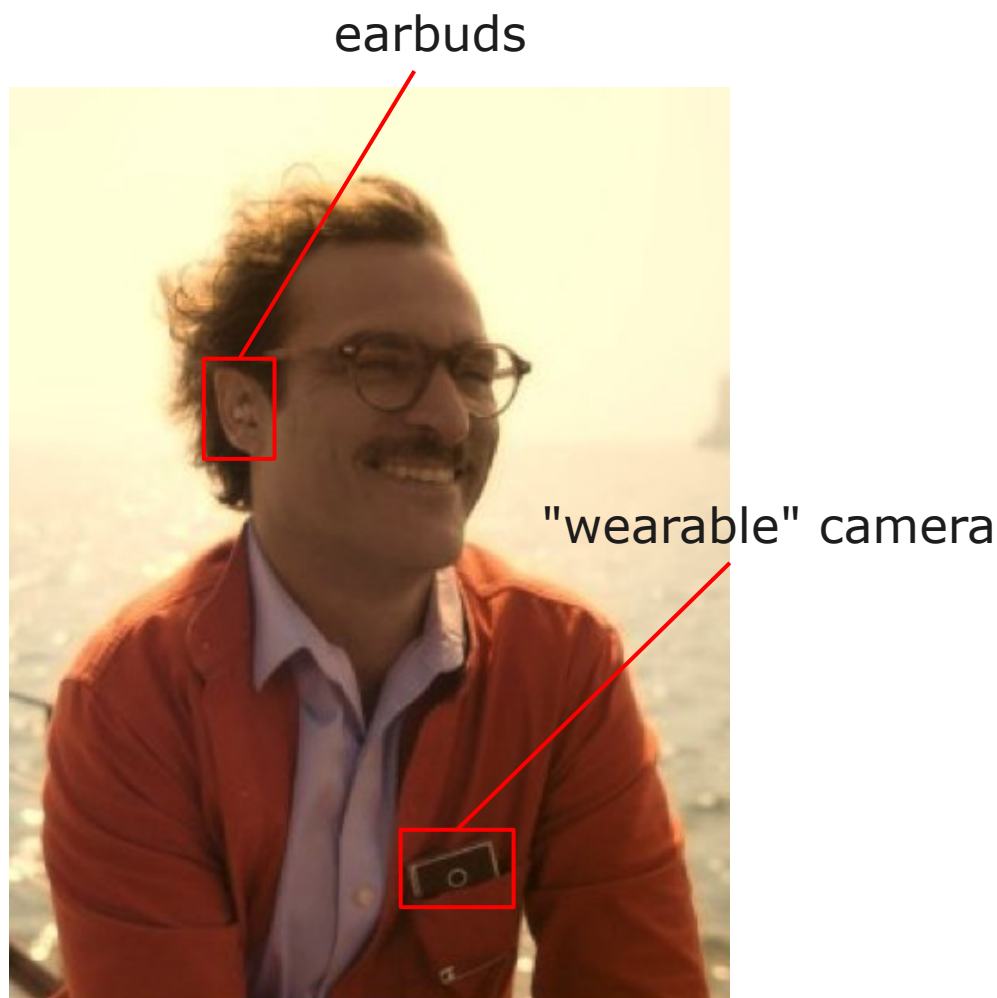
**Personal Computers:**  
computing for the mass, but not mobile and not context aware - dedicated access to computing



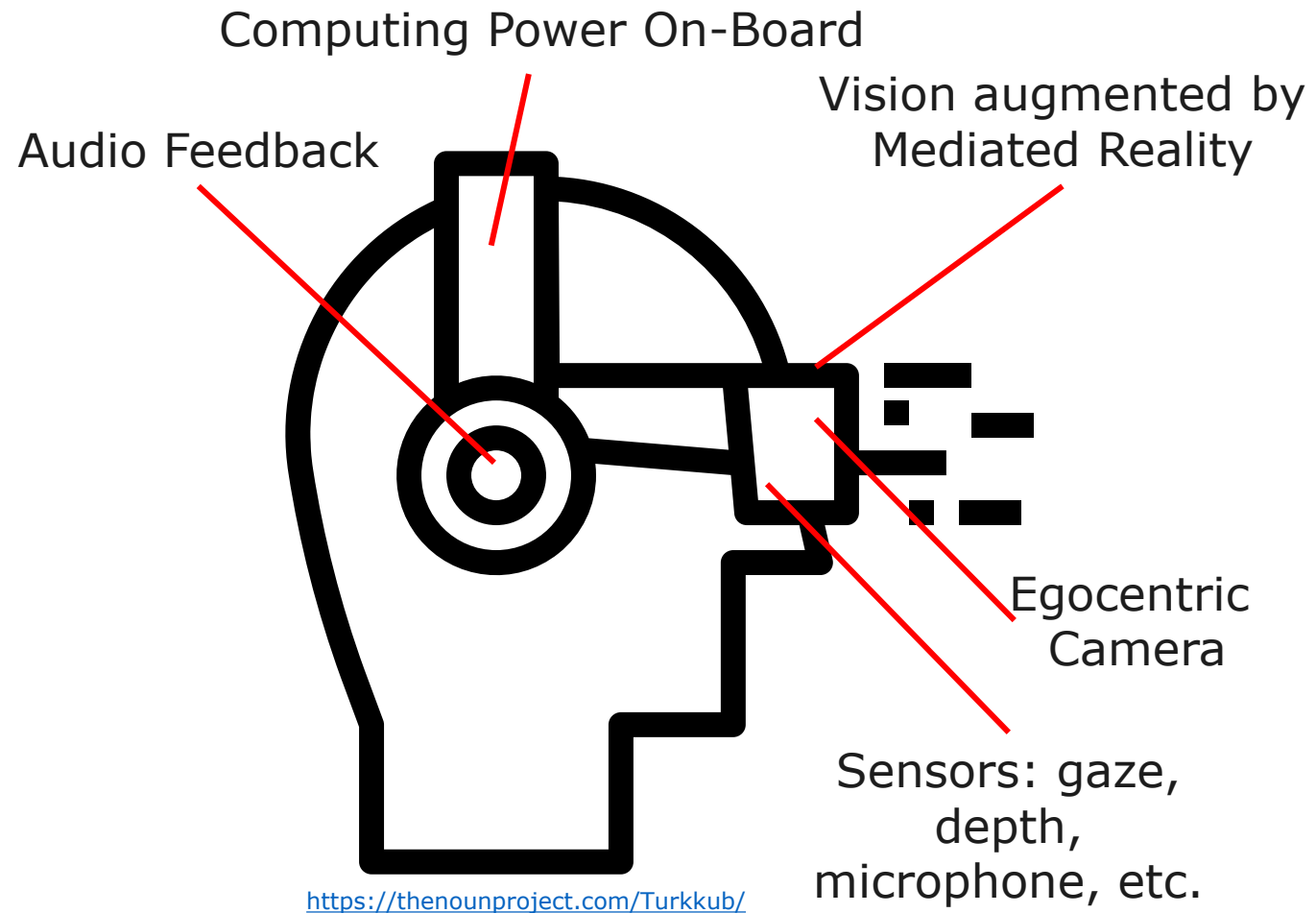
**Smartphones:** mobile computing is always accessible, but forces to switch between the digital and real world



**Eyeworn Devices:**  
computing everywhere with minimal switch between real and digital worlds

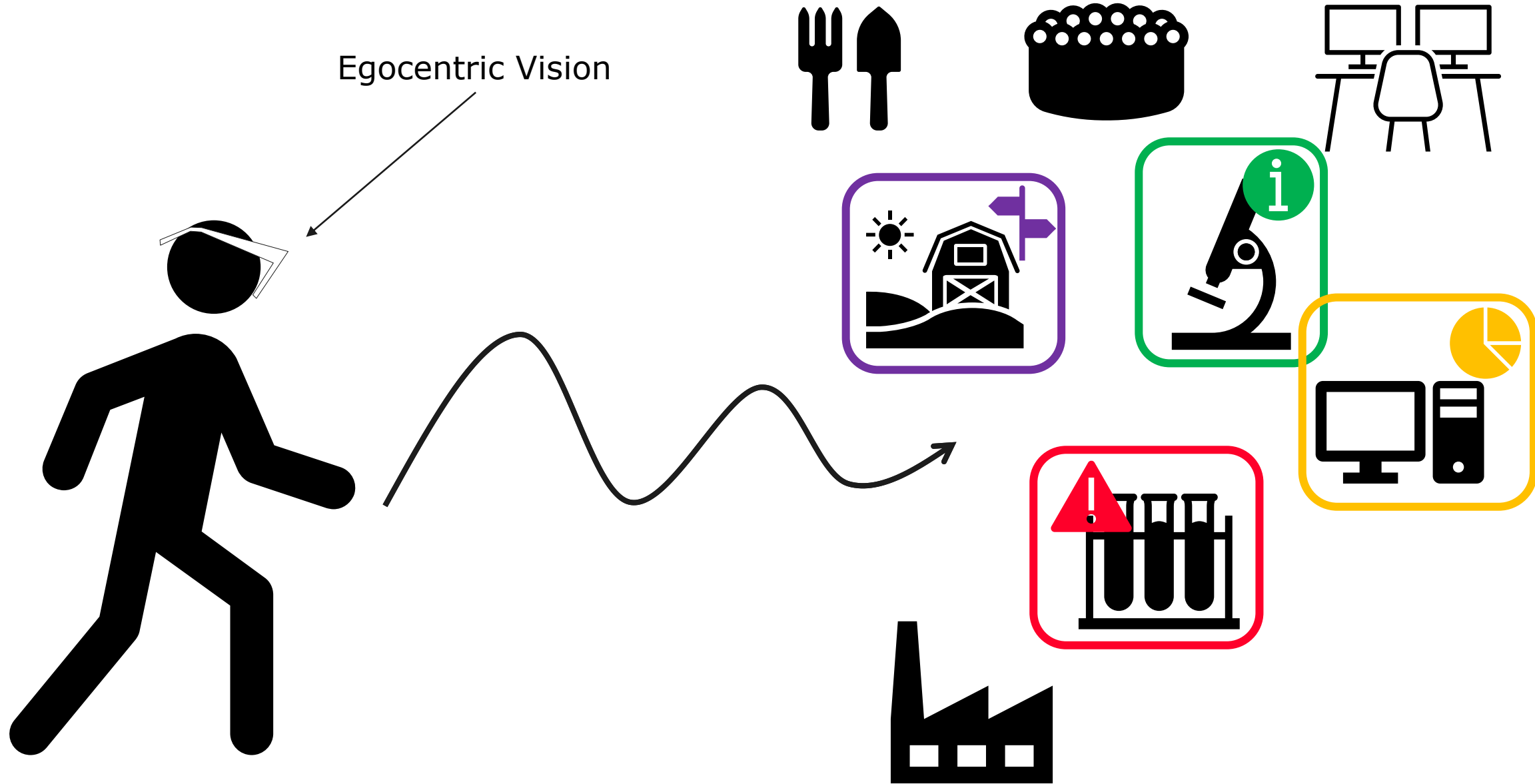


"her" 2013 movie



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!**



First Person Camera



Third Person Camera

## Wearable Camera



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

## Fixed Camera



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

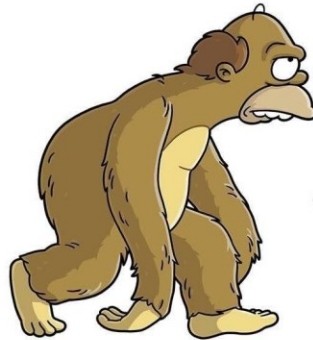
# When Did it All Begin?



MONKIUS EATALOTIS



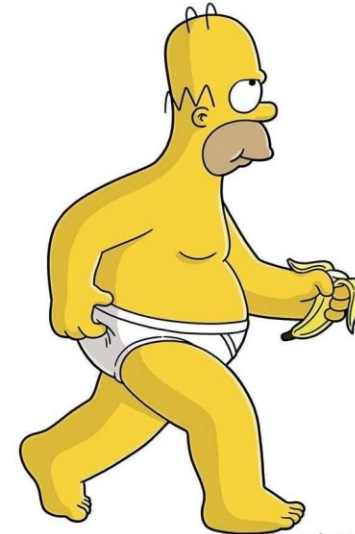
CHIMPUS IMBECILUS



APEIS STUPIDIUS



NEANDERSLOB

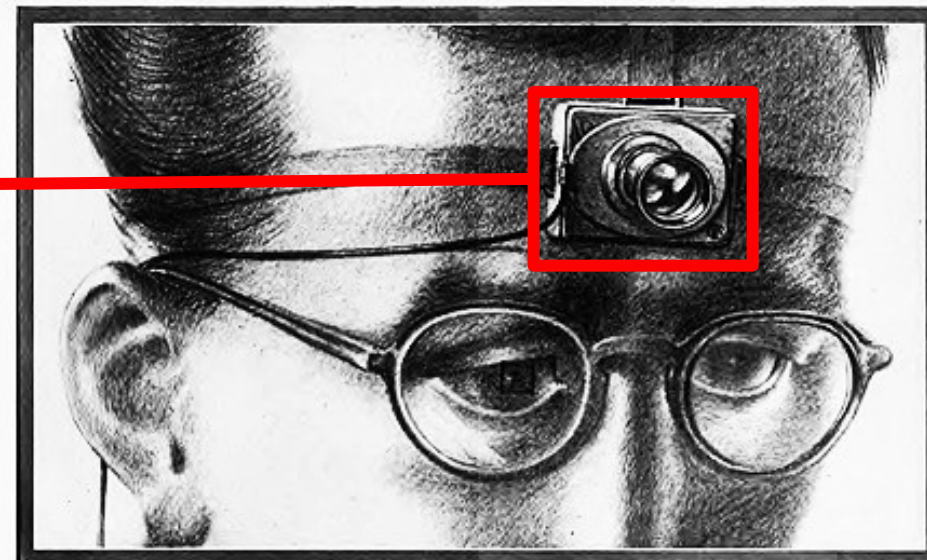
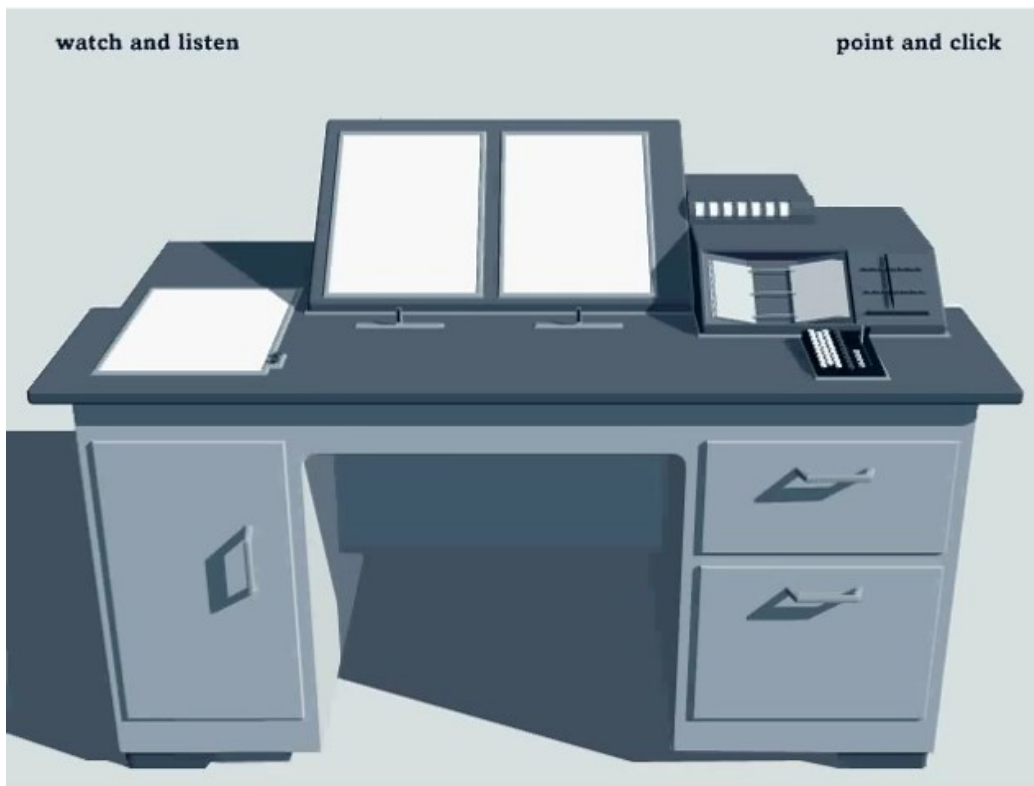


HOMERSAPIEN

HOMERSAPIEN

MAT 6000-10

“Certainly, progress in photography is not going to stop. [...] Let us project this trend ahead to a logical, if not inevitable, outcome. The camera hound of the future wears on his forehead a lump a little larger than a walnut.”



A SCIENTIST OF THE FUTURE RECORDS EXPERIMENTS WITH A TRY CAMERA, FITTED WITH UNIVERSAL-FOCUS LENS. THE SMALL SQUARE IN THE EPIGLASS AT THE LEFT SIGHTS THE OBJECT

## AS WE MAY THINK

A TOP U. S. SCIENTIST FORESEES A POSSIBLE FUTURE WORLD  
IN WHICH MAN-MADE MACHINES WILL START TO THINK

by VANNEVAR BUSH

DIRECTOR OF THE OFFICE OF SCIENTIFIC RESEARCH AND DEVELOPMENT  
Condensed from the *Atlantic Monthly*, July 1945

THIS has not been a scientists' war; it has been a war in which all have had a part. The scientists, buying their old professional competition in the demand of a common cause, have shared greatly and learned much. It has been exhilarating to work in effective partnership. What are the scientists to do next?

For the biologists, and particularly for the medical scientists, there can be little indication, for their war work has hardly required them to leave the old paths. Many indeed have been able to carry on their war research in their familiar peacetime laboratories. Their objectives remain much the same.

It is the physicists who have been thrown most violently off stride, who have had to devise new methods for their unanticipated assignments. They have done their part on the devices that made it possible to turn back the enemy. They have worked in combined effort with the physicists of our allies. They have felt within themselves the stir of achievement. They have been part of a great team. Now one asks where they will find objectives worthy of their love.

There is a growing mountain of research. But there is increased evidence that we are being bogged down today as specialization extends. The investigator is staggered by the findings and conclusions of thousands of other workers—conclusions which he cannot find time to grasp, much less to remember, as they appear. Yet specialization becomes increasingly necessary for prog-

ress, and the effort to bridge between disciplines is correspondingly superficial.

Professionally our methods of transmitting and reviewing the results of research are generations old and by now are totally inadequate for their purpose. If the aggregate time spent in writing scholarly works and in reading them could be evaluated, the ratio between these amounts of time might well be startling. Those who conscientiously attempt to keep abreast of current thought, even in restricted fields, by close and continuous reading might well shy away from an examination calculated to show how much of the previous month's efforts could be produced on call.

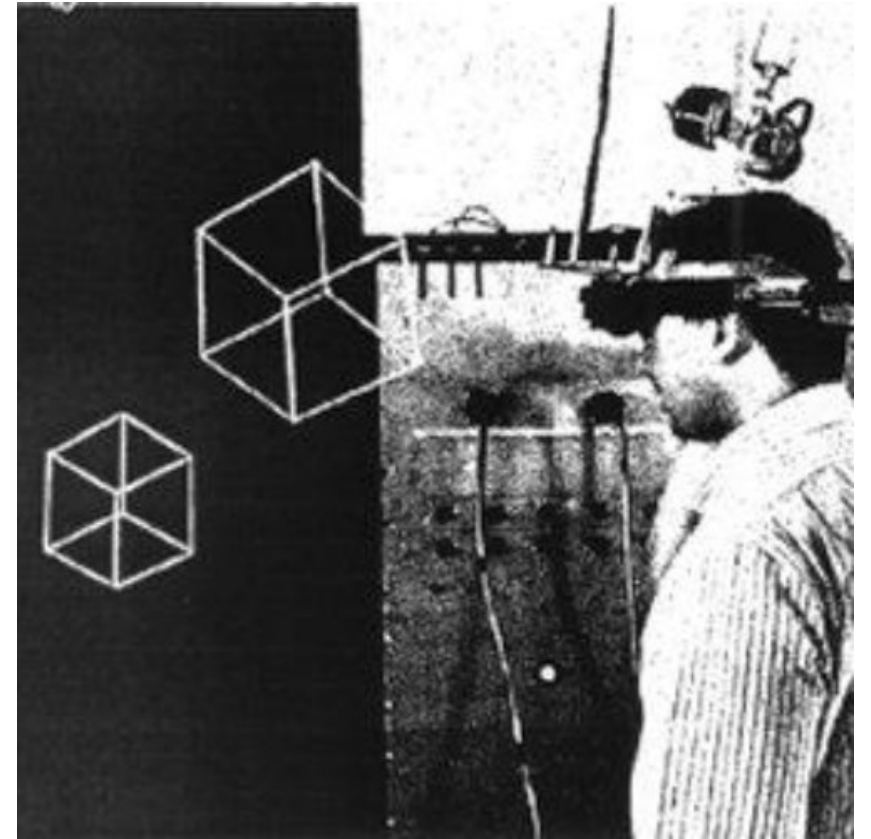
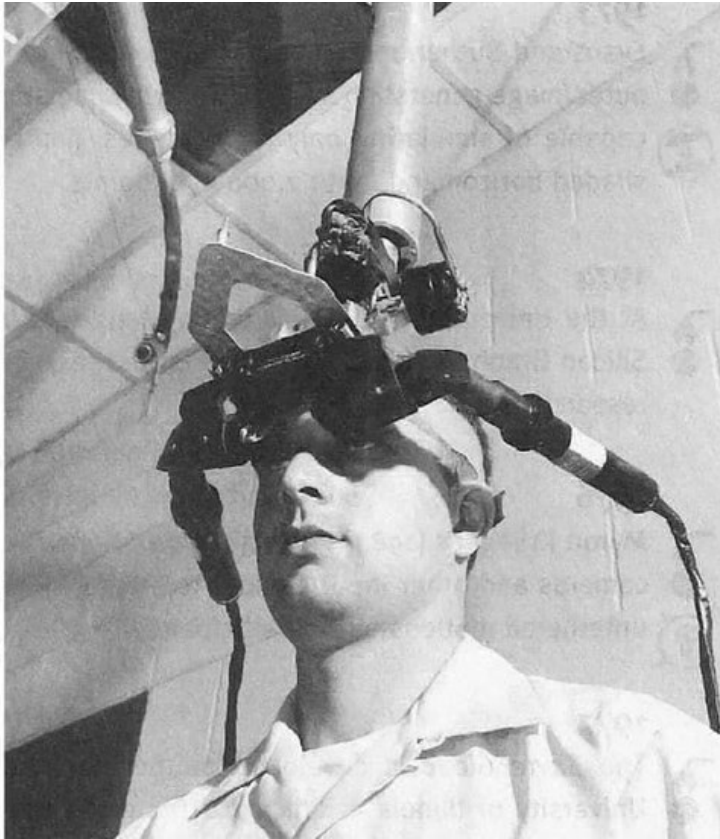
Mendel's concept of the laws of genetics was lost to the world for a generation because his publication did not reach the few who were capable of grasping and extending it. This sort of catastrophe is undoubtedly being repeated all about us as truly significant attainments become lost in the mass of the inconsequential.

Publication has been extended far beyond our present ability to make real use of the record. The summation of human experience is being expanded at a prodigious rate, and the means we use for threading through the consequent mass to the momentarily important item is the same as was used in the days of square-rigged ships.

But there are signs of a change as new and powerful instrumentalities come into use. Photochemicals capable of seeing things in a physical sense, advanced photography which can record what is seen or even what is not, thermionic tubes capable of controlling potent forces under the guidance of

# Head Mounted Display (1968)

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.



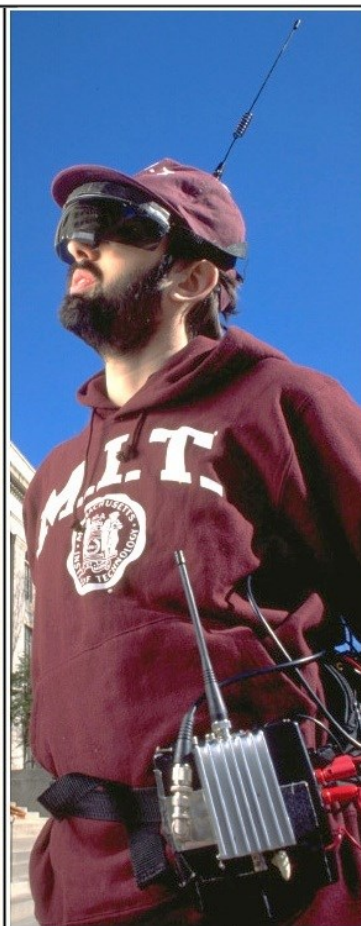
(a)  
**1980**



(b)  
**Mid 1980s**



(c)  
**Early 1990s**



(d)  
**Mid 1990s**



(e)  
**Late 1990s**

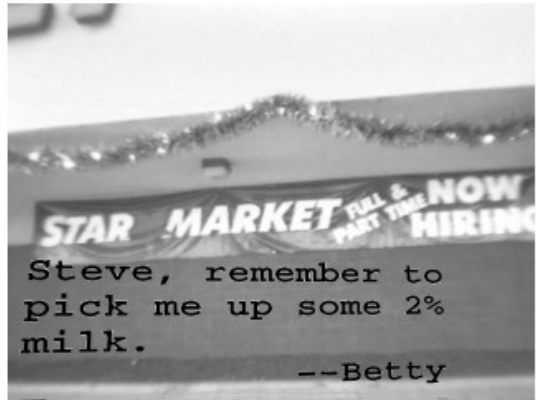
In the 80s and 90s Steve Mann (PhD in Media Arts and Sciences at MIT, 1997) invented a number of wearable computers featuring video capabilities, computing capabilities, and a wearable screen for feedback. **Steve Mann is often referred to as «the father of wearable computing»**



Visual Orbits



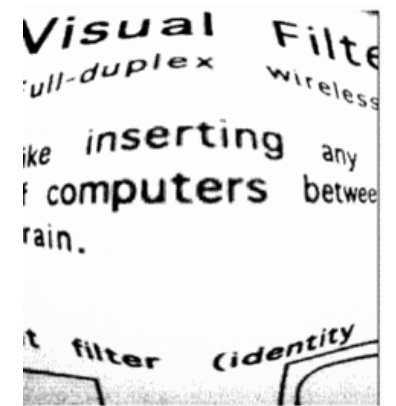
Meta-Vision



Spatialized Reminders



Spatialized Shopping List



Visual Filters

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.







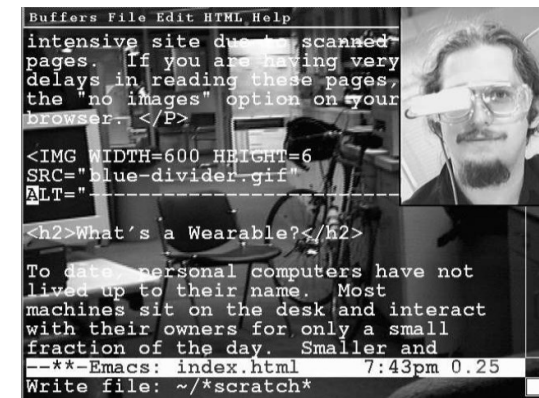
Clip from movie Terminator 2-Judgment day: <https://youtu.be/9MeaaCwBW28>

Ref: [https://www.redsharknews.com/vr\\_and\\_ar/item/3539-terminator-2-vision-the-augmented-reality-standard-for-25-years](https://www.redsharknews.com/vr_and_ar/item/3539-terminator-2-vision-the-augmented-reality-standard-for-25-years)

## 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

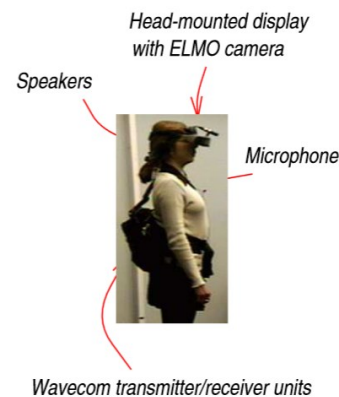
Thad Starner      Bernt Schiele      Alex Pentland  
Media Laboratory, Massachusetts Institute of Technology

(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

1999

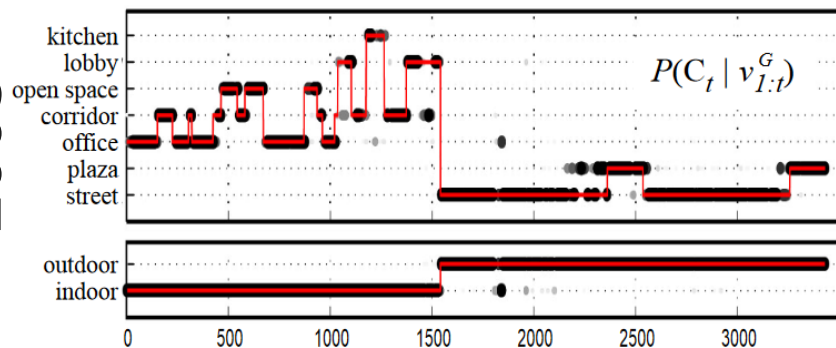
## Wearable Visual Robots

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



2000

2003



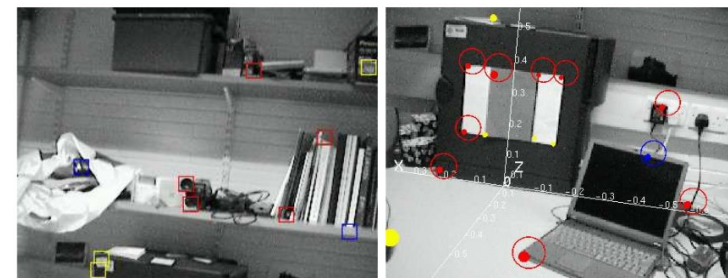
## Context-based vision system for place and object recognition

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

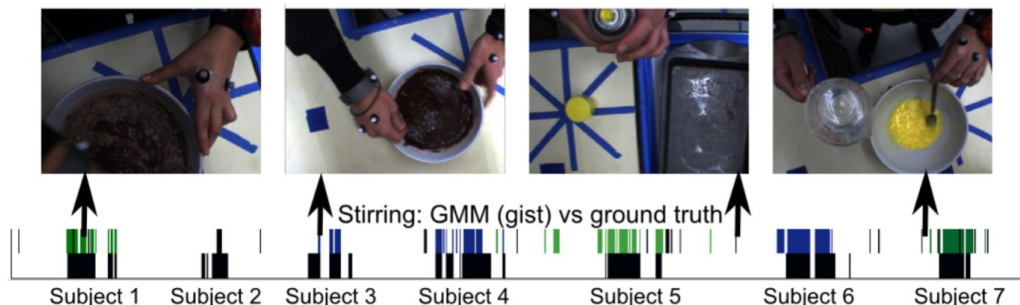
Department of Engineering Science  
University of Oxford

(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

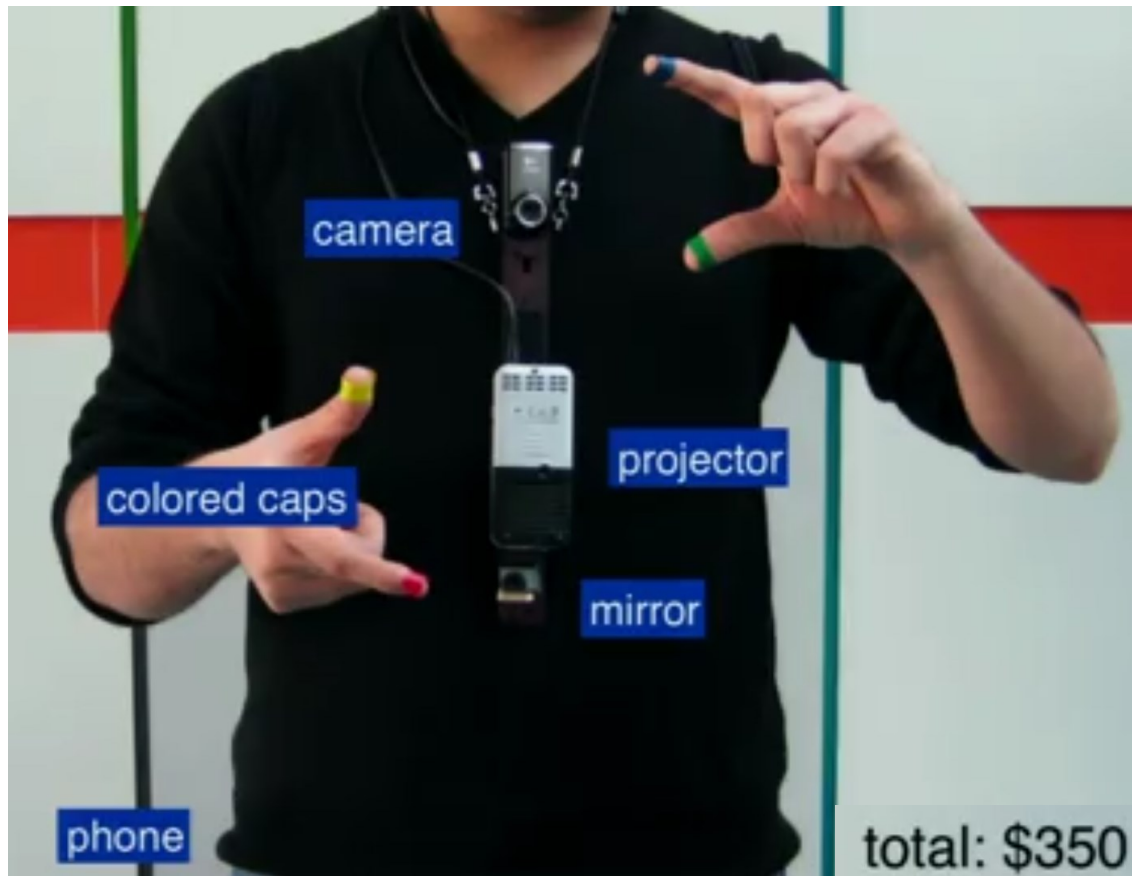
(handheld object recognition)



2010

## Neck worn camera with a projector and a gesture-based user interface.

«to give people access to information without requiring that the user changes any of their behavior»



Pattie Maes & Pranav Mistry (MIT) @ TED

[https://www.ted.com/talks/pattie\\_maes\\_demos\\_the\\_sixth\\_sense](https://www.ted.com/talks/pattie_maes_demos_the_sixth_sense)



## "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/>

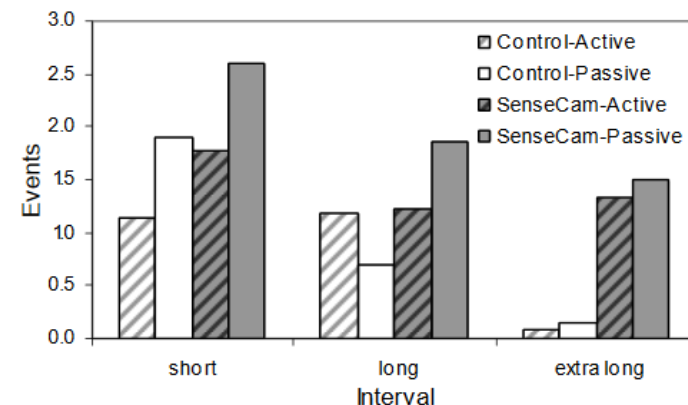
Bell, Gordon, and Jim Gemmell. *Your life, uploaded: The digital way to better memory, health, and productivity*. Penguin, 2010.



## 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

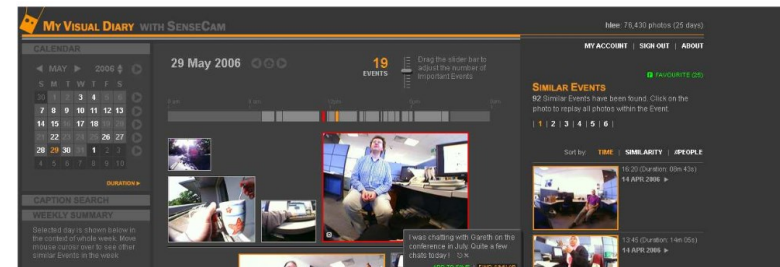
Michael Blighe and Noel E. O'Connor  
 Centre for Digital Video Processing, Adaptive Information Cluster  
 Dublin City University, Ireland  
 {blighem, oconnorn}@eeng.dcu.ie

(lifelogging, place recognition)

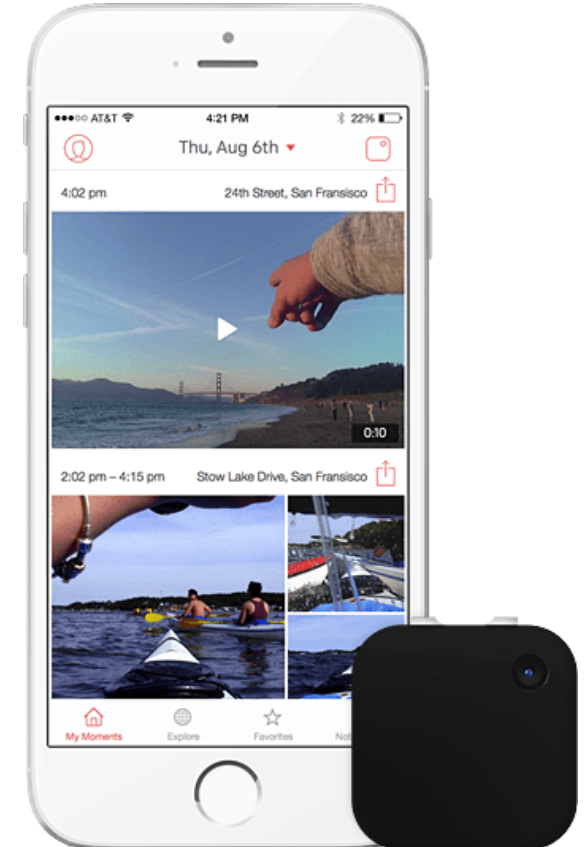
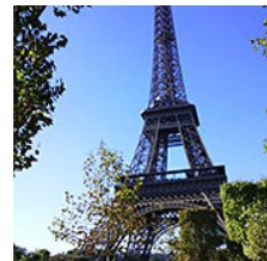
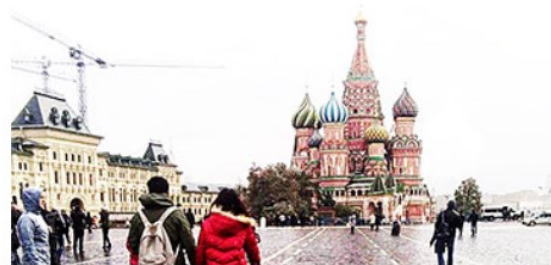
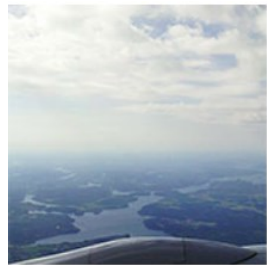
## Constructing a SenseCam Visual Diary as a Media Process

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<sup>a,\*</sup>, Mariella Dimiccoli<sup>a,b</sup>, Petia Radeva<sup>a,b</sup>  
(lifelogging, face tracking)



2016

2017

Day's Lifelog:



Event Segmentation

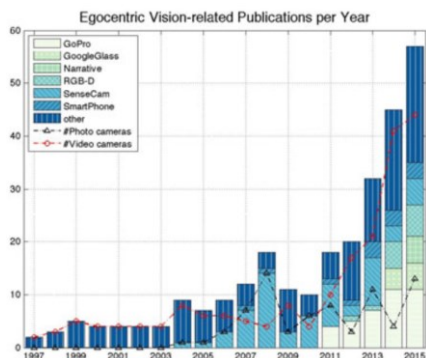
Multiple Events:



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

Mariella Dimiccoli<sup>a,c,1,\*</sup>, Marc Bolaños<sup>a,1,\*</sup>, Estefania Talavera<sup>a,b</sup>, Maedeh Aghaei<sup>a</sup>, Stavri G. Nikolov<sup>d</sup>, Petia Radeva<sup>a,c,\*</sup>

(lifelogging, event segmentation)



## Toward Storytelling From Visual Lifelogging: An Overview

Marc Bolaños, Mariella Dimiccoli, and Petia Radeva

(lifelogging, survey)

2017



## different wearing modalities



head-mounted



chest-mounted



wrist-mounted



helmet-mounted

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



## Fast Unsupervised Ego-Action Learning for First-Person Sports Videos

Kris M. Kitani  
UEC Tokyo  
Tokyo, Japan

Takahiro Okabe, Yoichi Sato  
University of Tokyo  
Tokyo, Japan

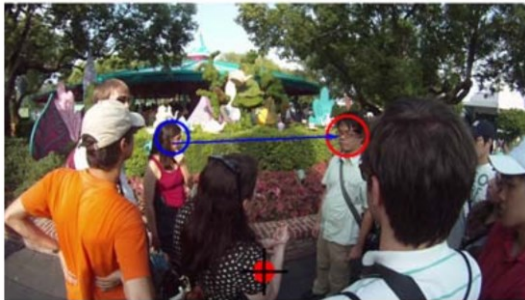
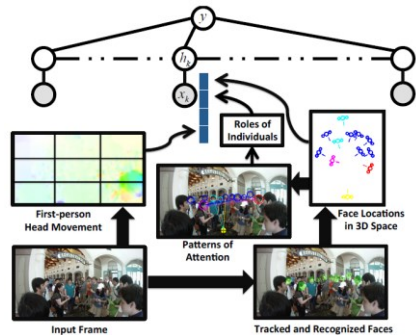
Akihiro Sugimoto  
National Institute of Informatics  
Tokyo, Japan

(unsupervised action recognition, video indexing)



2011

2012



## Social Interactions: A First-Person Perspective

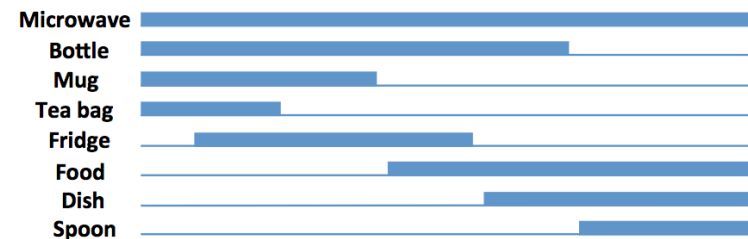
Alireza Fathi<sup>1</sup>, Jessica K. Hodgins<sup>2,3</sup>, James M. Rehg<sup>1</sup>

(detection and recognition of social interactions)

## Story-Driven Summarization for Egocentric Video

Zheng Lu and Kristen Grauman  
University of Texas at Austin

(egocentric video summarization)

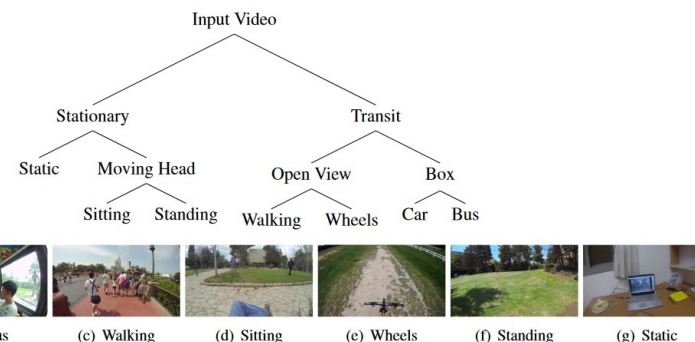


Our method

2013

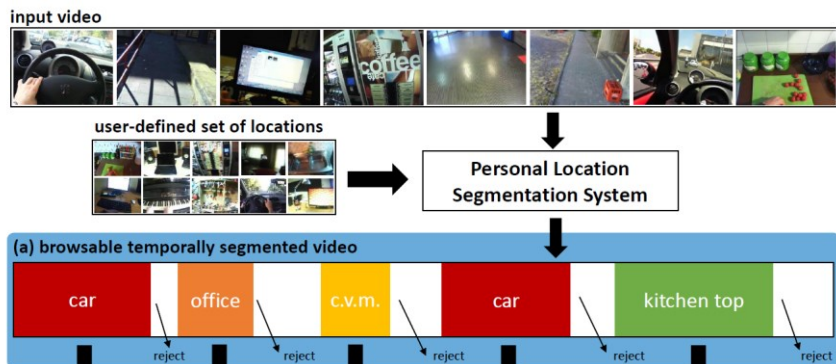
## Temporal Segmentation of Egocentric Videos

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



2014

2016



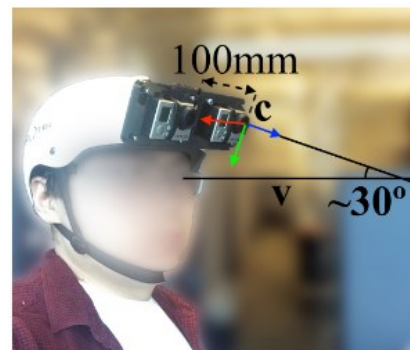
## Recognizing Personal Locations from Egocentric Videos

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

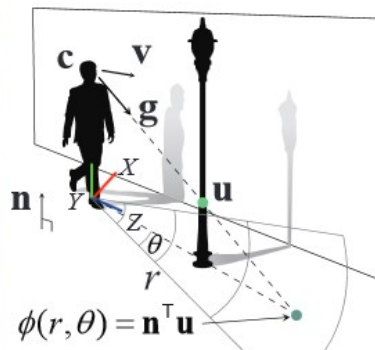
(localization, indexing, context-aware computing)

## Egocentric Future Localization

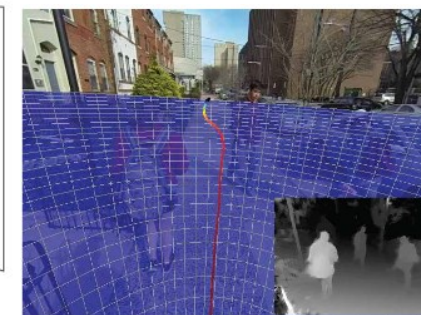
Hyun Soo Park      Jyh-Jing Hwang      Yedong Niu      Jianbo Shi  
 (future localization, navigation)



(a) Ego-stereo cameras



(b) Geometry

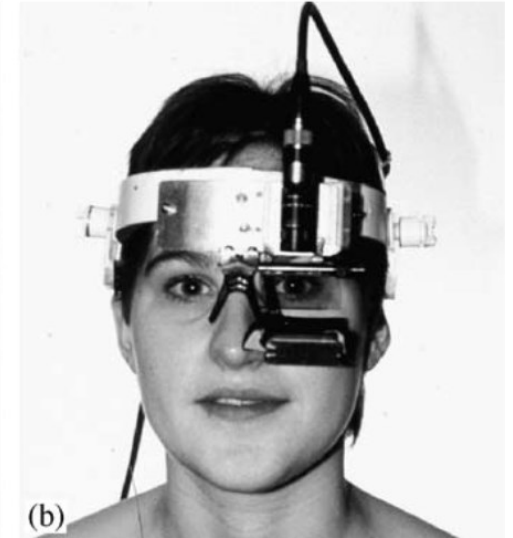
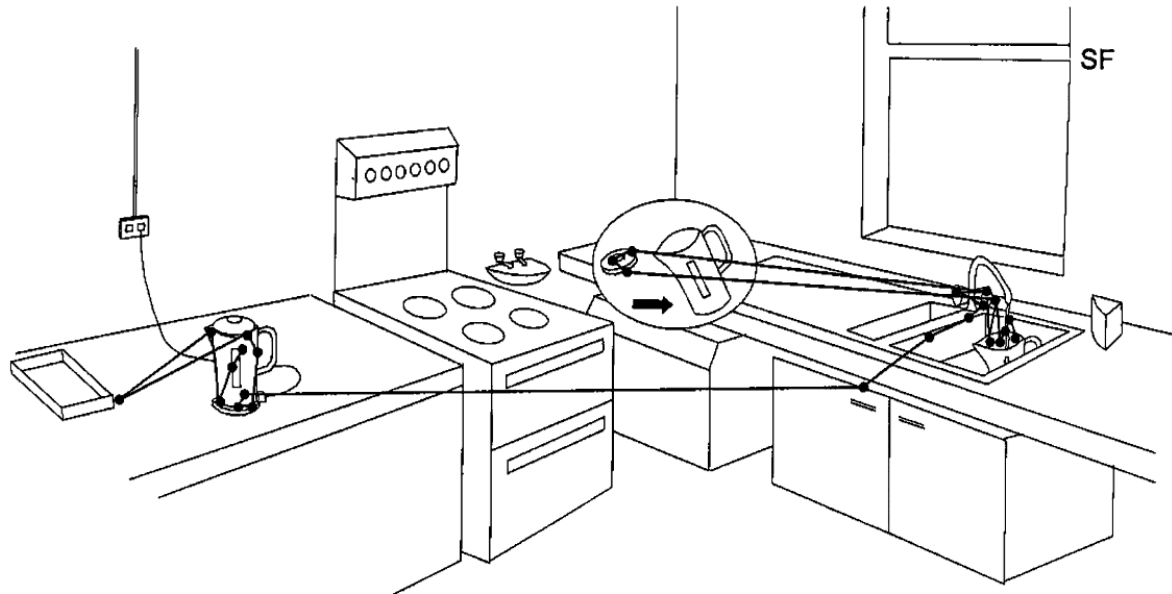


(c) Egocentric RGBD image

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)

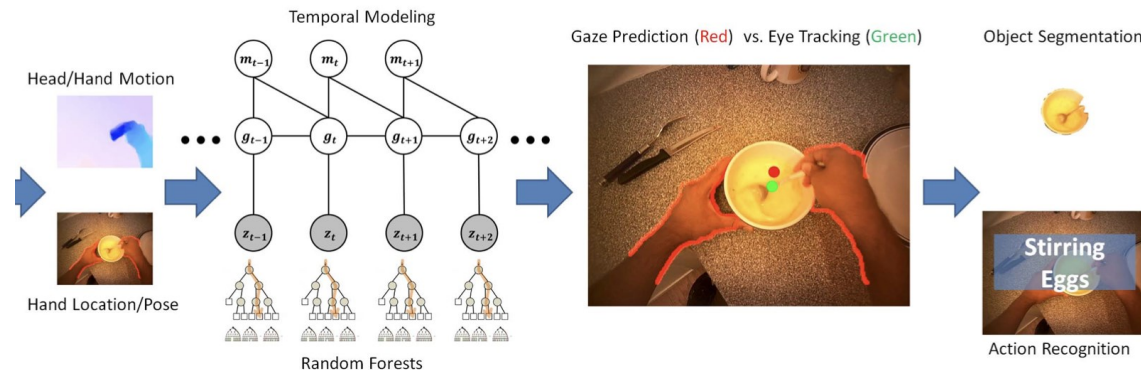


Pupil Eye Trackers (2014 - )



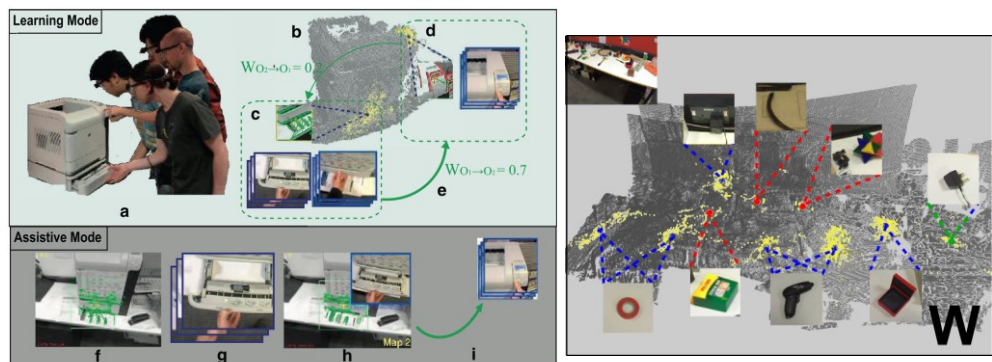
## Learning to Predict Gaze in Egocentric Video

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

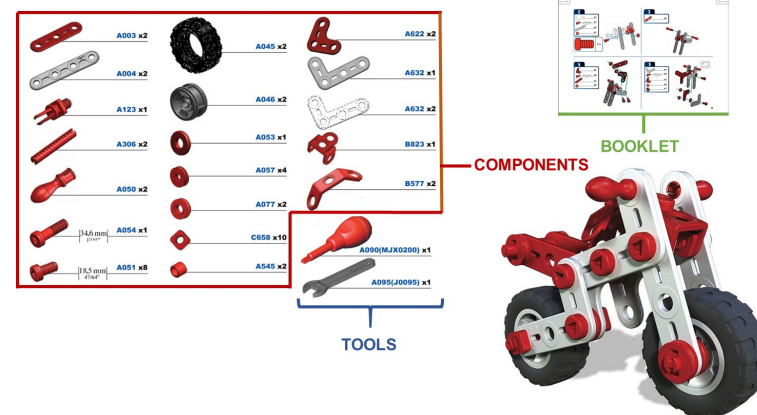
Dima Damen\*, Teesid Leelasawassuk, Walterio Mayol-Cuevas

(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



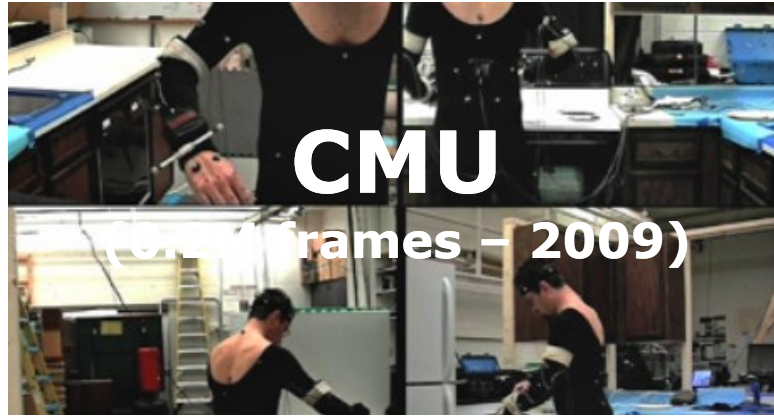
## Workshop on Egocentric (First Person) Vision

# ACVR



**ONE DOES NOT SIMPLY**

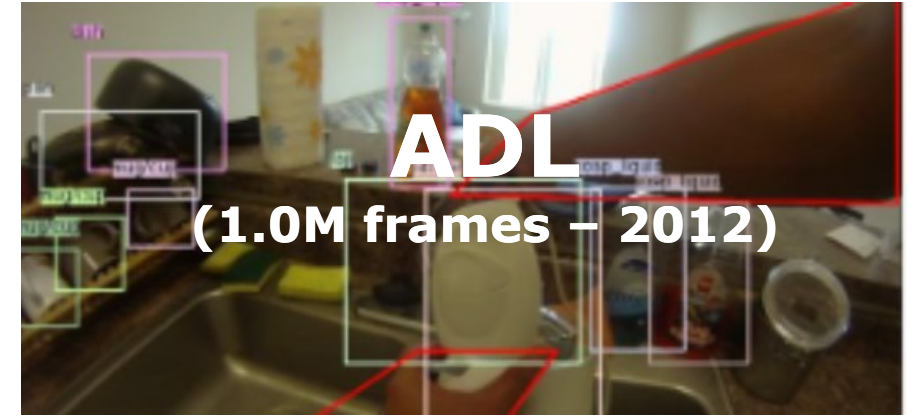
**PUBLISH AN EGOVISION  
PAPER IN A TOP CONFERENCE**



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

# 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



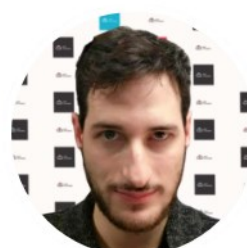
**Michael Wray**  
(Apr 2017 - )  
University of Bristol



**Hazel Doughty**  
(Apr 2017 - )  
University of Bristol



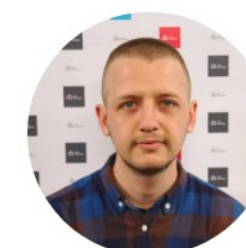
**Toby Perrett**  
(Apr 2017 - )  
University of Bristol



**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



- [Semi-Supervised Video Object Segmentation Challenge](#)
- [Hand-Object Segmentation Challenge](#)
- [TREK-150 Object Tracking Challenge](#)
- [EPIC-SOUNDS Audio-Based Interaction Recognition](#)
- [Action Recognition](#)
- [Action Detection](#)
- [Action Anticipation](#)
- [UDA for Action Recognition](#)
- [Multi-Instance Retrieval](#)

## EPIC-KITCHENS-100- 2022 Challenges Report

## EPIC-KITCHENS-100- 2023 Challenges Report

Dima Damen, Jacob Chalk, Ahmad Darkhalil, Toby Perrett, Daniel Whettam,  
Saptarshi Sinha, Michael Wray, Bin Zhu  
University of Bristol, UK

Antonino Furnari, Francesco Ragusa, Giovanni Maria Farinella      Dandan Shan, David Fouhey  
University of Catania, Italy      University of Michigan, US

Matteo Dunnhofer, Christian Micheloni      Jaesung Huh, Andrew Zisserman  
University of Udine, Italy      University of Oxford, UK

### Abstract

*This report presents the findings from the 5th EPIC-KITCHENS-100 challenges, opened from Jan 2023 and concluded on the 1st of June 2023. It serves as an introduction to all technical reports that were submitted to the 11th EPIC@CVPR2023 workshop, and an official announcement of the winners.*

*The report covers 4 new challenges, announced for the first time in the 2023 round as well as 5 recurring challenges*

### 1. Datasets

The challenges cover three datasets publicly available,

The 5 recurring challenges are based on the publicly available EPIC-KITCHENS-100 dataset. In summary, EPIC-KITCHENS-100 provides 20M frames of egocentric footage, captured in an unscripted manner, with carefully collated annotations of 90K fine-grained actions. Details of how the dataset was collected and annotated are available in our IJCV paper [8]. The challenges are: **Action Recognition**, **Action Anticipation**, **Action Detection**, **Unsupervised**

The TREK-150 Object Tracking is based on the TREK-150 dataset [20] released in 2021. The challenge focuses on single object tracking in egocentric videos.

Finally, one challenge is based on EPIC-SOUNDS [27]. Released last year as well, EPIC-SOUNDS annotates 78.4k categorised segments of audible events and actions, distributed across 44 classes. A challenge on audio-only action recognition was run this year as: **EPIC-SOUNDS Audio-Based Interaction Recognition**.

Each challenge we released codebase with pre-trained models, features and evaluation scripts:

- **Action Recognition** at <https://github.com/epic-kitchens/C1-Action-Recognition>: Five pre-trained models were made available using the codebases: TSN, TRN, TBN, TSM and SlowFast, as well as evaluation script.
- **Action Detection** at <https://github.com/epic-kitchens/C2-Action-Detection>: with pre-extracted features, a baseline using BMN model and evaluation script.
- **Action Anticipation** at <https://github.com/epic-kitchens/C3-Action-Anticipation> with pre-extracted features, RULSTM base model and

## 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

Eleventh International Workshop on Egocentric Perception, Interaction and Computing  
held in conjunction with the 3rd Ego4D Workshop

June 19<sup>th</sup>, 2023



A meme featuring Darth Vader in his iconic black armor and helmet, standing in a crowd of people. The background is slightly blurred, showing other individuals in various attire. The text is overlaid in a bold, white, sans-serif font with a black outline.

**IMPRESSIVE, MOST  
IMPRESSIVE**

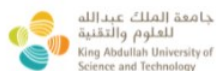
**BUT YOU ARE NOT A JEDI YET**



**Can We Scale?**



## Consortium



INDIANA UNIVERSITY  
BLOOMINGTON



UNIVERSITY  
OF MINNESOTA



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

Kristen Grauman<sup>1,2</sup>, Andrew Westbury<sup>1</sup>, Eugene Byrne<sup>\*1</sup>, Zachary Chavis<sup>\*3</sup>, Antonino Furnari<sup>\*4</sup>, Rohit Girdhar<sup>\*1</sup>, Jackson Hamburger<sup>\*1</sup>, Hao Jiang<sup>\*5</sup>, Miao Liu<sup>\*6</sup>, Xingyu Liu<sup>\*7</sup>, Miguel Martin<sup>\*1</sup>, Tushar Nagarajan<sup>\*1,2</sup>, Ilija Radosavovic<sup>\*8</sup>, Santhosh Kumar Ramakrishnan<sup>\*1,2</sup>, Fiona Ryan<sup>\*6</sup>, Jayant Sharma<sup>\*3</sup>, Michael Wray<sup>\*9</sup>, Mengmeng Xu<sup>\*10</sup>, Eric Zhongcong Xu<sup>\*11</sup>, Chen Zhao<sup>\*10</sup>, Siddhant Bansal<sup>17</sup>, Dhruv Batra<sup>1</sup>, Vincent Cartillier<sup>1,6</sup>, Sean Crane<sup>7</sup>, Tien Do<sup>3</sup>, Morrie Doulaty<sup>13</sup>, Akshay Erapalli<sup>13</sup>, Christoph Feichtenhofer<sup>1</sup>, Adriano Fragomeni<sup>9</sup>, Qichen Fu<sup>7</sup>, Christian Fuegen<sup>13</sup>, Abraham Gebreselasie<sup>12</sup>, Cristina González<sup>14</sup>, James Hillis<sup>5</sup>, Xuhua Huang<sup>7</sup>, Yifei Huang<sup>15</sup>, Wenqi Jia<sup>6</sup>, Weslie Khoo<sup>16</sup>, Jachym Kolar<sup>13</sup>, Satwik Kottur<sup>13</sup>, Anurag Kumar<sup>5</sup>, Federico Landini<sup>13</sup>, Chao Li<sup>5</sup>, Zhenqiang Li<sup>15</sup>, Karttikeya Mangalam<sup>1,8</sup>, Raghava Modhugu<sup>17</sup>, Jonathan Munro<sup>9</sup>, Tullie Murrell<sup>1</sup>, Takumi Nishiyasu<sup>15</sup>, Will Price<sup>9</sup>, Paola Ruiz Puentes<sup>14</sup>, Mery Ramazanova<sup>10</sup>, Leda Sari<sup>5</sup>, Kiran Somasundaram<sup>5</sup>, Audrey Southerland<sup>6</sup>, Yusuke Sugano<sup>15</sup>, Ruijie Tao<sup>11</sup>, Minh Vo<sup>5</sup>, Yuchen Wang<sup>16</sup>, Xindi Wu<sup>7</sup>, Takuma Yagi<sup>15</sup>, Yunyi Zhu<sup>11</sup>, Pablo Arbeláez<sup>†14</sup>, David Crandall<sup>†16</sup>, Dima Damen<sup>†9</sup>, Giovanni Maria Farinella<sup>†4</sup>, Bernard Ghanem<sup>†10</sup>, Vamsi Krishna Ithapu<sup>†5</sup>, C. V. Jawahar<sup>†17</sup>, Hanbyul Joo<sup>†1</sup>, Kris Kitani<sup>†7</sup>, Haizhou Li<sup>†11</sup>, Richard Newcombe<sup>†5</sup>, Aude Oliva<sup>†18</sup>, Hyun Soo Park<sup>†3</sup>, James M. Rehg<sup>†6</sup>, Yoichi Sato<sup>†15</sup>, Jianbo Shi<sup>†19</sup>, Mike Zheng Shou<sup>†11</sup>, Antonio Torralba<sup>†18</sup>, Lorenzo Torresani<sup>†1,20</sup>, Mingfei Yan<sup>†5</sup>, Jitendra Malik<sup>1,8</sup>

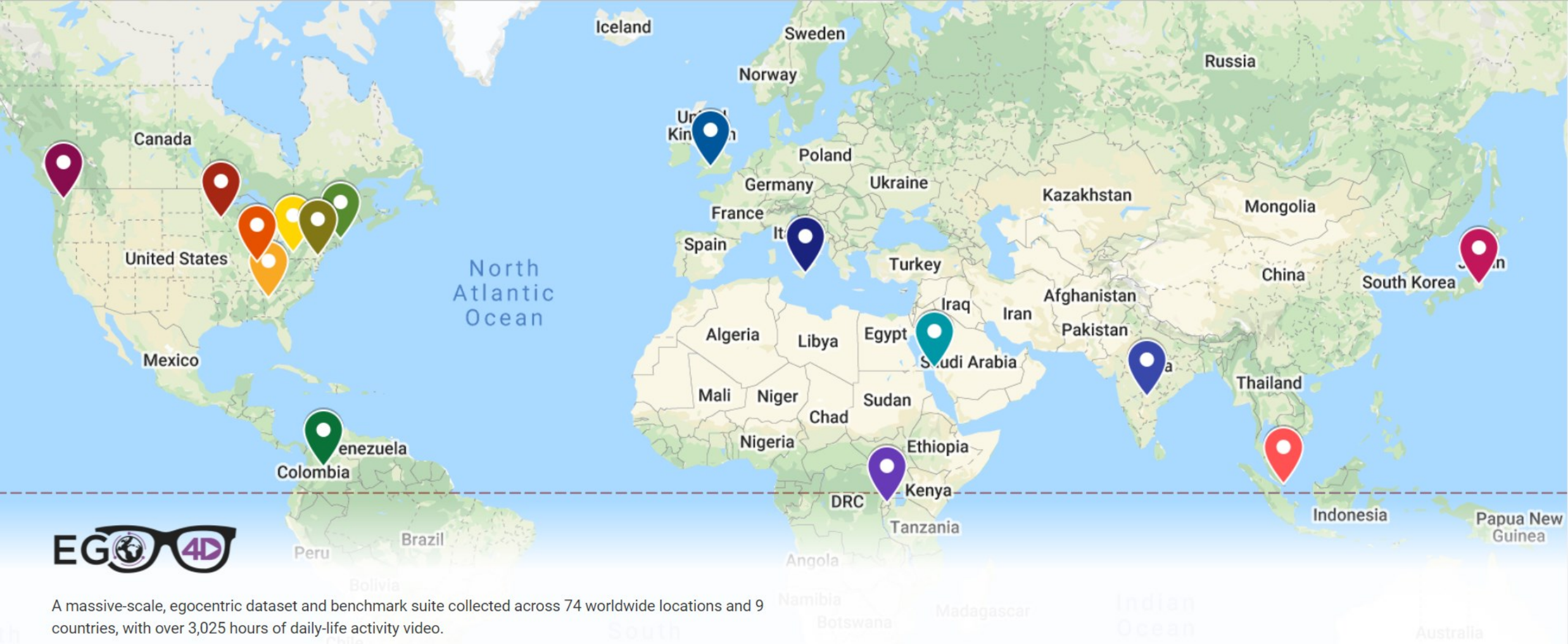
<sup>1</sup>Facebook AI Research (FAIR), <sup>2</sup>University of Texas at Austin, <sup>3</sup>University of Minnesota, <sup>4</sup>University of Catania,

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

<sup>10</sup>King Abdullah University of Science and Technology, <sup>11</sup>National University of Singapore,

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

<sup>17</sup>International Institute of Information Technology, Hyderabad, <sup>18</sup>MIT, <sup>19</sup>University of Pennsylvania, <sup>20</sup>Dartmouth



A massive-scale, egocentric dataset and benchmark suite collected across 74 worldwide locations and 9 countries, with over 3,025 hours of daily-life activity video.



855 Subjects



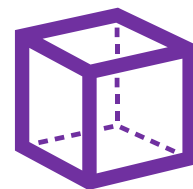
74 Locations



9 Countries



3025 Hours



3D Scans



Audio



Gaze



**Episodic Memory**



**Hand-Object  
Interactions**



**AV Diarization**



**Social**



**Forecasting**

## 1st Ego4D Workshop @ CVPR 2022

Held in conjunction with [10th EPIC Workshop](#)

**19 and 20 June 2022**

## 2nd International Ego4D Workshop @ ECCV 2022

**24 October 2022**

## 3rd International Ego4D Workshop @ CVPR 2023

Held in conjunction with 11th EPIC Workshop

**19 June 2023**





**EGO-EXO4D**

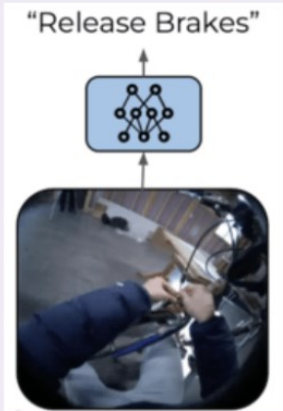




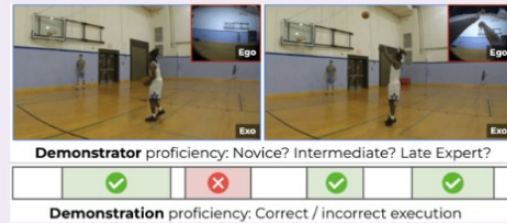
## **Ego-Exo4D: Understanding Skilled Human Activity from First- and Third-Person Perspectives**

Kristen Grauman<sup>1,2</sup>, Andrew Westbury<sup>1</sup>, Lorenzo Torresani<sup>1</sup>, Kris Kitani<sup>1,3</sup>, Jitendra Malik<sup>1,4</sup>, Triantafyllos Afouras<sup>\*1</sup>, Kumar Ashutosh<sup>\*1,2</sup>, Vijay Baiyya<sup>\*5</sup>, Siddhant Bansal<sup>\*6,7</sup>, Bikram Boote<sup>\*8</sup>, Eugene Byrne<sup>\*1,9</sup>, Zach Chavis<sup>\*10</sup>, Joya Chen<sup>\*11</sup>, Feng Cheng<sup>\*1</sup>, Fu-Jen Chu<sup>\*1</sup>, Sean Crane<sup>\*9</sup>, Avijit Dasgupta<sup>\*7</sup>, Jing Dong<sup>\*5</sup>, Maria Escobar<sup>\*12</sup>, Cristhian Forigua<sup>\*12</sup>, Abrrham Gebreselasie<sup>\*9</sup>, Sanjay Haresh<sup>\*13</sup>, Jing Huang<sup>\*1</sup>, Md Mohaiminul Islam<sup>\*14</sup>, Suyog Jain<sup>\*1</sup>, Rawal Khirodkar<sup>\*9</sup>, Devansh Kukreja<sup>\*1</sup>, Kevin J Liang<sup>\*1</sup>, Jia-Wei Liu<sup>\*11</sup>, Sagnik Majumder<sup>\*1,2</sup>, Yongsen Mao<sup>\*13</sup>, Miguel Martin<sup>\*1</sup>, Effrosyni Mavroudi<sup>\*1</sup>, Tushar Nagarajan<sup>\*1</sup>, Francesco Ragusa<sup>\*15</sup>, Santhosh Kumar Ramakrishnan<sup>\*2</sup>, Luigi Seminara<sup>\*15</sup>, Arjun Somayazulu<sup>\*2</sup>, Yale Song<sup>\*1</sup>, Shan Su<sup>\*16</sup>, Zihui Xue<sup>\*1,2</sup>, Edward Zhang<sup>\*16</sup>, Jinxu Zhang<sup>\*16</sup>, Angela Castillo<sup>12</sup>, Changan Chen<sup>2</sup>, Xinzhu Fu<sup>11</sup>, Ryosuke Furuta<sup>17</sup>, Cristina González<sup>12</sup>, Prince Gupta<sup>5</sup>, Jiabo Hu<sup>18</sup>, Yifei Huang<sup>17</sup>, Yiming Huang<sup>16</sup>, Weslie Khoo<sup>19</sup>, Anush Kumar<sup>10</sup>, Robert Kuo<sup>18</sup>, Sach Lakhavani<sup>5</sup>, Miao Liu<sup>18</sup>, Mi Luo<sup>2</sup>, Zhengyi Luo<sup>3</sup>, Brighid Meredith<sup>18</sup>, Austin Miller<sup>18</sup>, Oluwatumininu Oguntola<sup>14</sup>, Xiaqing Pan<sup>5</sup>, Penny Peng<sup>18</sup>, Shraman Pramanick<sup>20</sup>, Merey Ramazanova<sup>21</sup>, Fiona Ryan<sup>22</sup>, Wei Shan<sup>14</sup>, Kiran Somasundaram<sup>5</sup>, Chenan Song<sup>11</sup>, Audrey Southerland<sup>22</sup>, Masatoshi Tateno<sup>17</sup>, Huiyu Wang<sup>1</sup>, Yuchen Wang<sup>19</sup>, Takuma Yagi<sup>17</sup>, Mingfei Yan<sup>5</sup>, Xitong Yang<sup>1</sup>, Zecheng Yu<sup>17</sup>, Shengxin Cindy Zha<sup>18</sup>, Chen Zhao<sup>21</sup>, Ziwei Zhao<sup>19</sup>, Zhifan Zhu<sup>6</sup>, Jeff Zhuo<sup>14</sup>, Pablo Arbeláez<sup>†12</sup>, Gedas Bertasius<sup>†14</sup>, David Crandall<sup>†19</sup>, Dima Damen<sup>†6</sup>, Jakob Engel<sup>†5</sup>, Giovanni Maria Farinella<sup>†15</sup>, Antonino Furnari<sup>†15</sup>, Bernard Ghanem<sup>†21</sup>, Judy Hoffman<sup>†22</sup>, C. V. Jawahar<sup>†7</sup>, Richard Newcombe<sup>†5</sup>, Hyun Soo Park<sup>†10</sup>, James M. Rehg<sup>†8</sup>, Yoichi Sato<sup>†17</sup>, Manolis Savva<sup>†13</sup>, Jianbo Shi<sup>†16</sup>, Mike Zheng Shou<sup>†11</sup>, and Michael Wray<sup>†6</sup>

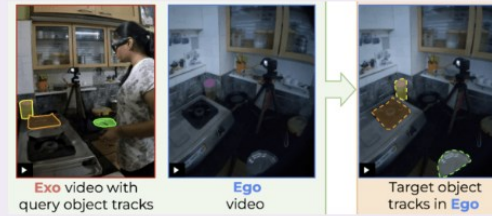
<https://ego-exo4d-data.org/>



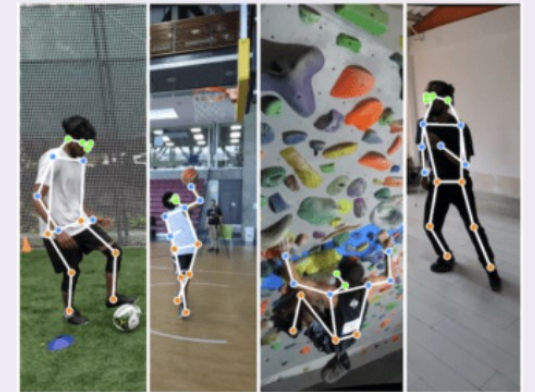
**Keystep Recognition**



**Proficiency Estimation**



**Relation**



**Pose Estimation**

# First Joint Egocentric Vision (EgoVis) Workshop

## Held in Conjunction with CVPR 2024

17 June 2024 - Seattle, USA



Ego-Exo4D



Ego4D



EPIC-Kitchens

## MECCANO: A Multimodal Egocentric Dataset for Humans Behavior Understanding in the Industrial-like Domain

F. Ragusa<sup>1,2</sup>, A. Furnari<sup>1,2</sup>, G. M. Farinella<sup>1,2</sup>

<sup>1</sup>FPV@IPLab, Department of Mathematics and Computer Science - University of Catania, Italy

<sup>2</sup>Next Vision s.r.l., Spin-off of the University of Catania, Italy

**Running ICIAP competition with Prize!**

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

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

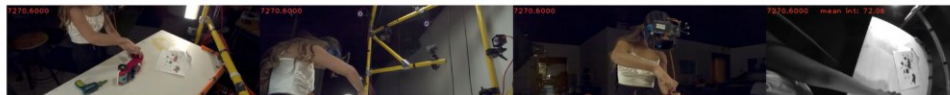
Fadime Sener<sup>1</sup>    Dibyadip Chatterjee<sup>2</sup>    Daniel Shelepov<sup>1</sup>    Kun He<sup>1</sup>  
Dipika Singhanian<sup>2</sup>    Robert Wang<sup>1</sup>    Angela Yao<sup>2</sup>

<sup>1</sup>Reality Labs at Meta

<sup>2</sup>National University of Singapore

CVPR 2022

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



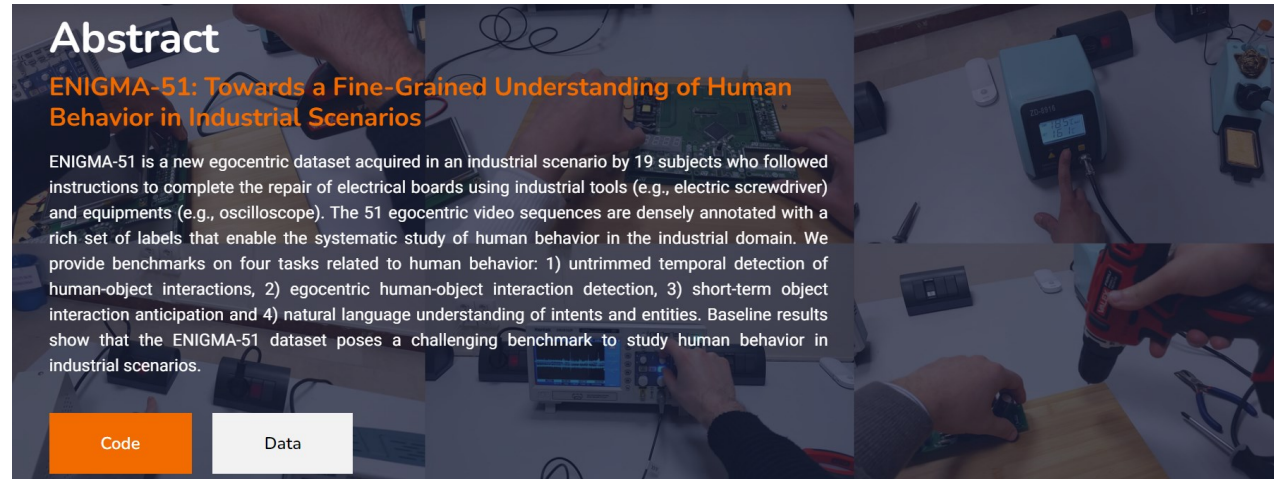
### Abstract

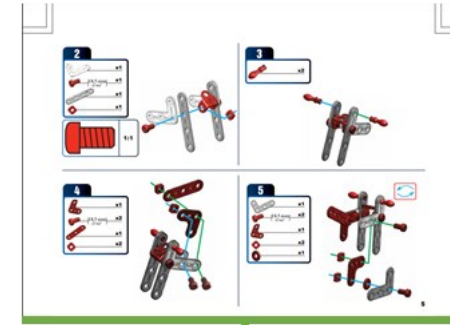
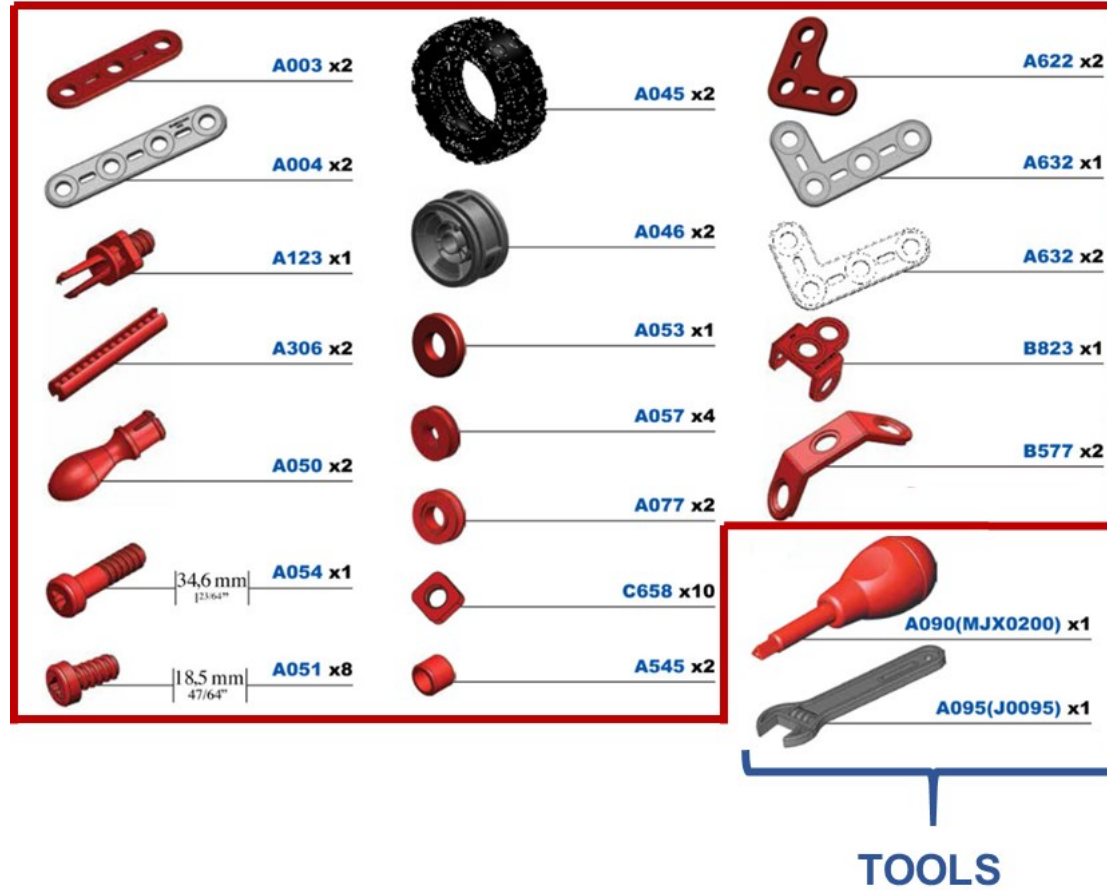
#### ENIGMA-51: Towards a Fine-Grained Understanding of Human Behavior in Industrial Scenarios

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](#)





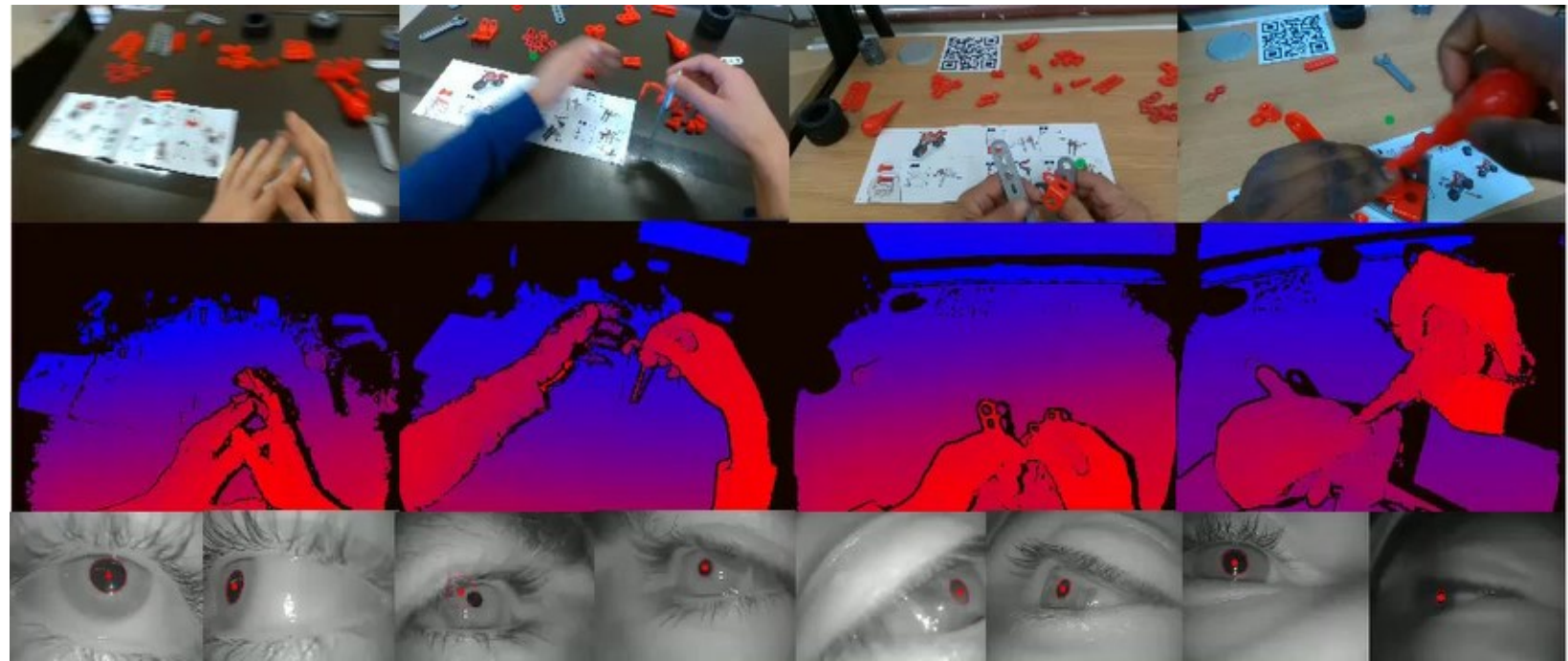
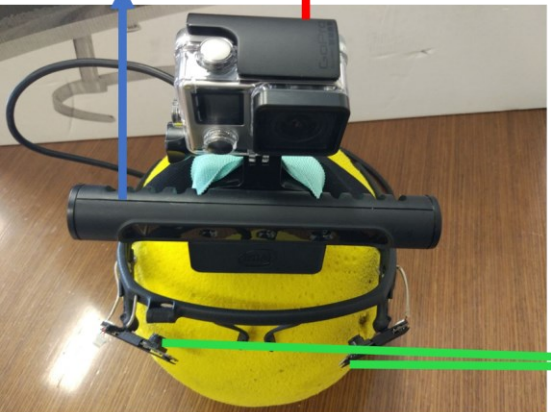
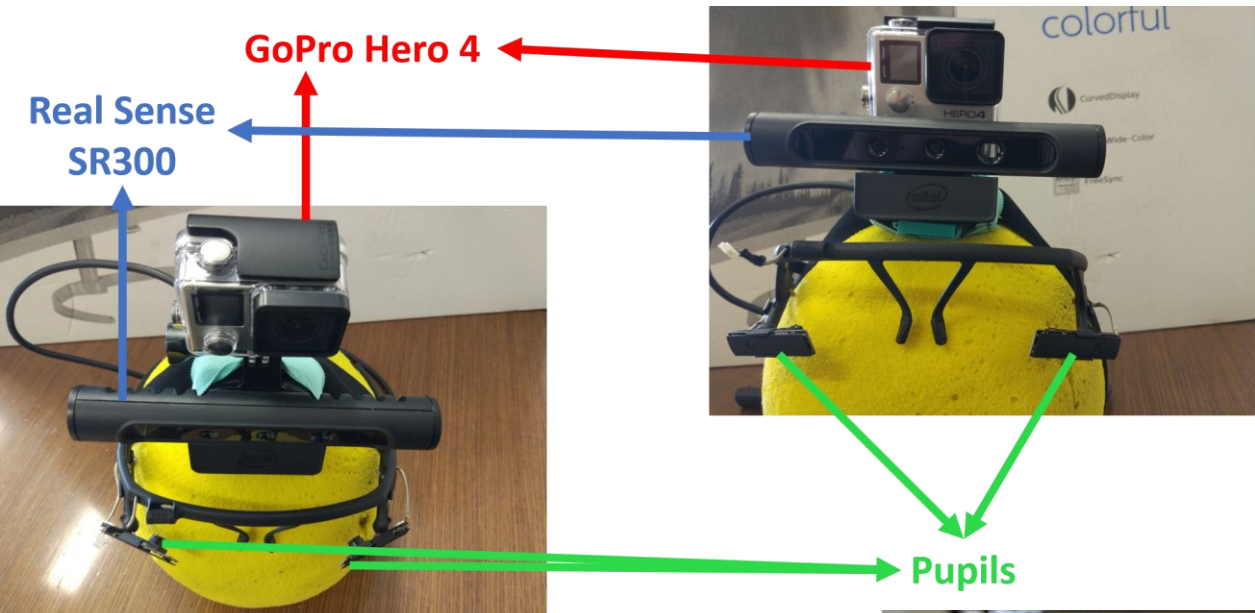
BOOKLET

COMPONENTS

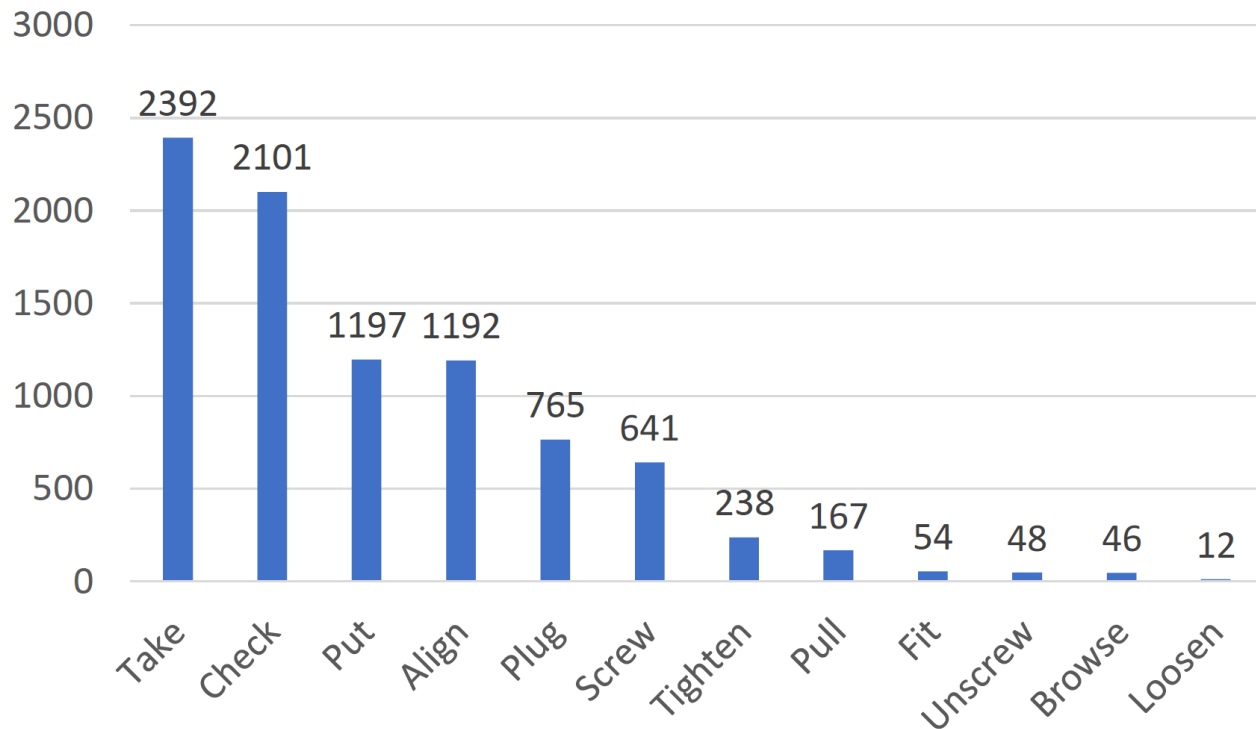


TOOLS

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

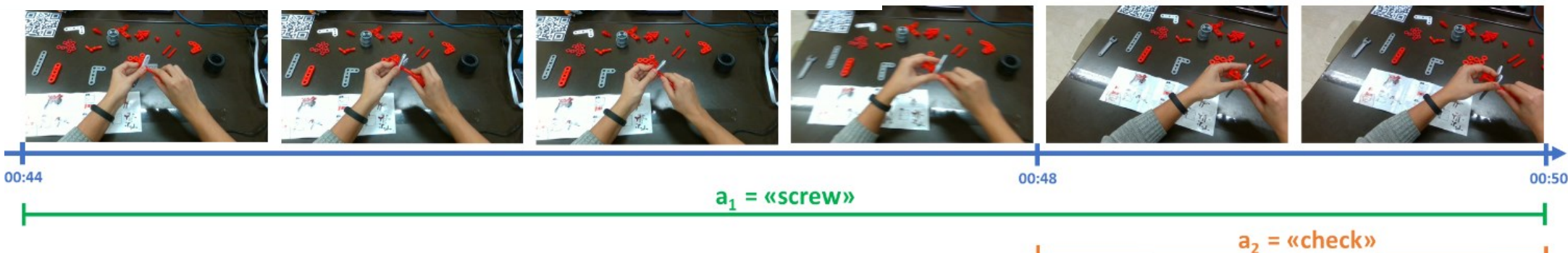


Verbs Classes



**8857 video segments**

**1401 overlap segments (15.82%)**





# Data Annotation: Active Object Bounding Boxes



red\_perforated\_bar



gray\_bar



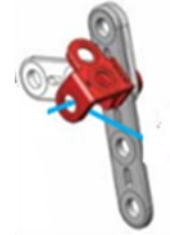
wheels\_axle



bar



handlebar



partial\_model



gray\_angled\_bar



bolt



red\_3\_junction\_bar



wrench



tire



rim



washer



white\_bar



instruction\_booklet



cylinder



red\_angled\_bar



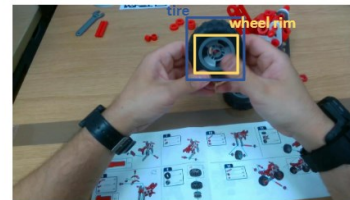
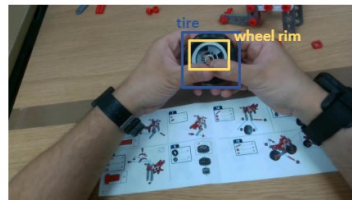
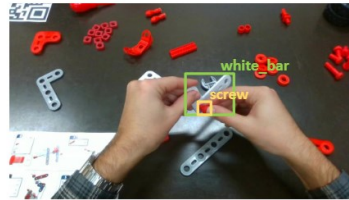
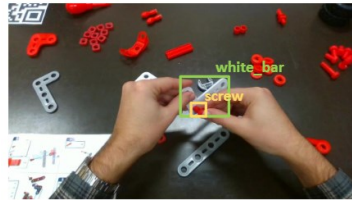
screw



red\_4\_junction\_bar

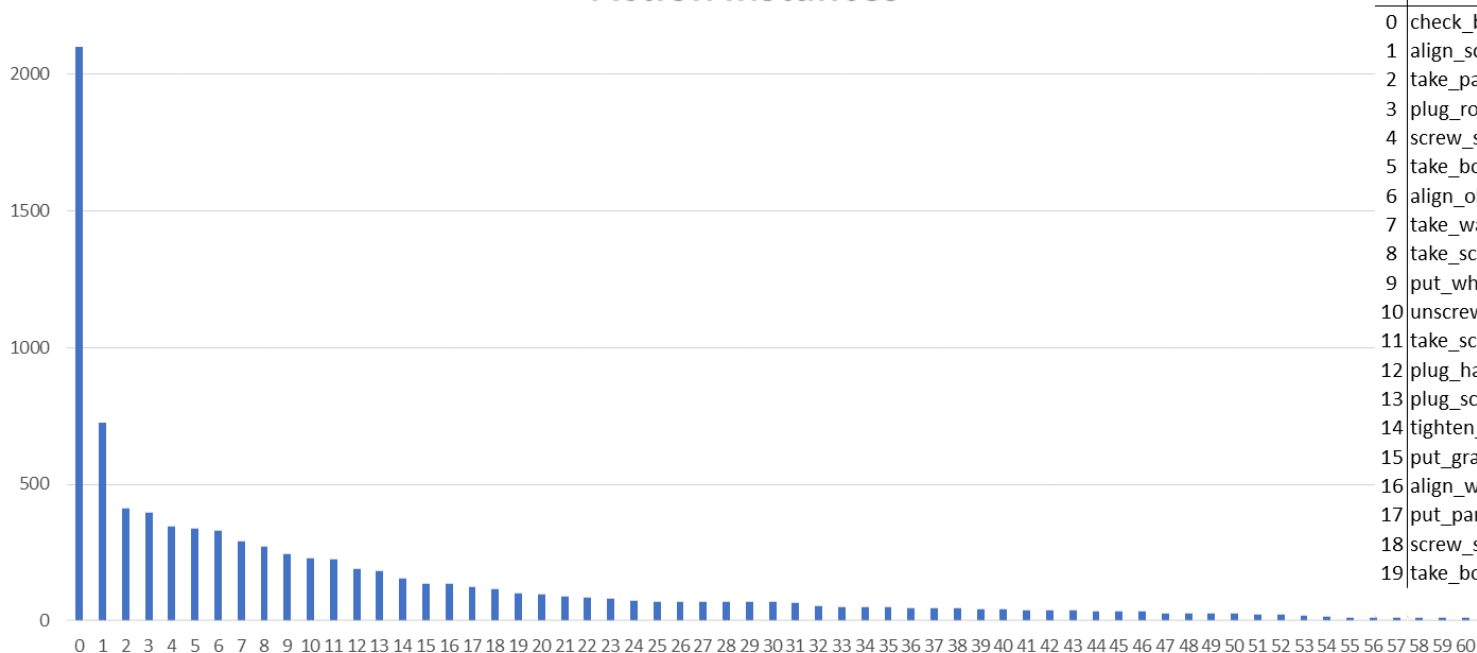


screwdriver



**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

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

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

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

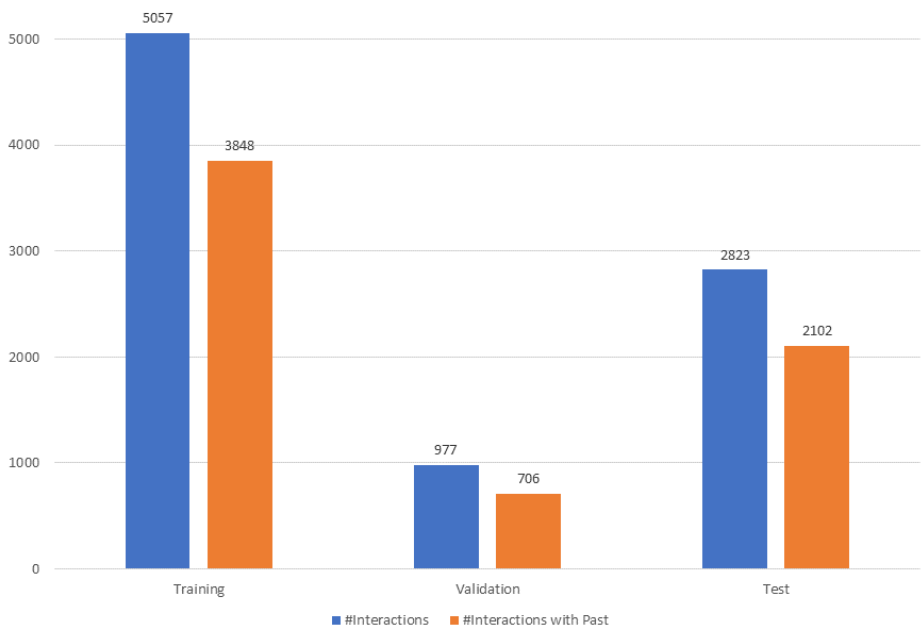


<take, screwdriver>

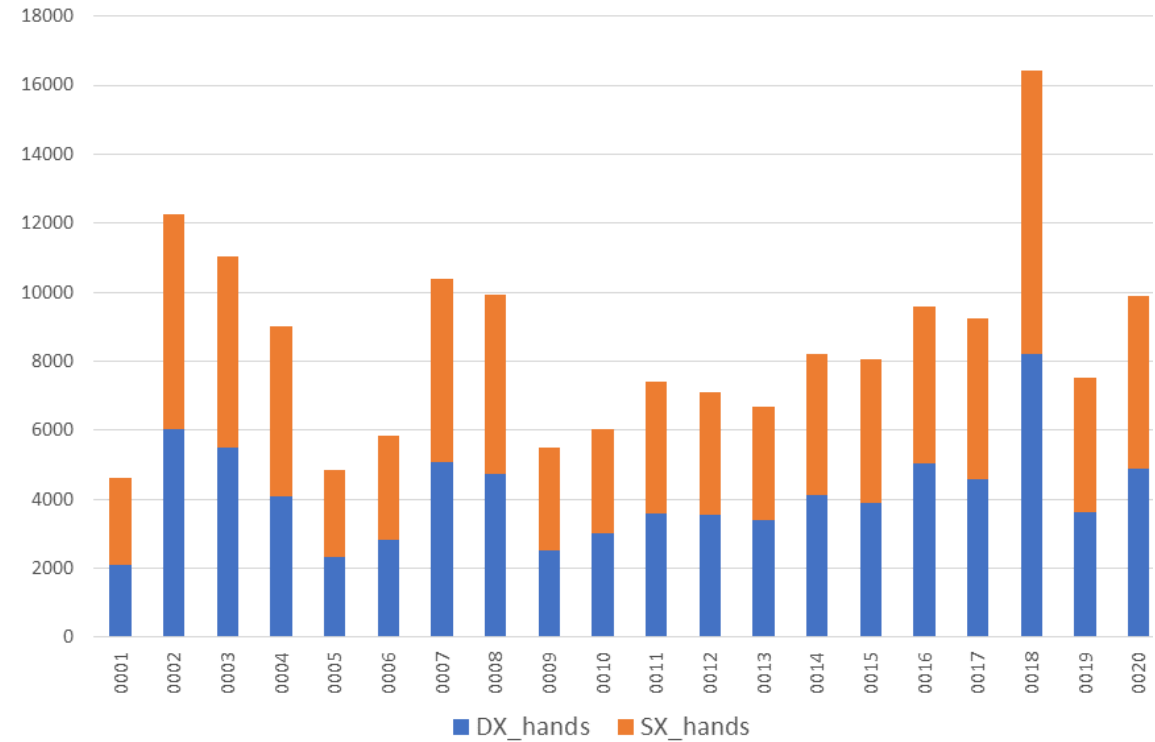
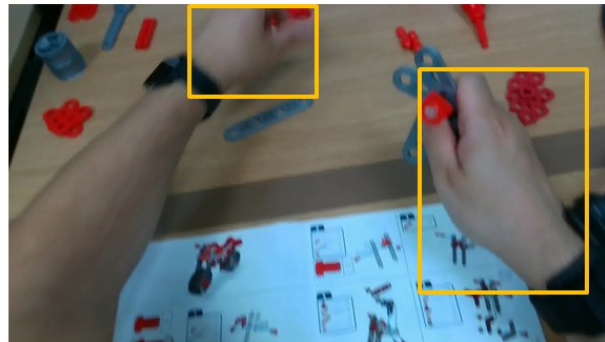
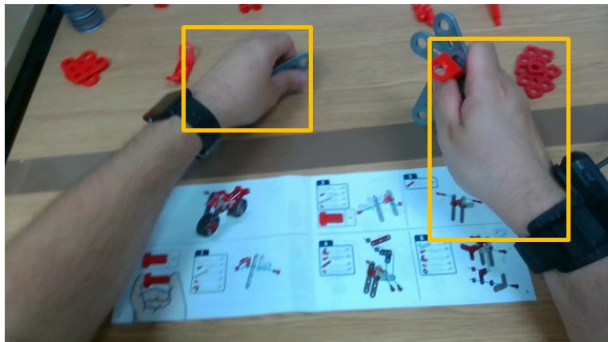
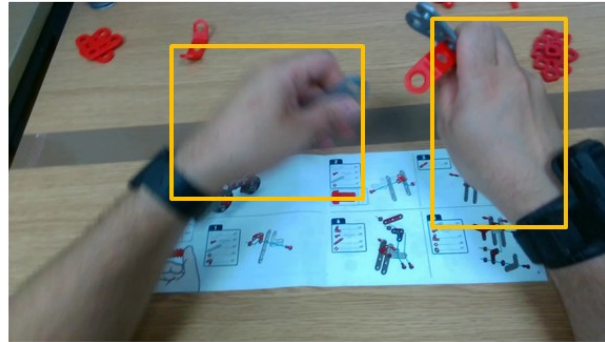
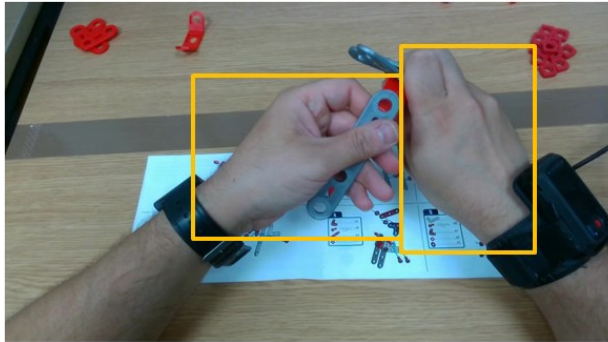


<screw, {screwdriver, screw, partial\_model}>

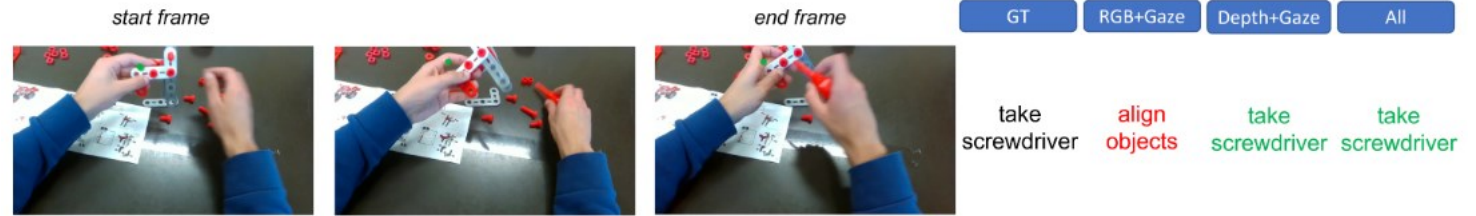
# Data Annotation: Next Active Object Annotations



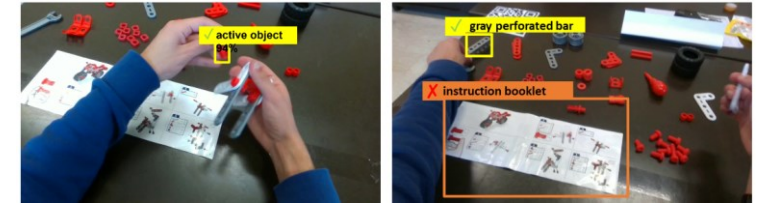
Video	Interactions	Interactions with past
0001	319	257
0002	586	452
0003	573	429
0004	485	372
0005	251	200
0006	307	234
0007	493	367
0008	550	384
0009	289	289
0010	304	194
0011	400	310
0012	384	258
0013	313	244
0014	434	297
0015	425	324
0016	576	436
0017	484	339
0018	788	603
0019	400	294
0020	496	373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



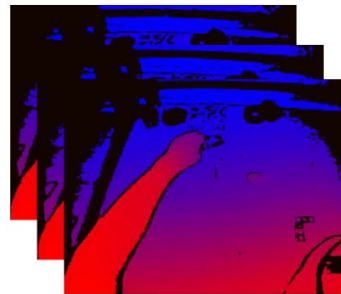
# 1) Action Recognition



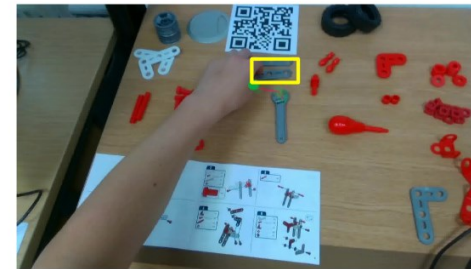
# 2) Active Object Detection and Recognition



# 3) EHOI Detection







<take>



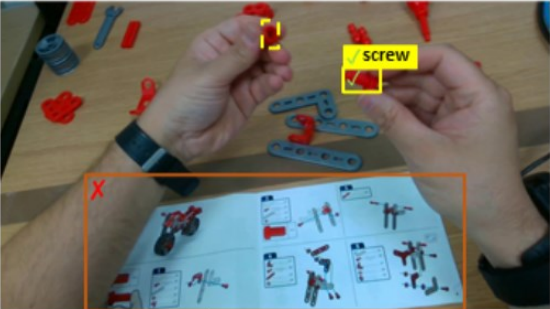
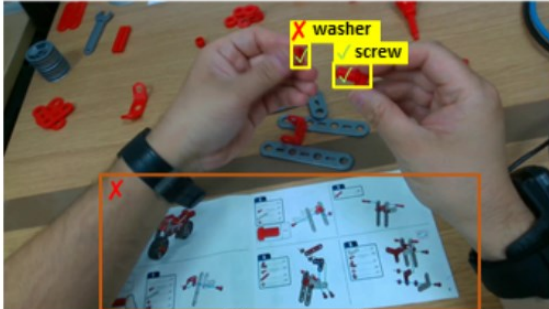
<gray perforated bar>

## 4) Action Anticipation

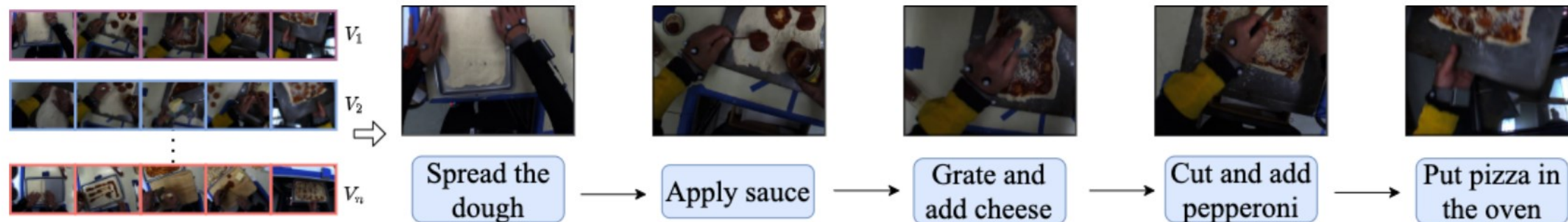
Ground Truth action: **take bolt**

$\tau_a = 2.00$	$\tau_a = 1.50$	$\tau_a = 1.00$	$\tau_a = 0.25$
			
<i>take bolt, align objects, tighten <b>bolt</b>, plug screw, check booklet</i>	<i>take bolt, align objects, plug screw, tighten <b>bolt</b>, check booklet</i>	<i>take bolt, align objects, plug screw, check booklet, tighten <b>bolt</b></i>	<i>take bolt, align objects, plug screw, check booklet, take screwdriver</i>

## 5) Next-Active Object (NAO) Detection

 <p>Time to start = 1.6s</p>	 <p>Time to start = 0.8s</p>
---	--

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





Spin-off of the University of Catania

## Leaderboard 2023

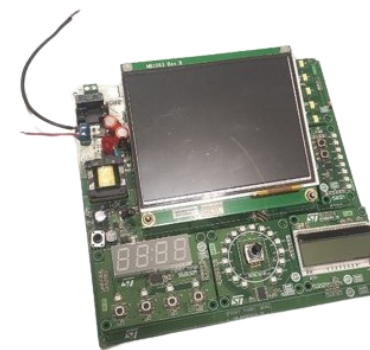
Rank	Team	Top-1 Accuracy	Top-5 Accuracy	Technical Report
🏆 1	UCF	52.82	83.85	↓
🌟 2	UNIBZ	52.57	81.53	↓
🌟 3	LUBECK	51.82	83.35	↓
4	MACAU	50.30	78.46	↓
5	<i>Baseline (RGB-Depth-Gaze)</i>	49.66	77.82	
6	TORONTO	49.52	74.21	↓
7	<i>Baseline (RGB-Depth)</i>	49.49	77.61	
8	CUNY	24.69	52.46	↓

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

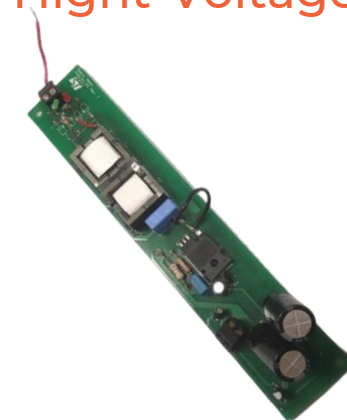


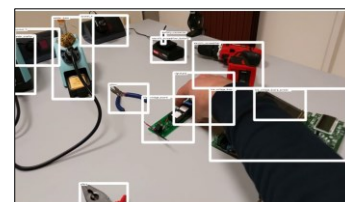
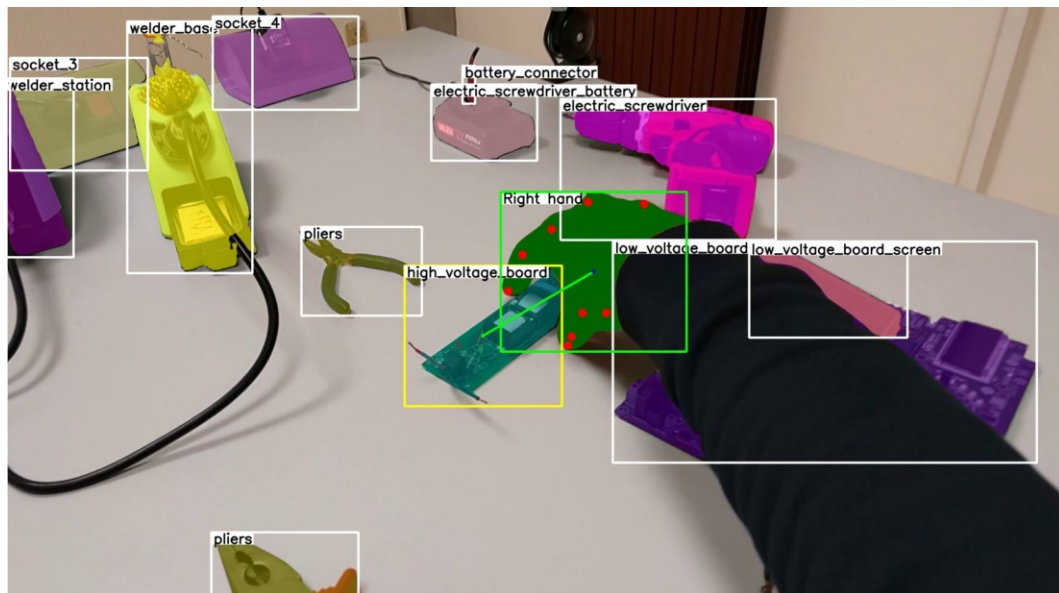
We designed two procedures consisting of instructions that involve humans interacting with the objects present in the laboratory to achieve the goal of repairing two electrical boards

Low-Voltage

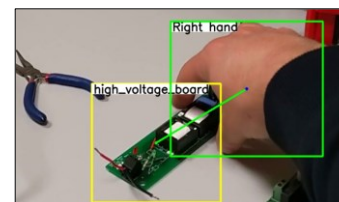


High-Voltage





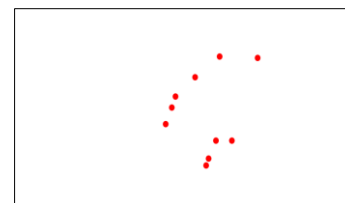
Hand-Object boxes



Human-Object Interactions



Hand-Object Masks



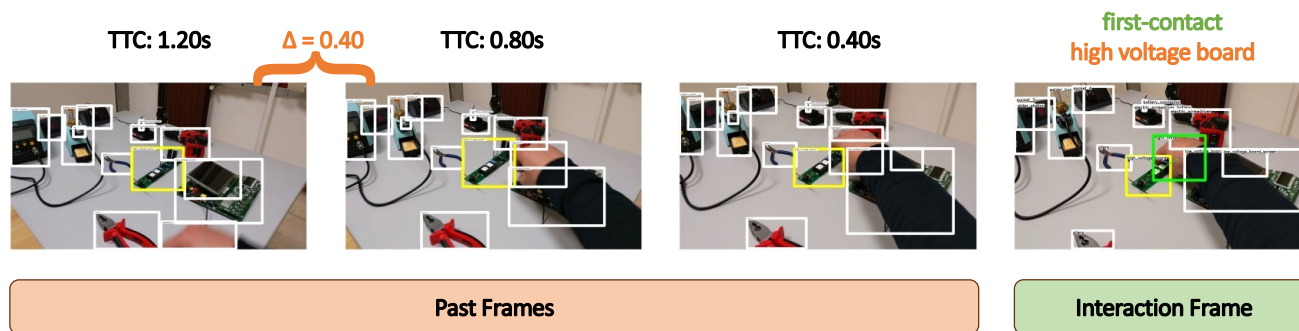
Hand Keypoints



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



- Google envisioned a future in which smart glasses replace smartphones;
- The goal of Google Glass was to make computation available to the user when they need it and get out of the way when they don't.

<https://www.youtube.com/watch?v=YAXTQL3jPFk>

<https://www.youtube.com/watch?v=ClvI9fZaz6M>



**Google Glass failed because of the lack of clear use cases + privacy issues.**

## Is this it?

SenseCam



2004

Vicon Revue



2010

Autographer



2013

Looxcie



2010

Google Glass



2012



**Success Cases**



## Moverio BT-40

- USB-C connectivity
- Full HD 1080p
- Second screen privacy

OUR PRICE:

**£579.00**

incl. VAT (£482.50 ex. VAT)

In Stock

[Learn more ▶](#)

[Buy Now ▶](#)

[FIND A DEALER ▶](#)

[REQUEST A CALLBACK ▶](#)

[SUPPORT ▶](#)



## Moverio BT-40S

- Intelligent Controller
- Full HD 1080p
- Commercial applications

OUR PRICE:

**£1,002.00**

incl. VAT (£835.00 ex. VAT)

In Stock

[Learn more ▶](#)

[Buy Now ▶](#)

[FIND A DEALER ▶](#)

[REQUEST A CALLBACK ▶](#)

[SUPPORT ▶](#)



## Moverio BT-45CS

- Centred 8MP camera
- Rugged design
- Intelligent Controller

OUR PRICE:

**£1,836.00**

incl. VAT (£1,530.00 ex. VAT)

In Stock

[Learn more ▶](#)

[Buy Now ▶](#)

[FIND A DEALER ▶](#)

[REQUEST A CALLBACK ▶](#)

[SUPPORT ▶](#)

## focused application scenarios

[https://www.epson.co.uk/en\\_GB/search/allproducts?text=smart+glasses](https://www.epson.co.uk/en_GB/search/allproducts?text=smart+glasses)





## Manufacturing Solutions

---

LEARN MORE

## Warehouse Solutions

---

LEARN MORE

## Field Service & Remote Assist Solutions

---

LEARN MORE

## Tele-Medicine Solutions

---

LEARN MORE

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



**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**

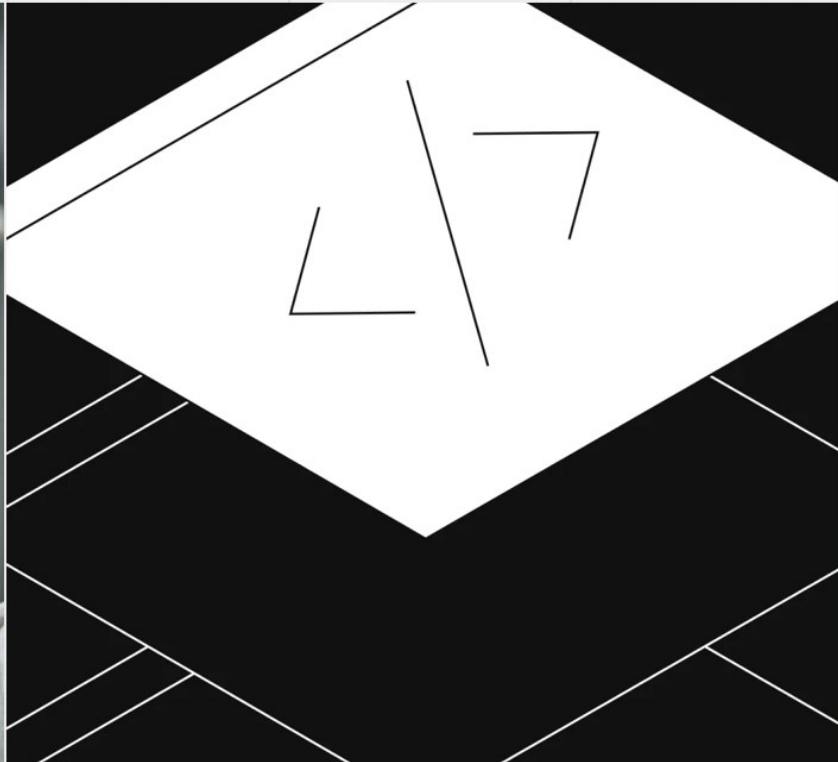


<https://www.magicleap.com/magic-leap-2>



## Scalable

Magic Leap 2 is built to support scalable augmented reality (AR) solutions necessitating multiple simultaneous users.



## Integrative

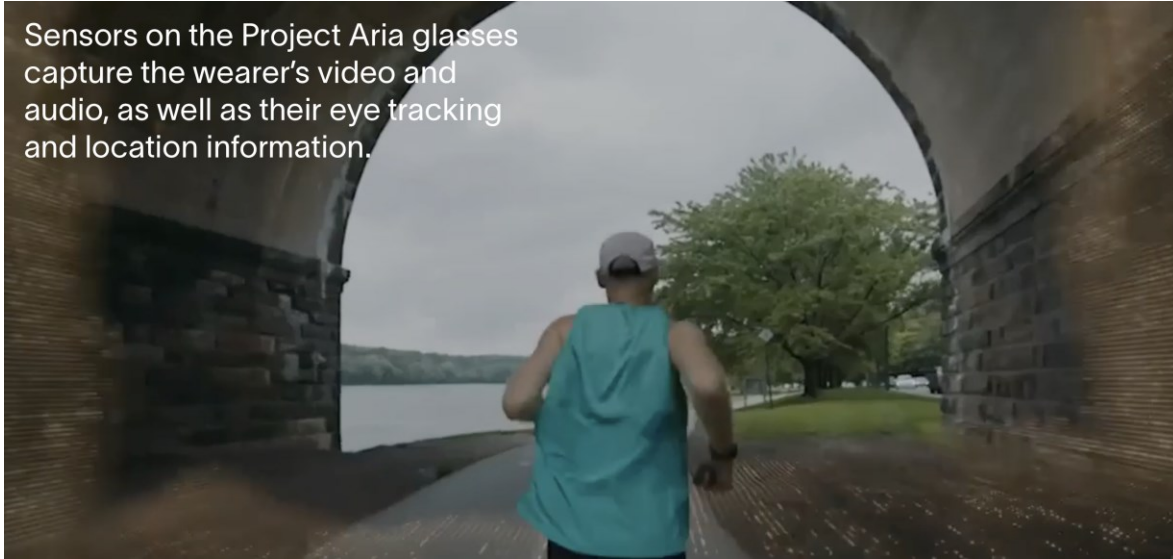
Magic Leap 2 is purpose-built on an open platform to integrate with leading enterprise multi-device management (MDM) systems.



## Secure

Store your data anywhere and use any preferred cloud setup. Magic Leap 2 lets users retain control of their data and is compatible with leading enterprise security protocols.

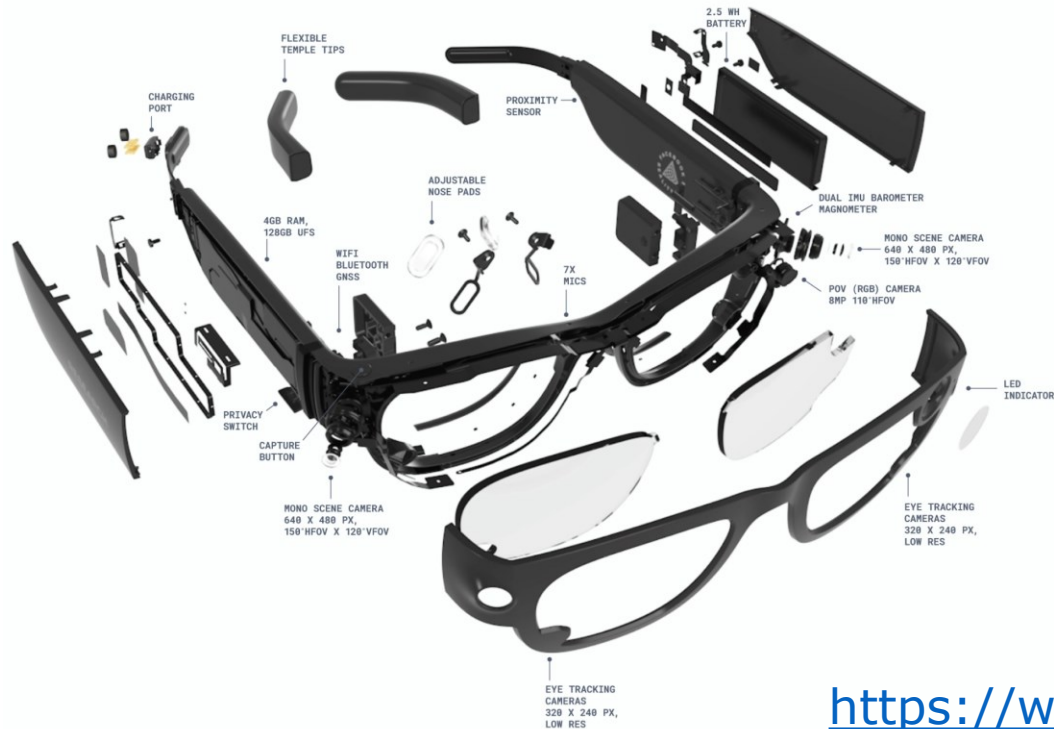
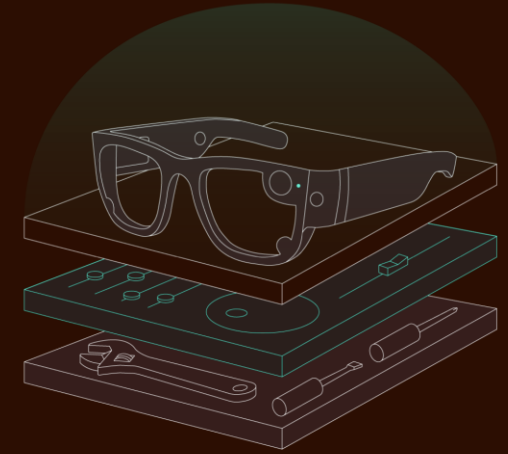
Sensors on the Project Aria glasses capture the wearer's video and audio, as well as their eye tracking and location information.



## 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](#)





**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

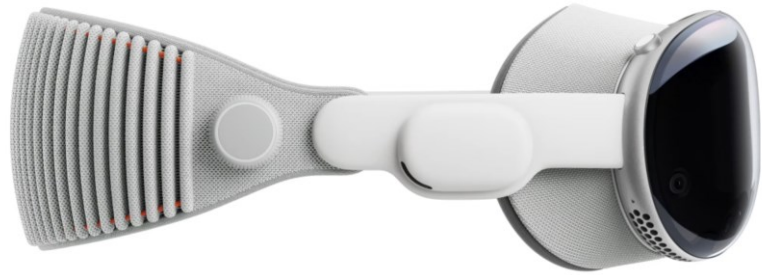
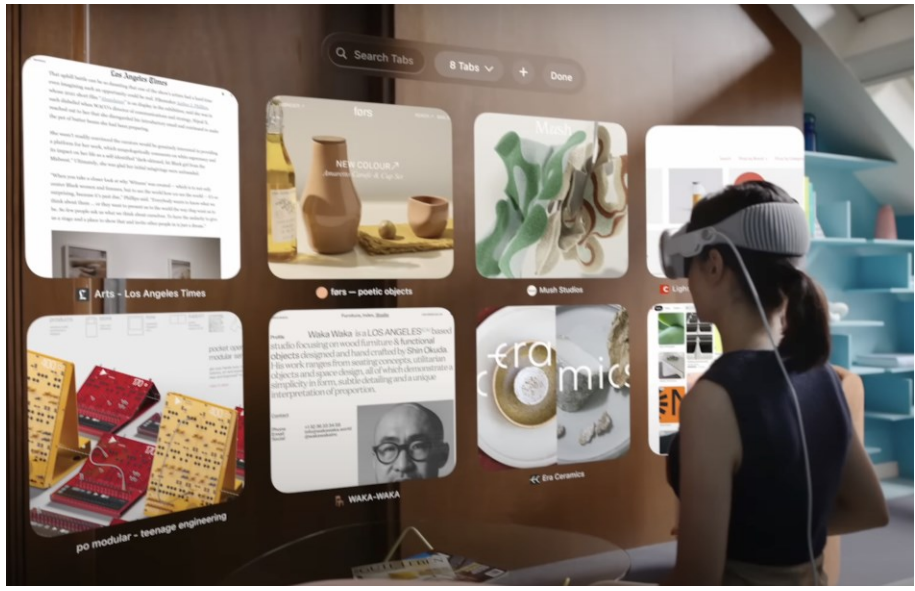
Allowing 3D surface and object detection for realistic AR integration with the real world.

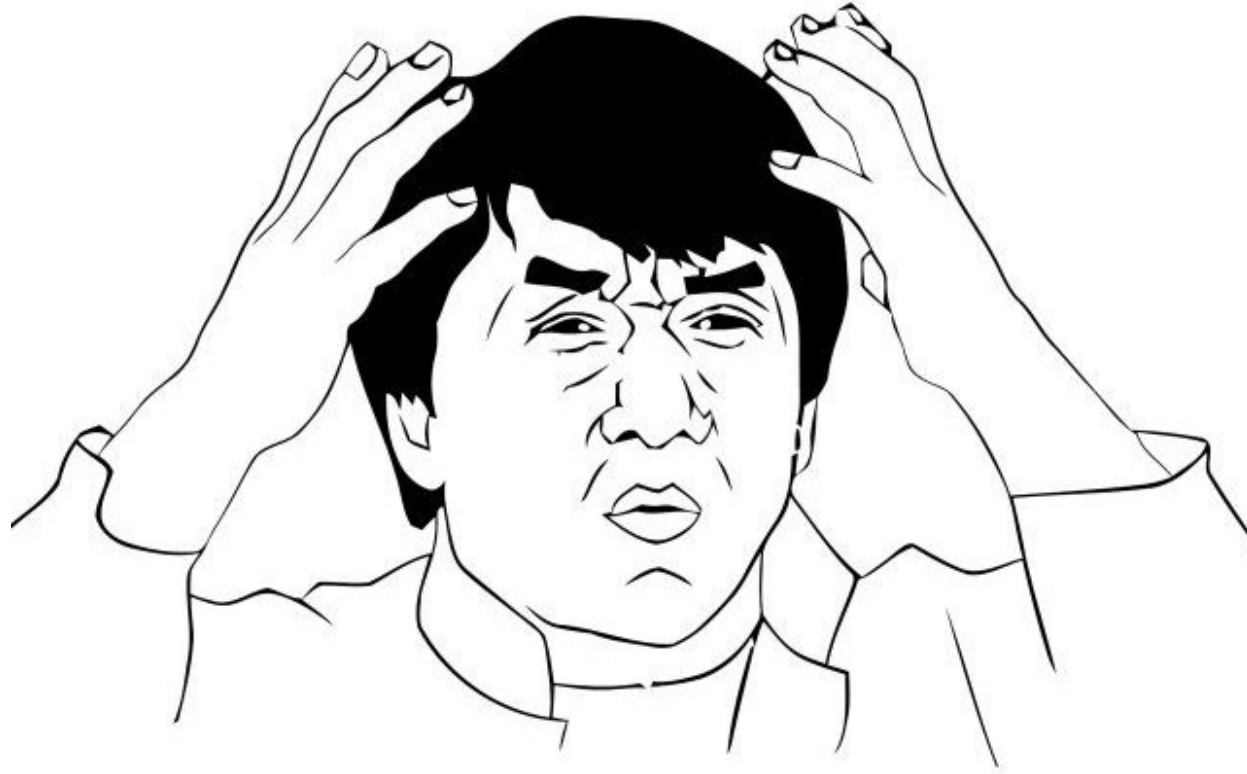
### Optimized Rendering

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





“The Snapdragon Spaces XR Developer Platform reduces developer friction by providing a uniform set of augmented reality features independent of device manufacturers. This allows developers to seamlessly blend the lines between our physical and digital realities and transform the world around us in ways limited only by our imaginations.”

# What's Next?



# **An Outlook into the Future**

# What's Relevant in Egovision? A top-down approach



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 helps with cooking
- Sam is impressed by how fun it is to cook with his 3D friend
- Toaster reminder
- 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)

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

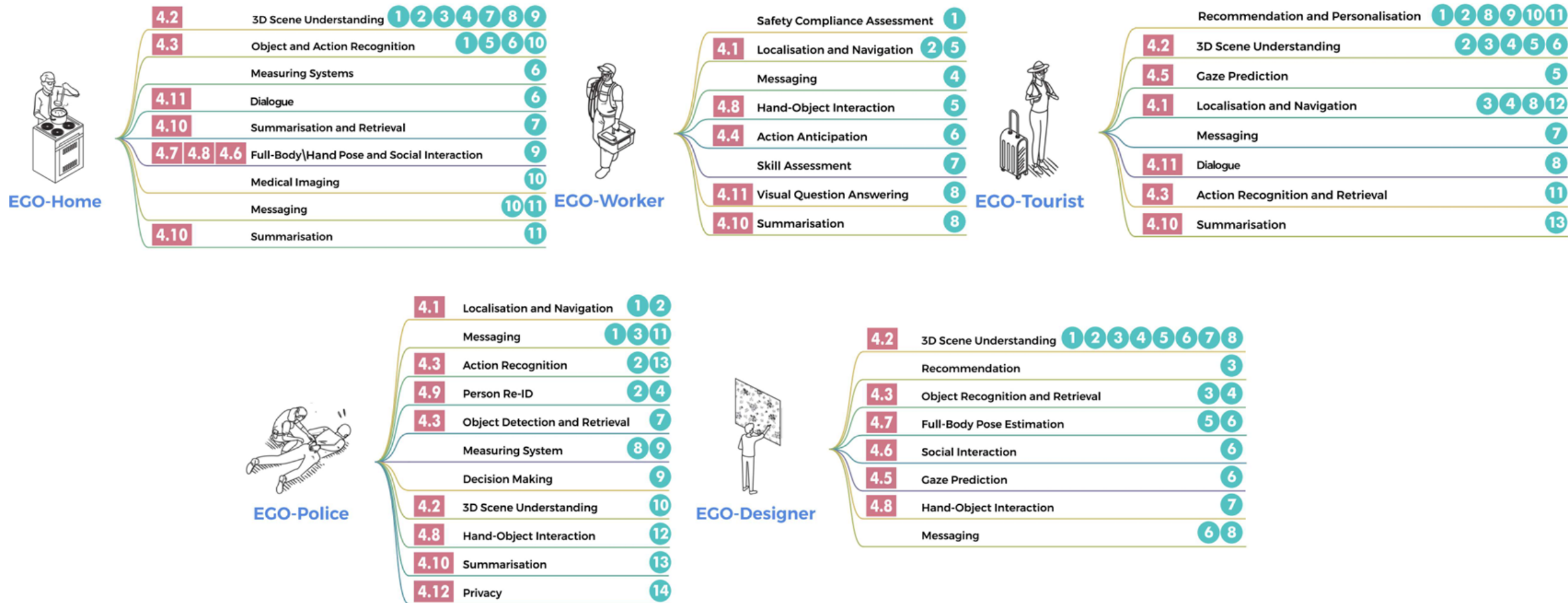
EgoAI assists make-up artists with advanced 3D modelling techniques to project guidelines on the actor's face while applying make-up

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 parts 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 ( <a href="#">Ragusa et al 2023b</a> )			✓	✓	✓			✓					
ADL ( <a href="#">Pirsiavash and Ramanan 2012</a> )			✓	✓						✓			
HOI4D ( <a href="#">Liu et al 2022c</a> )								✓					
EGTEA Gaze+ ( <a href="#">Li et al 2021a</a> )			✓	✓	✓			✓					
UTE ( <a href="#">Lee et al 2012</a> )								✓		✓			
EGO-CH ( <a href="#">Ragusa et al 2020a</a> )	✓												
FPSI ( <a href="#">Fathi et al 2012a</a> )							✓			✓		✓	
KrishnaCam ( <a href="#">Singh et al 2016a</a> )				✓									
EPIC-KITCHENS-100 ( <a href="#">Damen et al 2022</a> )		✓	✓	✓				✓			✓	✓	
Assembly101 ( <a href="#">Sener et al 2022</a> )			✓					✓					
Ego4D ( <a href="#">Grauman et al 2022</a> )			✓	✓	✓	✓		✓		✓	✓		



It's an exciting time for wearable devices & egocentric vision!

Hardware is increasingly available as big tech gets interested.



Large datasets and pre-defined challenges can help get started to explore the field







**VISIGRAPP 2024**

19<sup>th</sup> International Joint Conference on Computer Vision, Imaging  
and Computer Graphics Theory and Applications

Rome, Italy 27 - 29 February, 2024

GRAPP HUCAPP IVAPP VISAPP



Università  
di Catania

**NEXT VISION**

Spin-off of the University of Catania



# THANK YOU!

## First Person (Egocentric) Vision: History and Applications

### Francesco Ragusa

First Person Vision@Image Processing Laboratory - <http://iplab.dmi.unict.it/fpv>

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

Department of Mathematics and Computer Science - University of Catania

[francesco.ragusa@unict.it](mailto:francesco.ragusa@unict.it) - <https://francescoragusa.github.io/>



1) Part I: History and motivations [09.00 - 10.30]

- a) Agenda of the tutorial;
- b) Definitions, motivations, history and research trends of First Person (egocentric) Vision;
- c) Seminal works in First Person (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;

**Coffee Break [10.30 – 10.45]**

**Keynote presentation: Gerhard Rigoll [10.45 – 12.00]**

1) Part II: Fundamental tasks for First Person Vision systems [12.00 – 13.00]

- a) Localization;
- b) Hand/Object Detection;
- c) Action/Activity Recognition;
- d) Egocentric Human-Object Interaction;
- e) Anticipation;
- f) Industrial Applications;
- g) Conclusion.