# Egocentric Vision:
# Emerging Trends and Human-Centric Applications

# 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/

**LIVE Group @ UNICT**

# The 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/iciap2025

1) **Part I: History and motivations [14.30 - 15.30]**

   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?**

   i) **Industrial Applications**

1) Part I: History and motivations [14.30 -  15.30]

   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;
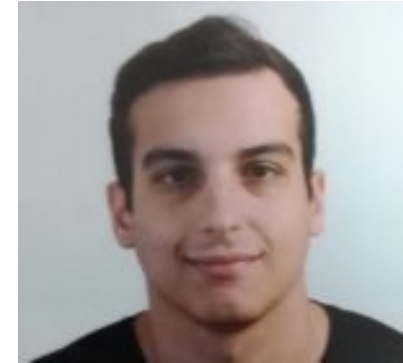
   h) What's next?

   i) Industrial Applications



**Coffee Break [15.30 – 15.50]**



**Coffee Breaks** are organized autonomously by each workshop to best suit their schedule and format. Participants will receive coupons to enjoy coffee and snacks at the campus café
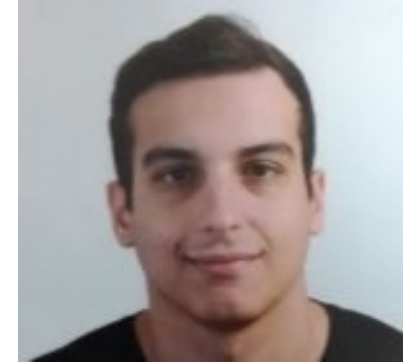
2) **Part II: Hand-Object Interactions in Egocentric Vision [15.50 – 16.50]**

    a)   **Introduction to Hand-Object Interactions Detection**

    b)  **Datasets and Benchmarks for Hand-Object Interactions in Egocentric Vision**

    c)  **Models and Architectures for Hand-Object Interactions Detection**

    d)  **Open Challenges**

2) Part II: Hand-Object Interactions in Egocentric Vision [15.50 – 16.50]

   a) Introduction to Hand-Object Interactions Detection

   b) Datasets and Benchmarks for Hand-Object Interactions in Egocentric Vision

   c) Models and Architectures for Hand-Object Interactions Detection

   d) Open Challenges

**Short Break [16.50 – 17.00]**

3) **Part III: Gaze Understanding and Visual-Language Benchmarks [17.00 – 18.00]**

    a) **Gaze Signal Fundamentals**

    b) **Gaze-Based Dataset**

    c) **Gaze signal in computer vision**

    d) **Building procedural assistant with VLLM**

    e) **Open Challenges and Future Directions**

# Part I

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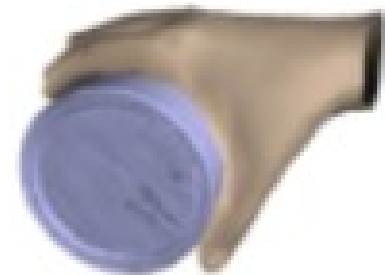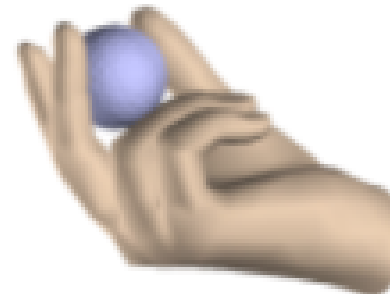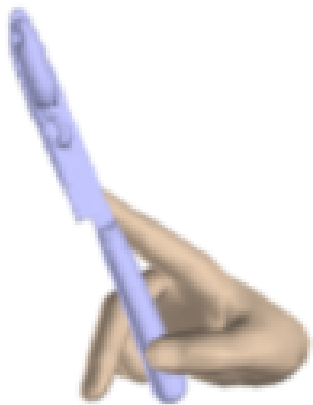
Perception is the process of **receiving**, **organize** and **interpret** information in order to give meaning to the surrounding world.

I'm in the kitchen!

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



E. Corona, A. Pumarola, G. Aleny, M. N. Francesc, R. Gregory. GanHand: Predicting Human Grasp Affordances in Multi-Object Scenes, CVPR, 2020.
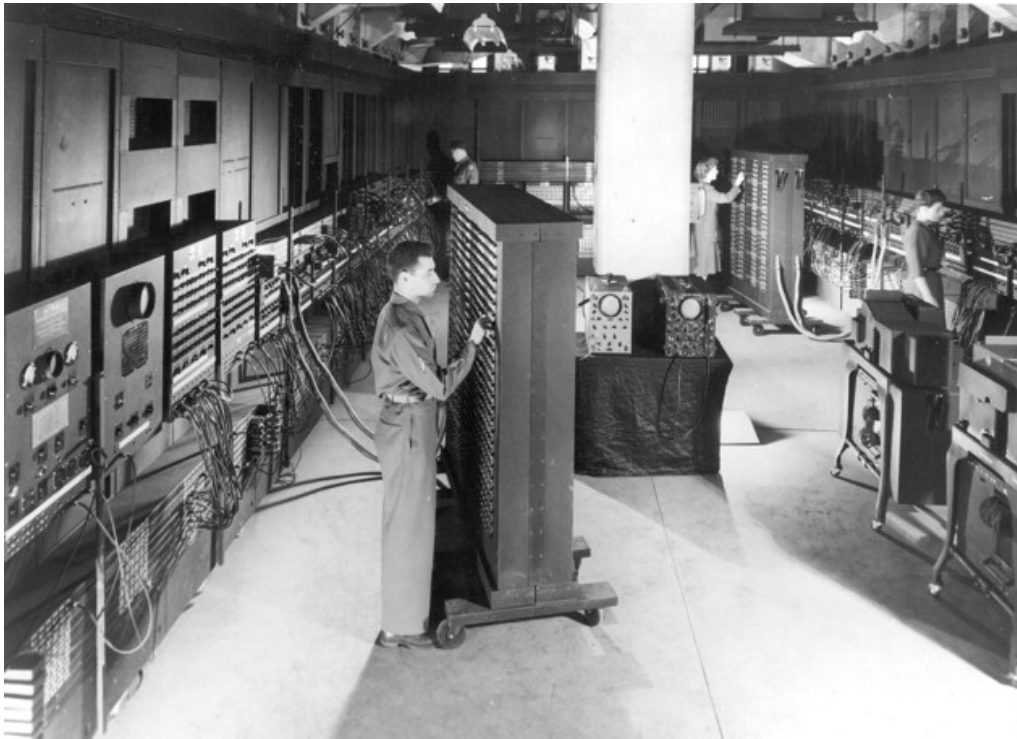
Perception is the process of **receiving**, **organize** and **interpret** information in order to give meaning to the surrounding world.

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)
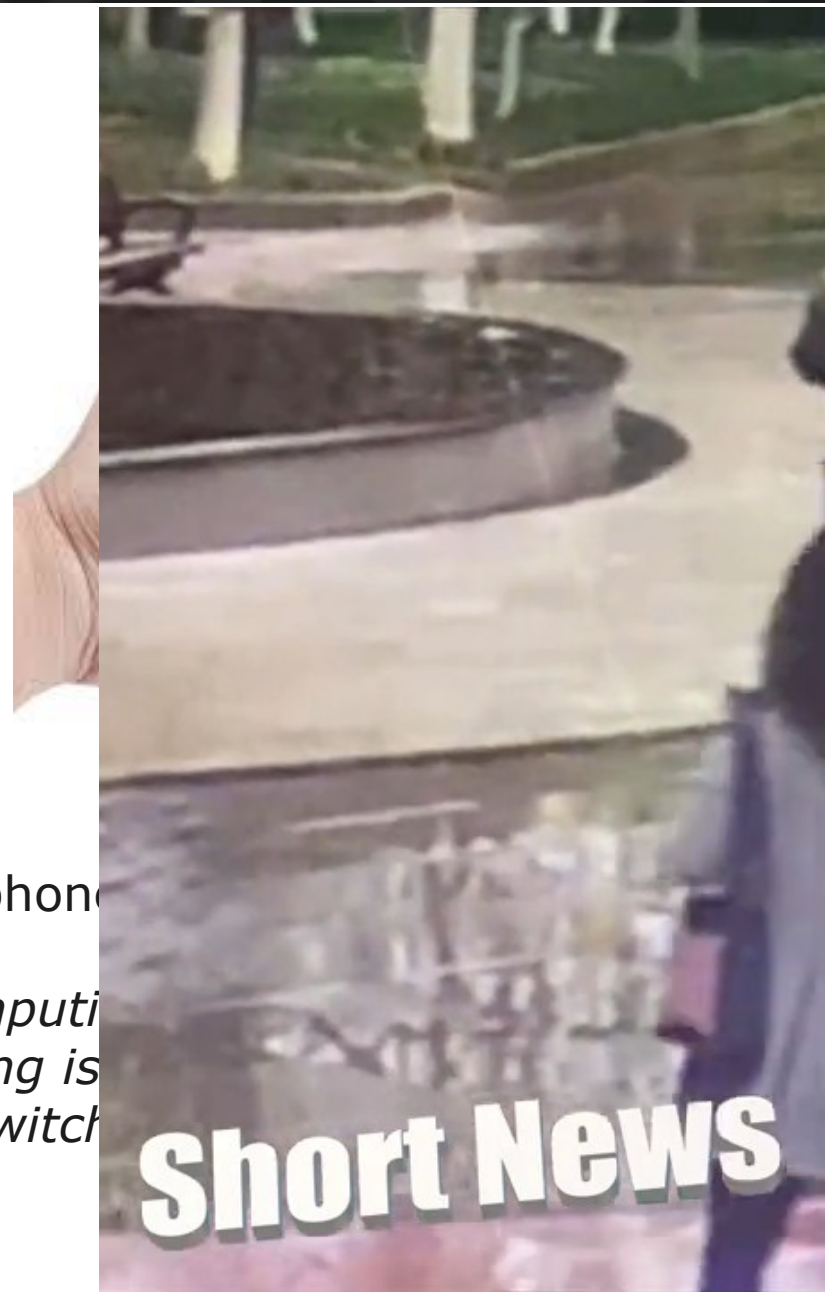
*Centralized, inaccessible, institutional*



Personal Computer Era (1980s–1990s)

*Desktop computing enters homes and offices*

Università di Catania



Laptop Era (1990s–2000s)

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

Smartphone

*Computi... Computing is... forces to switch...*

Short News

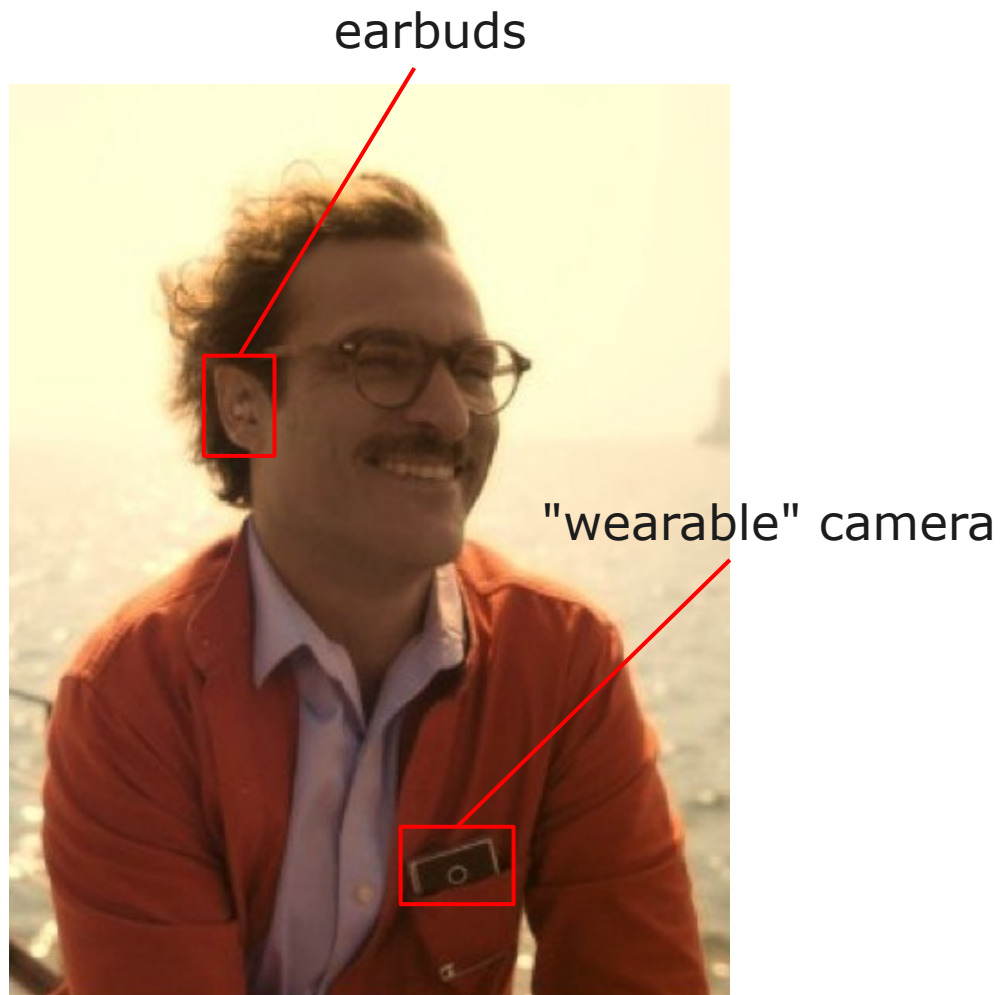**Smartphone Era (2007–present)**

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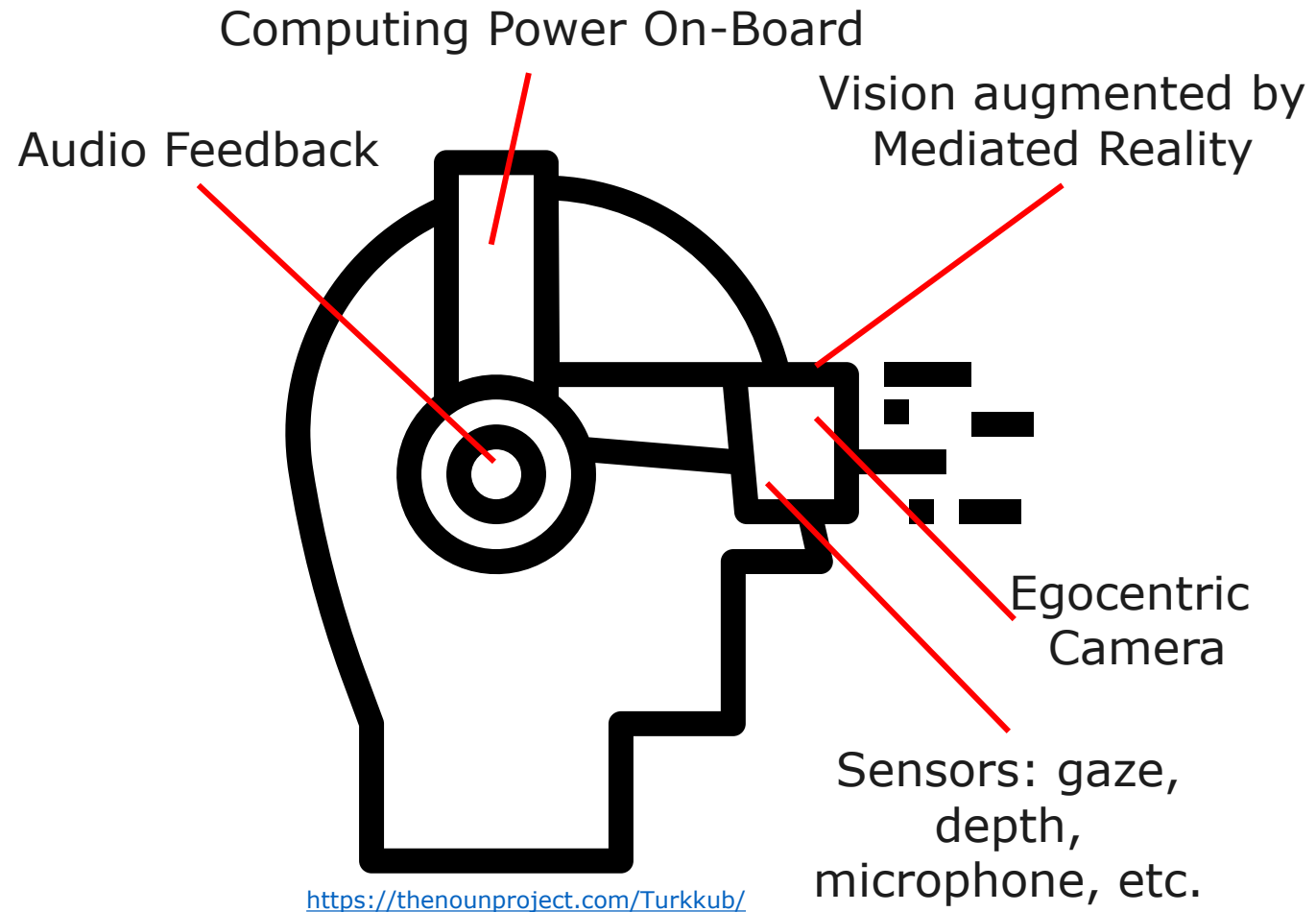forces to switch between the digital and
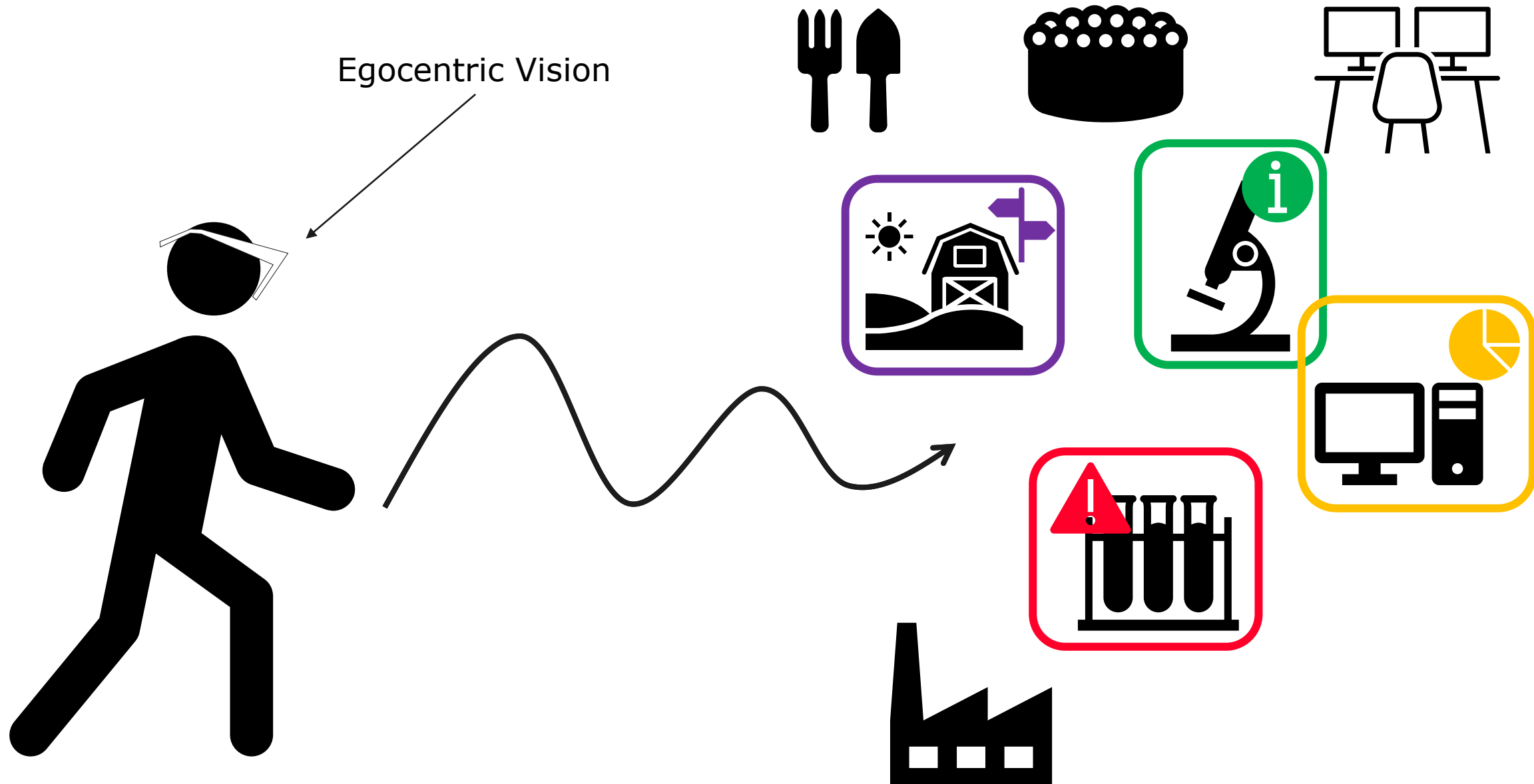real world*

**Smartglasses Era (Now and Future)**

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

earbuds

"wearable" camera

"her" 2013 movie

Computing Power On-Board

Audio Feedback

Vision augmented by Mediated Reality

Egocentric Camera

Sensors: gaze, depth, microphone, etc.

https://thenounproject.com/Turkkub/

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

Egocentric Vision

# (Egocentric) Computer Vision is Fundamental!

## Exocentric

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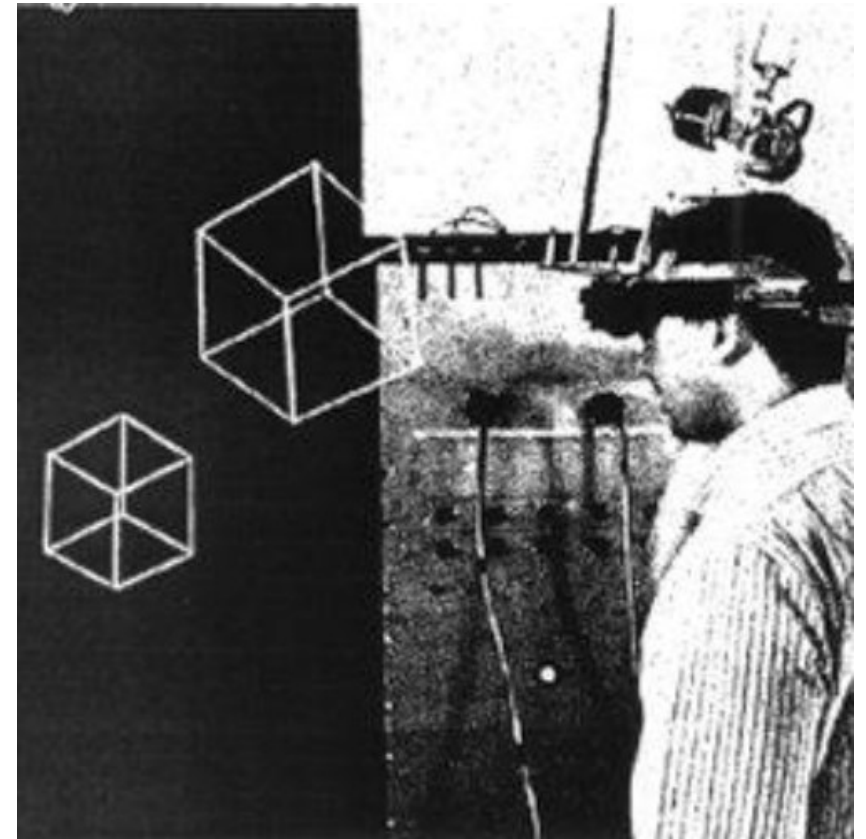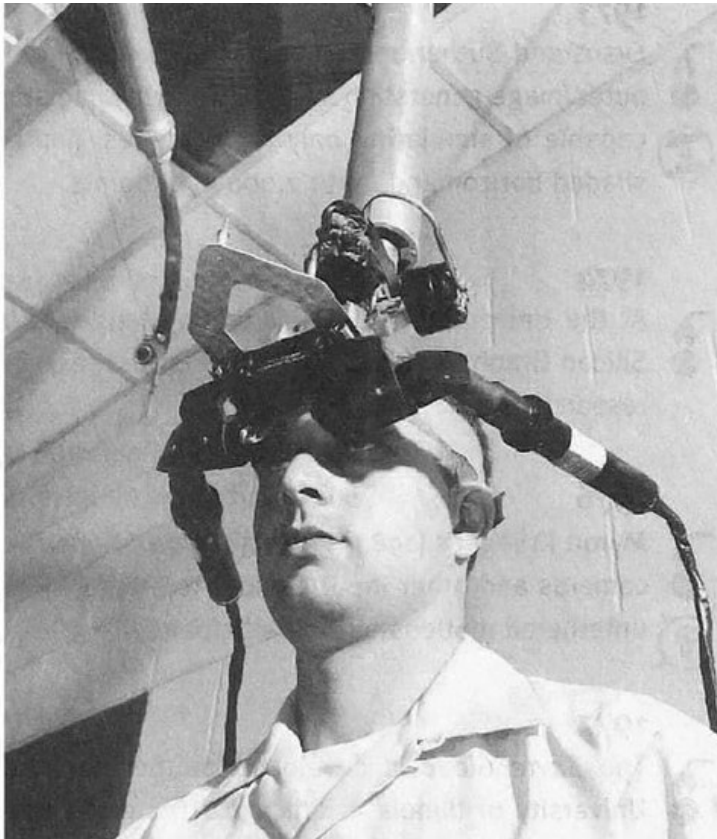× Doesn't always see everything
× Not really portable

## 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 <u>stereoscopic</u> 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.



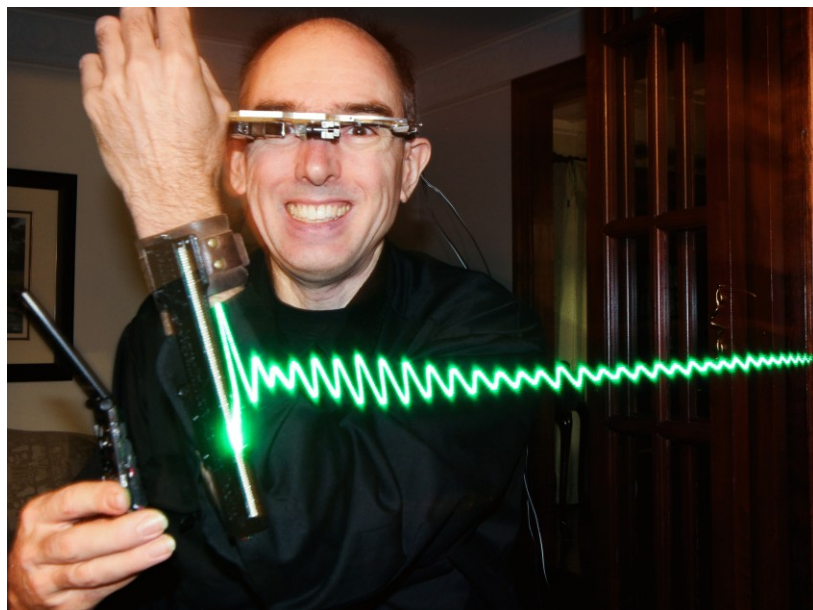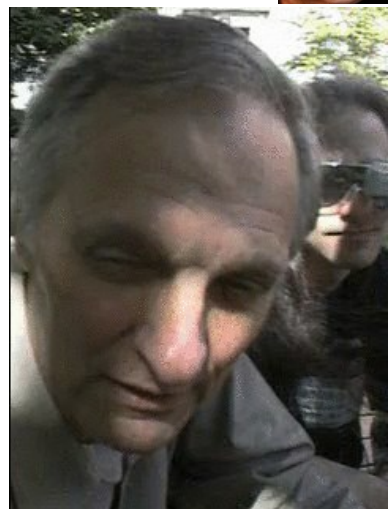(a) **1980**  (b) **Mid 1980s**  (c) **Early 1990s**  (d) **Mid 1990s**  (e) **Late 1990s**

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



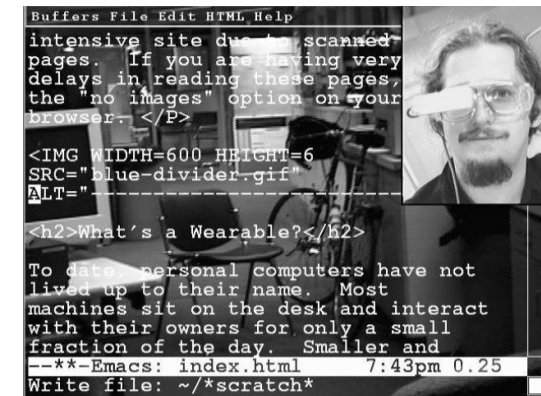Steve Mann. "Compositing multiple pictures of the same scene." *Proc. IS&T Annual Meeting, 1993.*
Steve Mann, "Wearable computing: a first step toward personal imaging," in *Computer*, vol. 30, no. 2, pp. 25-32, Feb. 1997.

Università di Catania

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)



Water me

CYBORG Central
Meeting tonight
4 AM

```
Buffers File Edit HTML Help
intensive site due to scanned
pages. If you are having very
delays in reading these pages,
the "no images" option on your
browser. </P>

<IMG WIDTH=600 HEIGHT=6
SRC="blue-divider.gif"
ALT="------------------

<h2>What's a Wearable?</h2>

To date, personal computers have not
lived up to their name. Most
machines sit on the desk and interact
with their owners for only a small
fraction of the day. Smaller and
--**-Emacs: index.html    7:43pm 0.25
Write file: ~/*scratch*
```

1998

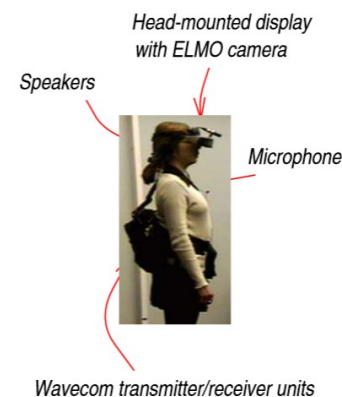## Visual Contextual Awareness in Wearable Computing

Thad Starner          Bernt Schiele          Alex Pentland

Media Laboratory, Massachusetts Institute of Technology
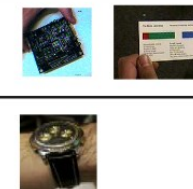
(location and task recognition)

1999

## 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)

Head-mounted display with ELMO camera

Speakers

Microphone

Wavecom transmitter/receiver units
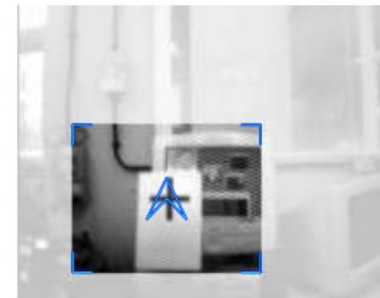
VISUAL TRIGGER

ASSOCIATED SEQUENCE

? 

GARBAGE NO PLAY-BACK

## Wearable Visual Robots

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

(active vision)

$$P(C_t \mid v_{1:t}^G)$$

## 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
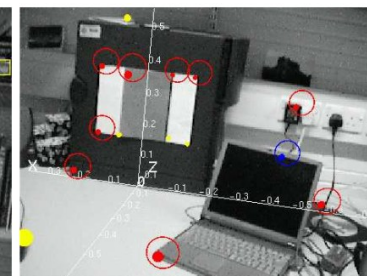Lincoln Labs
Lexington, MA 02420

(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)

## Wearable Hand *Activity* Recognition for Event Summarization

W.W. Mayol

D.W. Murray

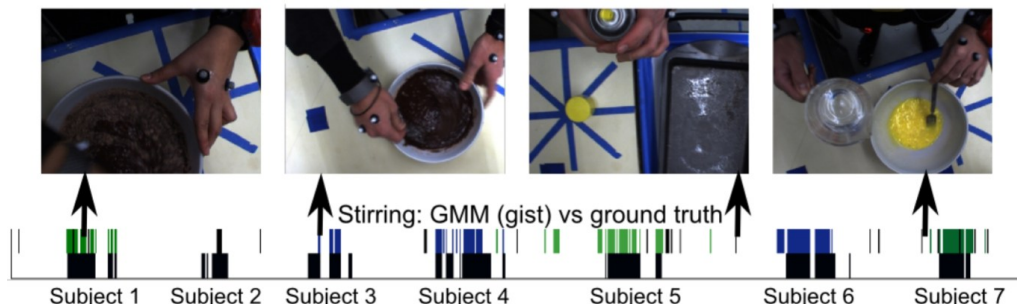Department of Computer Science

Department of Engineering Science

University of Bristol

University of Oxford

(hand activity recognition)

**2005**

## Temporal Segmentation and Activity Classification from First-person Sensing

Ekaterina H. Spriggs, Fernando De La Torre, Martial Hebert

Carnegie Mellon University.

(activity classification)

Stirring: GMM (gist) vs ground truth

Subject 1  Subject 2  Subject 3  Subject 4  Subject 5  Subject 6  Subject 7

**2009**

## Figure-Ground Segmentation Improves Handled Object Recognition in Egocentric Video

Xiaofeng Ren

Chunhui Gu

Intel Labs Seattle

University of California at Berkeley

1100 NE 45th Street, Seattle, WA 98105

Berkeley, CA 94720

(handheld object recognition)

**2010**

# A COMMON HARDWARE PLATFORM WAS MISSING!

# "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
7 JJ Thomson Ave, Cambridge, UK, CB3 0FB
*Behavioural & Clinical Neuroscience Institute
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
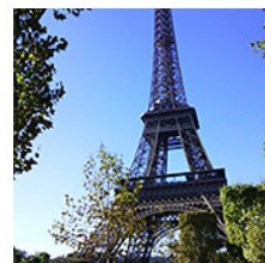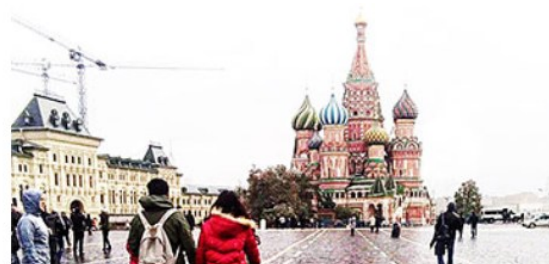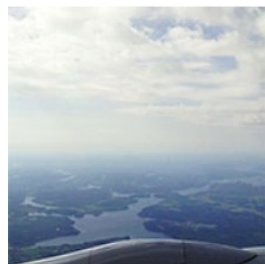Centre for Digital Video Processing & Adaptive Information Cluster,
Dublin City University

(lifelogging, multimedia retrieval)

**2008**

http://getnarrative.com/

## 2016

**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)

## 2017

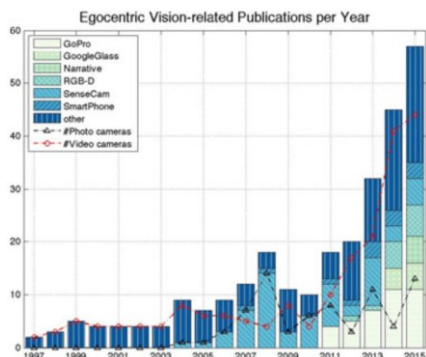Day's Lifelog:

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)

## 2017

Egocentric Vision-related Publications per Year

Event Segmentation

Working    Ride a bike    Dinner time

**Toward Storytelling From Visual Lifelogging: An Overview**

Marc Bolaños, Mariella Dimiccoli, and Petia Radeva

(lifelogging, survey)

## different wearing modalities

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



head-mounted

chest-mounted

wrist-mounted

helmet-mounted

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

Kris M. Kitani
UEC Tokyo
Tokyo, Japan

Takahiro Okabe, Yoichi Sato
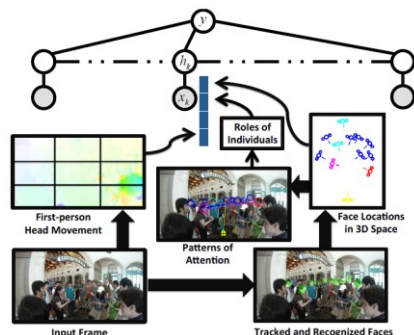University of Tokyo
Tokyo, Japan

Akihiro Sugimoto
National Institute of Informatics
Tokyo, Japan

(unsupervised action recognition, video indexing)
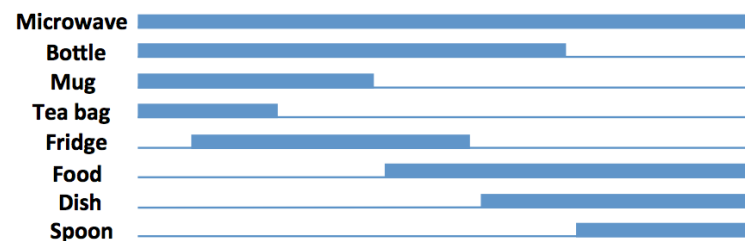
## Social Interactions: A First-Person Perspective

Alireza Fathi[1], Jessica K. Hodgins[2,3], James M. Rehg[1]

(detection and recognition of social interactions)

## Story-Driven Summarization for Egocentric Video

Zheng Lu and Kristen Grauman
University of Texas at Austin

(egocentric video sumarization)

**Our method**

## 2014

### Temporal Segmentation of Egocentric Videos

Yair Poleg          Chetan Arora*          Shmuel Peleg

(egocentric video indexing)



(a) Car   (b) Bus   (c) Walking   (d) Sitting   (e) Wheels   (f) Standing   (g) Static

## 2021

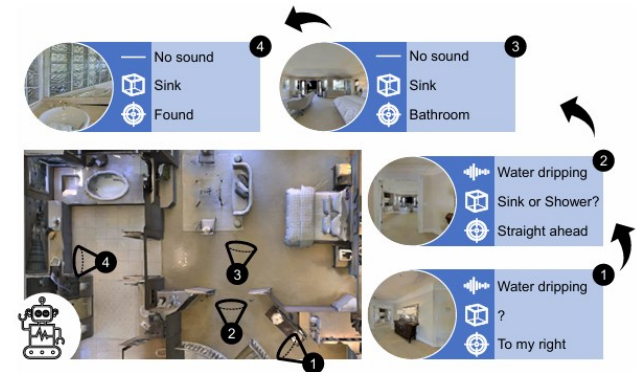### Semantic Audio-Visual Navigation

Changan Chen[1,2]    Ziad Al-Halah[1]    Kristen Grauman[1,2]

[1]UT Austin    [2]Facebook AI Research



## 2025

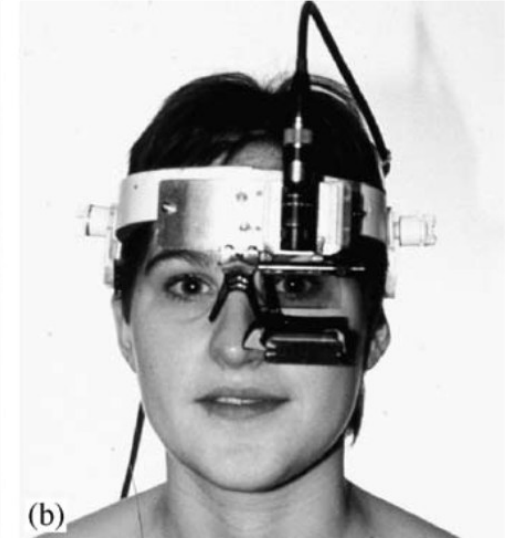### EgoLife: Towards Egocentric Life Assistant
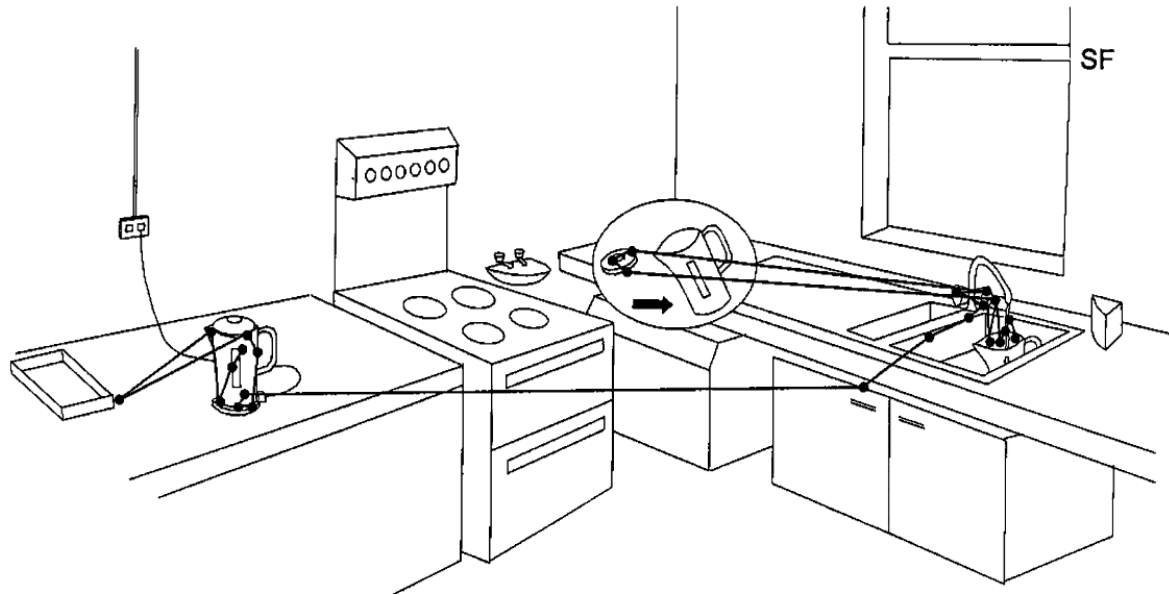
The EgoLife Team

https://egolife-ai.github.io/

(video understanding, egocentric assistant)

Eye movements and the control of actions in everyday life

Michael F. Land



Prototype by Land (1993)

**Gaze is important in Egocentric Vision!**
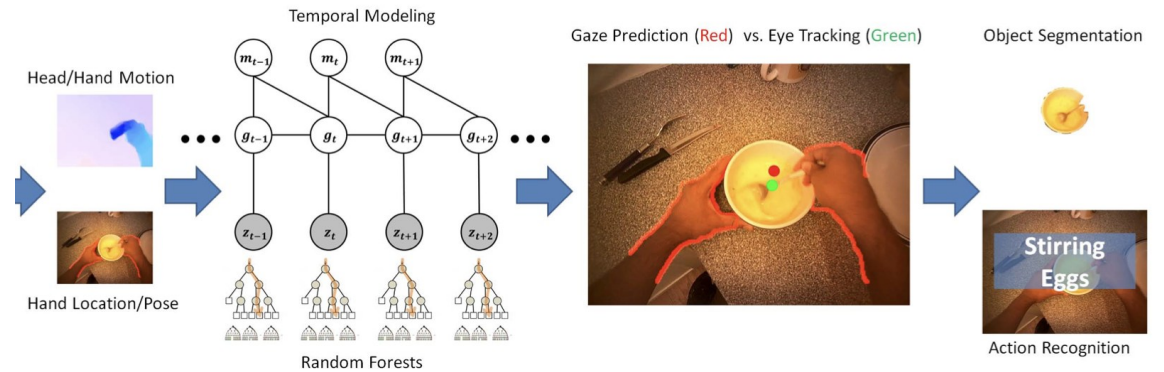
Tobii Pro Glasses 2 (2014)     Microsoft HoloLens 2 (2016)
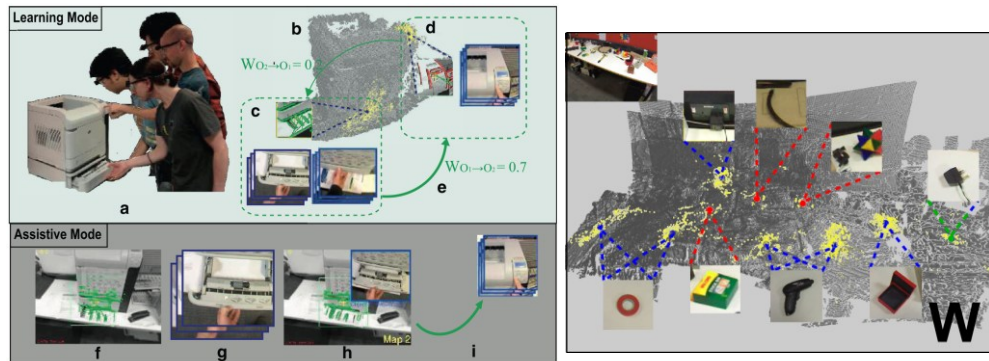
Mobile Eye-XG (2013)     Pupil Eye Trackers (2014 - )

## Learning to Predict Gaze in Egocentric Video

Yin Li, Alireza Fathi, James M. Rehg

(gaze prediciton, action recognition)

**2012**

You-Do, I-Learn: Egocentric unsupervised discovery of objects and their modes of interaction towards video-based guidance
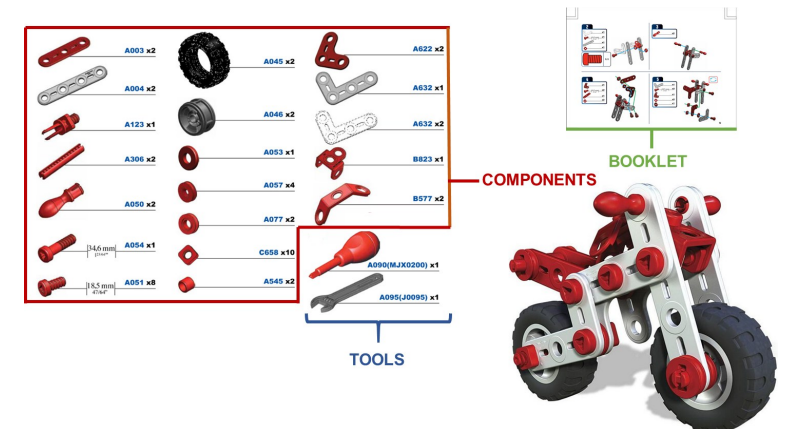
Dima Damen*, Teesid Leelasawassuk, Walterio Mayol-Cuevas

(object usage discovery, assistance)

**2016**

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

Francesco Ragusa *, Antonino Furnari, Giovanni Maria Farinella

(gaze prediciton, procedural video)

**2023**

**Health, assistive technologies**

https://www.orcam.com/

https://www.orcam.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.
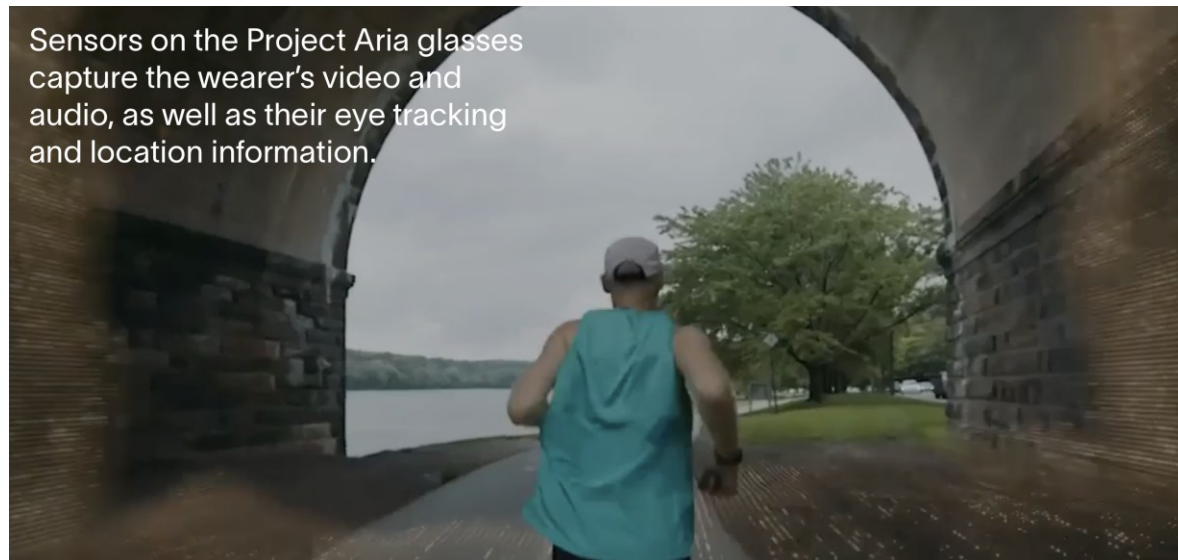
$5,199

https://www.microsoft.com/en-us/hololens/buy
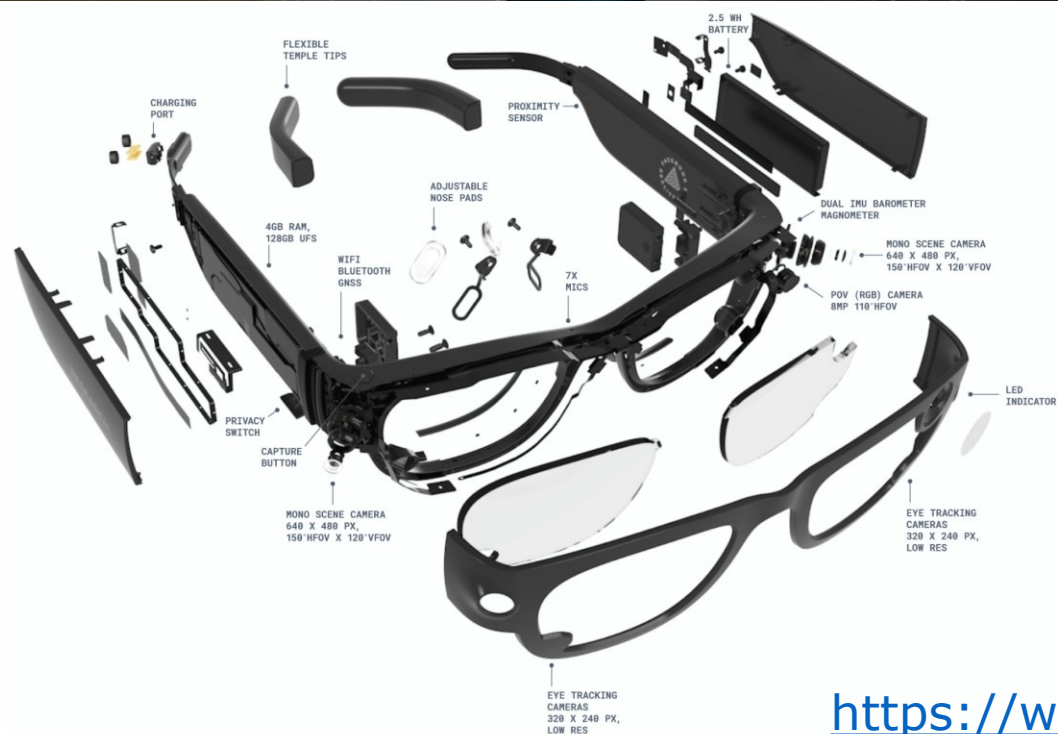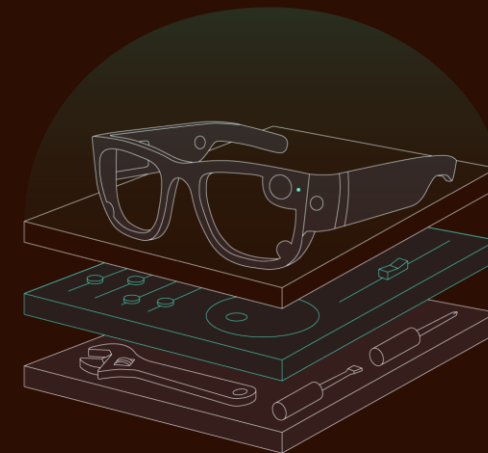
https://www.magicleap.com/magic-leap-2

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

https://www.xreal.com/

https://www.apple.com/apple-vision-pro/

# Too Many Devices?

towards standardization...

Unified API supported by many AR and VR devices



https://www.khronos.org/openxr/

**Workshop on Egocentric (First Person) Vision**

# Digital Information

**CMU**
(0.2M frames – 2009)

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

**GTEA Gaze+**
(0.4M frames – 2012)

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

**ADL**
(1.0M frames – 2012)

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

**Charades-ego**
(2.3M frames – 2018)

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

**EGTEA Gaze+**
(2.4M frames – 2018)

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

# The EPIC series


EPIC-Kitchens 55


EPIC-Kitchens 100


EPIC-Kitchens VISOR


EPIC-SOUNDS


EPIC-FIELDS


HD-EPIC

# EPIC-KITCHENS
# TEAM

**University of BRISTOL**

**UNIVERSITY OF TORONTO**

**UNIVERSITÀ degli STUDI di CATANIA** · 1434 ·

**Dima Damen**
**Principal Investigator**
University of Bristol
United Kingom

**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

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

https://epic-kitchens.github.io/

| | 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) |

https://epic-kitchens.github.io/

- 272K manual sparse masks for hands and active objects;

- Hand-object contact relations;

- 1477 unique entities;

- 22 categories.

https://epic-kitchens.github.io/VISOR

spray

44 classes...

- 74.8K categorised audio segments;

- Material-based collision sounds;

- Repetitive sounds;

- 44 classes.

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

and goes 3D

frame_000000307

- 19M registered frames;

- Camera poses;

- 3D reconstruction;

- Paired with VISOR annotations.

https://epic-kitchens.github.io/epic-fields

**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 Recognition](#)
- [Action Recognition](#)
- [Action Detection](#)
- [UDA for Action Recognition](#)
- [Multi-Instance Retrieval](#)

https://epic-kitchens.github.io/

# Can We Scale?

# EGO 4D

## Consortium

Carnegie Mellon University

Università di Catania

NUS National University of Singapore

King Abdullah University of Science and Technology

東京大学 THE UNIVERSITY OF TOKYO

University of BRISTOL

INDIANA UNIVERSITY BLOOMINGTON

UNIVERSITY OF MINNESOTA

MiT

Penn UNIVERSITY of PENNSYLVANIA

Georgia Institute of Technology

INTERNATIONAL INSTITUTE OF INFORMATION TECHNOLOGY HYDERABAD

Carnegie Mellon University Africa

Universidad de los Andes Colombia

FACEBOOK AI

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

### 84 authors

Kristen Grauman[1,2], Andrew Westbury[1], 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 Bansal[17], Dhruv Batra[1], Vincent Cartillier[1,6], Sean Crane[7], Tien Do[3], Morrie Doulaty[13], Akshay Erapalli[13], Christoph Feichtenhofer[1], Adriano Fragomeni[9], Qichen Fu[7], Christian Fuegen[13], Abrham Gebreselasie[12], Cristina González[14], James Hillis[5], Xuhua Huang[7], Yifei Huang[15], Wenqi Jia[6], Weslie Khoo[16], Jachym Kolar[13], Satwik Kottur[13], Anurag Kumar[5], Federico Landini[13], Chao Li[5], Zhenqiang Li[15], Karttikeya Mangalam[1,8], Raghava Modhugu[17], Jonathan Munro[9], Tullie Murrell[1], Takumi Nishiyasu[15], Will Price[9], Paola Ruiz Puentes[14], Merey Ramazanova[10], Leda Sari[5], Kiran Somasundaram[5], Audrey Southerland[6], Yusuke Sugano[15], Ruijie Tao[11], Minh Vo[5], Yuchen Wang[16], Xindi Wu[7], Takuma Yagi[15], Yunyi Zhu[11], Pablo Arbeláez†[14], David Crandall†[16], Dima Damen†[9], Giovanni Maria Farinella†[4], Bernard Ghanem†[10], Vamsi Krishna Ithapu†[5], C. V. Jawahar†[17], Hanbyul Joo†[1], Kris Kitani†[7], Haizhou Li†[11], Richard Newcombe†[5], Aude Oliva†[18], Hyun Soo Park†[3], James M. Rehg†[6], Yoichi Sato†[15], Jianbo Shi†[19], Mike Zheng Shou†[11], Antonio Torralba†[18], Lorenzo Torresani†[1,20], Mingfei Yan†[5], Jitendra Malik[1,8]

[1]Facebook AI Research (FAIR), [2]University of Texas at Austin, [3]University of Minnesota, [4]University of Catania, [5]Facebook Reality Labs, [6]Georgia Tech, [7]Carnegie Mellon University, [8]UC Berkeley, [9]University of Bristol, [10]King Abdullah University of Science and Technology, [11]National University of Singapore, [12]Carnegie Mellon University Africa, [13]Facebook, [14]Universidad de los Andes, [15]University of Tokyo, [16]Indiana University, [17]International Institute of Information Technology, Hyderabad, [18]MIT, [19]University of Pennsylvania, [20]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**

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/

EPIC-Kitchens-100

Animation by Michael Wray – https://mwray.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**

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**

**First Joint Egocentric Vision (EgoVis) Workshop**

**Held in Conjunction with CVPR 2024**

17 June 2024 - Seattle, USA

Room: Summit 428

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

Kristen Grauman[1,2], Andrew Westbury[1], Lorenzo Torresani[1], 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], Abrham 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 Castillo[12], Changan Chen[2], Xinzhu Fu[11], Ryosuke Furuta[17], Cristina González[12], Prince Gupta[5], Jiabo Hu[18], Yifei Huang[17], Yiming Huang[16], Weslie Khoo[19], Anush Kumar[10], Robert Kuo[18], Sach Lakhavani[5], Miao Liu[18], Mi Luo[2], Zhengyi Luo[3], Brighid Meredith[18], Austin Miller[18], Oluwatumininu Oguntola[14], Xiaqing Pan[5], Penny Peng[18], Shraman Pramanick[20], Merey Ramazanova[21], Fiona Ryan[22], Wei Shan[14], Kiran Somasundaram[5], Chenan Song[11], Audrey Southerland[22], Masatoshi Tateno[17], Huiyu Wang[1], Yuchen Wang[19], Takuma Yagi[17], Mingfei Yan[5], Xitong Yang[1], Zecheng Yu[17], Shengxin Cindy Zha[18], Chen Zhao[21], Ziwei Zhao[19], Zhifan Zhu[6], Jeff Zhuo[14], 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**

# Second Joint Egocentric Vision (EgoVis) Workshop
## Held in Conjunction with CVPR 2025
## 11 or 12 June 2025 - Nashville, USA

Ego-Exo4D          Ego4D          EPIC-Kitchens

Università di Catania



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

F. Ragusa[1,3], A. Furnari[1], S. Livatino[2], G. M. Farinella[1]

[1]IPLab, Department of Mathematics and Computer Science - University of Catania, IT
[2]University of Hertfordshire, Hatfield, Hertfordshire, U.K.
[3]Xenia 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 Sener[1]  Dibyadip Chatterjee[2]  Daniel Shelepov[1]  Kun He[1]
Dipika Singhania[2]  Robert Wang[1]  Angela Yao[2]
[1]Reality Labs at Meta
[2]National 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. Schoonbeek[1], Tim Houben[1], Hans Onvlee[2], Peter H.N. de With[1], Fons van der Sommen[1],
[1]Eindhoven University of Technology, [2]ASML Research
Published in: WACV 2024

📄 Paper   𝕏 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

**COMPONENTS**

**BOOKLET**

**TOOLS**

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

F. Ragusa, A. Furnari, G. M. Farinella. MECCANO: A Multimodal Egocentric Dataset for Humans Behavior Understanding in the Industrial-like Domain. Computer Vision and Image Understanding (CVIU), 2023 (https://arxiv.org/abs/2209.08691).

GoPro Hero 4

Real Sense SR300

Pupils

**Verbs Classes**



**8857 video segments**

**1401 overlap segments (15.82%)**



F. Ragusa, A. Furnari, G. M. Farinella. MECCANO: A Multimodal Egocentric Dataset for Humans Behavior Understanding in the Industrial-like Domain. Computer Vision and Image Understanding (CVIU), 2023 (https://arxiv.org/abs/2209.08691).

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**

F. Ragusa, A. Furnari, G. M. Farinella. MECCANO: A Multimodal Egocentric Dataset for Humans Behavior Understanding in the Industrial-like Domain. Computer Vision and Image Understanding (CVIU), 2023 (https://arxiv.org/abs/2209.08691).

**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

F. Ragusa, A. Furnari, G. M. Farinella. MECCANO: A Multimodal Egocentric Dataset for Humans Behavior Understanding in the Industrial-like Domain. Computer Vision and Image Understanding (CVIU), 2023 (https://arxiv.org/abs/2209.08691).

## Egocentric Human-Object Interaction

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

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

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



**<take, screwdriver>**



**<screw, {screwdriver, screw, partial_model}>**

(«take, bolt»)

3 s before

0.2 s ... 0.2 s 0.2 s

past frames

start frame



| 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** |

Training: 5057, 3848
Validation: 977, 706
Test: 2823, 2102

■ #Interactions  ■ #Interactions with Past

# 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



active object — gray perforated bar — instruction booklet

# 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.



$V_1$
$V_2$
$\vdots$
$V_n$

Spread the dough → Apply sauce → Grate and add cheese → Cut and add pepperoni → Put pizza in the oven

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

4) MECCANO
5) EPIC-Tent

B. Siddhant, A. Chetan, C. V. Jawahar, My View is the Best View: Procedure Learning from Egocentric Videos. In European Conference on Computer Vision (ECCV), 2022.

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                    Hight-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

TTC: 1.20s   Δ = 0.40   TTC: 0.80s   TTC: 0.40s   first-contact high voltage board

Past Frames

Interaction Frame

**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

......

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).

Untrimmed temporal detection of human-object interactions

Egocentric human-object interaction detection

Short-term object interaction anticipation

Natural language understanding of intents and entities

# ENIGMA-360 (Extension)

Exocentric

Egocentric

35 subjects

7 tasks

- **Temporal Action Segmentation**
- **Keystep Recognition**
- **Hand Object Interaction Segmentation**

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)

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)

Ego-EXTRA: Egocentric dataset of EXpert-TRAinee assistance

# Multiple-Choice Question Answering



Video Clip

"Do I need to worry that the wheel might fall?"

Trainee's question

Input

MiniGPT4-video → **C**

Video-LLama → **A**

Expert (GT) → **A**

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."

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)

# What's Next?

# An Outlook into the Future

**A lot of data!**



Imagine the Future

↓

Write Stories in Different Scenarios

↓

Extract Important Tasks from the Stories

↓

Go in-depth with Tasks and Datasets

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

**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

\*: 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

### 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   PDF

*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: [all ▾]  Author: [everybody ▾]  Visible To: [all readers ▾]  Hidden From: [nobody ▾]    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. MPCHAT: Towards Multimodal Persona-Grounded Conversation, ACL 2023 (https://aclanthology.org/2023.acl-long.189/)

Add  Comment

### [−] Related work on egocentric full-body pose estimation

*Jiaxi 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

(1)

This way the tomato will cook evenly

(2)

A 3D projection of Remy helps Sam with cooking

(3) Audible 3D projection

Sam is impressed by how fun it is to cook with his 3D friend

(4)

(5) Toaster reminder

(6) EgoAI recommends some more spice

(8) Waves hitting the shore look and sound natural

(7) 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 (9)

(11)

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 (10)

EgoAI proposes a short clip from his day, but Sam decides not to share it

# EGO-WORKER

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

(2) EgoAI localises Marco and provides ruote instructions to reach his workstation for the day

Where should I go today in the factory?

(3) 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 (4)

Since the measuring device is a new brand, EgoAI guides Marco through the basic functionality and tools (5)

(6) EgoAI detects a risk and turns off the IoT electrical socket while promptly alerting Marco

(7) For the rest of the day, EgoAI validates Marco's work making sure that the procedures are properly and safely completed

(8) By the end of the day, EgoAI checks Marco's feedback for improving future sessions

# EGO-TOURIST

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

(2) EgoAI suggests an half-day visit to the Egyptian museum

(3) EgoAI activates the 3D projection of Cleopatra to guide and interact with Claire

Claire feels transported to ancient Egypt (6)

(7) Claire asks Cleopatra for a good place for a pizza

(5) Cleopatra leads Claire through the artworks and proposes her the most suited path

Claire observes virtual elements being added to the scene, which bring the artwork to life

(4)

(8) 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 (9)

EgoAI suggests Claire a proper Italian coffee at a nearby café, sided by a slice of bunet, Turin-based dessert (10)

(11) EgoAI offers a egocentric view from the chef who prepared her that delicacy

(13) EgoAI actively saved snapshots and videos of the day

(12) EgoAI retrieves the closest souvenir shop based on Claire's taste and budget

# EGO-POLICE

(1) 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 through the shortest safe path to target places (2)

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 (3)

(4) EgoAI detected and re-identified the man before he passed Judy

(5) Judy was able to swiftly arrest him

(7) EgoAI accesses the lost-and-found database of the airport

(8) EgoAI has both thermal and multi-spectral sensors

(9) Thanks to its sensors, EgoAI calculates a low risk for explosive content

(6) Judy also appreciated the help of EgoAI when she had to manage an abandoned backpack

EgoAI projects a clear red circle around the backpack with the minimal stand-off distance (10)

(11) EgoAI connects Judy with the bomb squad and live-shares the observed scene

Thanks to EgoAI, all the relevant events are saved and transformed into a document with related images and video recordings (13)

(12) EgoAI guides Judy with exact instructions to grasp the backpack and open it

(14) The sensitive information is properly identified and secured under admin rights to protect citizens' privacy

# EGO-DESIGNER

(1) 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

(2) 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

(4) EgoAI has access to the database of the equipment warehouse; Stanley can search for the available pieces of furniture

(3)

(5) 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 (6)

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

(8) 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

Plizzari, C., Goletto, G., Furnari, A., Bansal, S., Ragusa, F., Farinella, G. M., Damen., D. & Tommasi, T. (2023). An Outlook into the Future of Egocentric Vision. *IJCV, 2024*.

**EGO-Home**

| Section | Task | Stories |
|---|---|---|
| 4.2 | 3D Scene Understanding | 1 2 3 4 7 8 9 |
| 4.3 | Object and Action Recognition | 1 5 6 10 |
| | Measuring Systems | 6 |
| 4.11 | Dialogue | 6 |
| 4.10 | Summarisation and Retrieval | 7 |
| 4.7 4.8 4.6 | Full-Body\Hand Pose and Social Interaction | 9 |
| | Medical Imaging | 10 |
| | Messaging | 10 11 |
| 4.10 | Summarisation | 11 |

**EGO-Worker**

| Section | Task | Stories |
|---|---|---|
| | Safety Compliance Assessment | 1 |
| 4.1 | Localisation and Navigation | 2 5 |
| | Messaging | 4 |
| 4.8 | Hand-Object Interaction | 5 |
| 4.4 | Action Anticipation | 6 |
| | Skill Assessment | 7 |
| 4.11 | Visual Question Answering | 8 |
| 4.10 | Summarisation | 8 |

**EGO-Tourist**

| Section | Task | Stories |
|---|---|---|
| | Recommendation and Personalisation | 1 2 8 9 10 11 |
| 4.2 | 3D Scene Understanding | 2 3 4 5 6 |
| 4.5 | Gaze Prediction | 5 |
| 4.1 | Localisation and Navigation | 3 4 8 12 |
| | Messaging | 7 |
| 4.11 | Dialogue | 8 |
| 4.3 | Action Recognition and Retrieval | 11 |
| 4.10 | Summarisation | 13 |

**EGO-Police**

| Section | Task | Stories |
|---|---|---|
| 4.1 | Localisation and Navigation | 1 2 |
| | Messaging | 1 3 11 |
| 4.3 | Action Recognition | 2 13 |
| 4.9 | Person Re-ID | 2 4 |
| 4.3 | Object Detection and Retrieval | 7 |
| | Measuring System | 8 9 |
| | Decision Making | 9 |
| 4.2 | 3D Scene Understanding | 10 |
| 4.8 | Hand-Object Interaction | 12 |
| 4.10 | Summarisation | 13 |
| 4.12 | Privacy | 14 |

**EGO-Designer**

| Section | Task | Stories |
|---|---|---|
| 4.2 | 3D Scene Understanding | 1 2 3 4 5 6 7 8 |
| | Recommendation | 3 |
| 4.3 | Object Recognition and Retrieval | 3 4 |
| 4.7 | Full-Body Pose Estimation | 5 6 |
| 4.6 | Social Interaction | 6 |
| 4.5 | Gaze Prediction | 6 |
| 4.8 | Hand-Object Interaction | 7 |
| | Messaging | 6 8 |

*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

Plizzari, C., Goletto, G., Furnari, A., Bansal, S., Ragusa, F., Farinella, G. M., Damen., D. & Tommasi, T. (2023). An Outlook into the Future of Egocentric Vision. *IJCV, 2024*.

**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 |

Plizzari, C., Goletto, G., Furnari, A., Bansal, S., Ragusa, F., Farinella, G. M., Damen., D. & Tommasi, T. (2023). An Outlook into the Future of Egocentric Vision. *IJCV, 2024.*

**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 |

Plizzari, C., Goletto, G., Furnari, A., Bansal, S., Ragusa, F., Farinella, G. M., Damen., D. & Tommasi, T. (2023). An Outlook into the Future of Egocentric Vision. *IJCV, 2024.*

**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 | 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) | | | ✓ | ✓ | ✓ | ✓ | | ✓ | | ✓ | ✓ | |

Plizzari, C., Goletto, G., Furnari, A., Bansal, S., Ragusa, F., Farinella, G. M., Damen., D. & Tommasi, T. (2023). An Outlook into the Future of Egocentric Vision. *IJCV, 2024.*

# Industrial Applications

# NEXT VISION
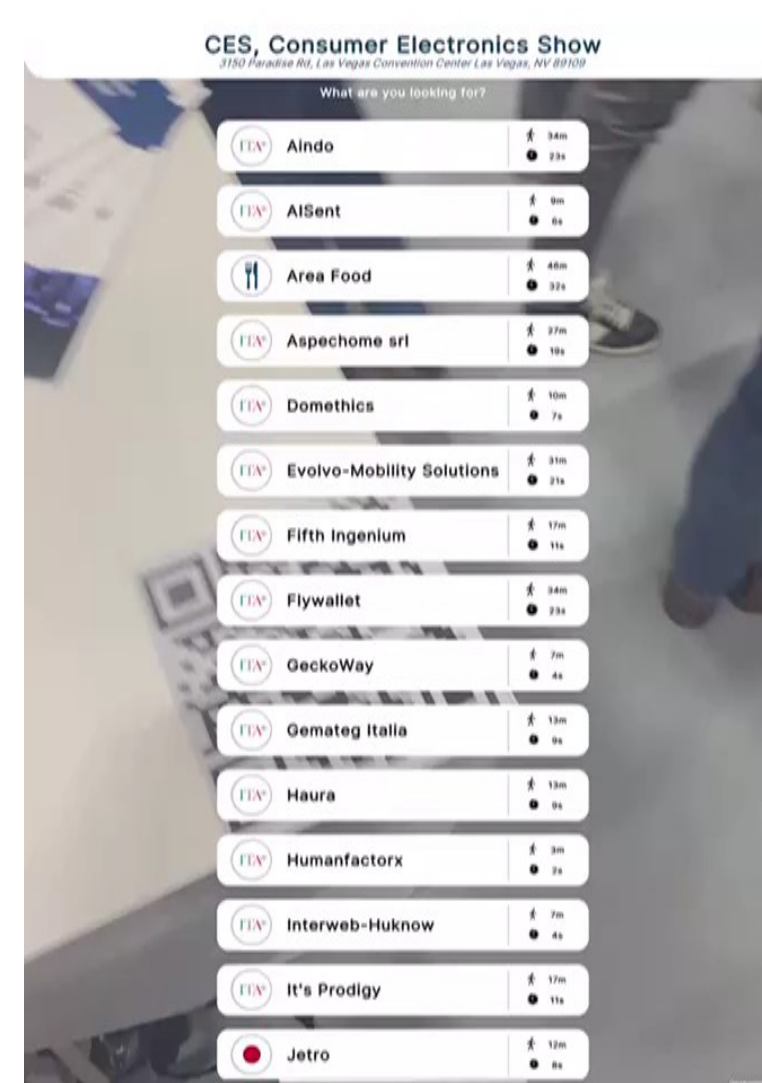
## Spin-off of the University of Catania

https://www.nextvisionlab.it/

https://drive.google.com/file/d/1lle4yF6b1kLp9P3yw
qKOi77koTvn5OuE/view?usp=share_link



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

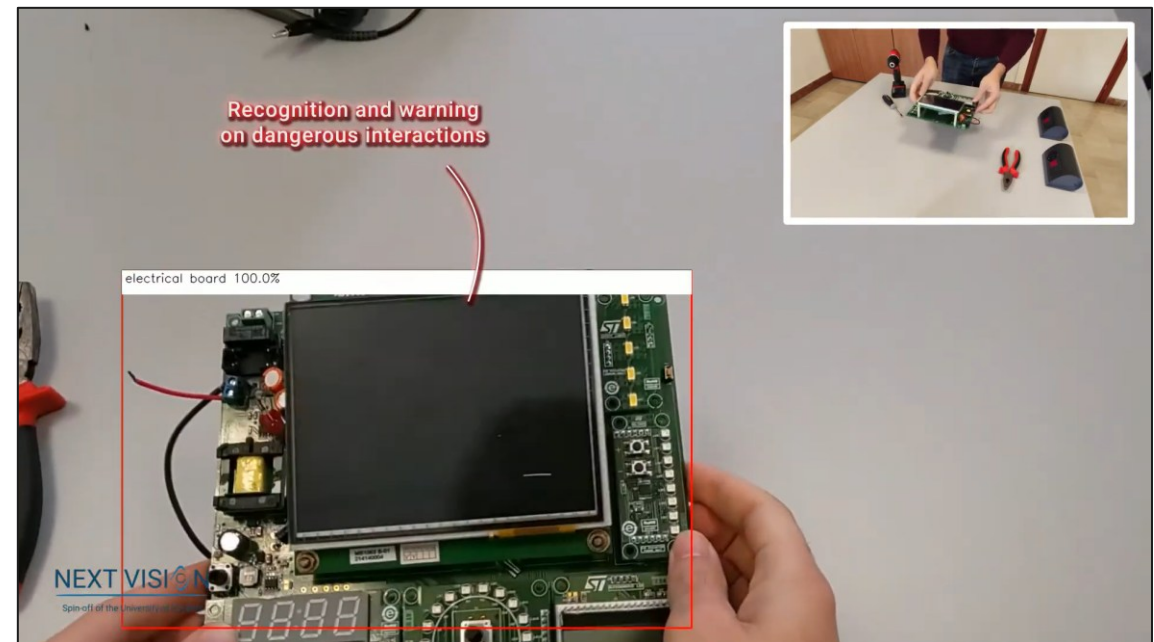Michele Mazzamuto, Francesco Ragusa, Antonino Furnari, Irene D'Ambra, Antonia Guarriera, Armando Sorbello, Giovanni Maria Farinella (2024). A Mixed Reality Application to Help Impaired People Rehabilitate Outside Clinical Environments. In IEEE International Conference on Metrology for eXtended Reality,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

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

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



Use cases



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

https://drive.google.com/file/d/1oOvhVbyyR7AZ35I-V90Zy7RyRTR7lkD4/view?usp=drive_link

# Skill Assessment



**Beginner**

**Expert**

Engine assembly

Duration: 00:09

Cumulative: 00:09

Body assembly

Duration: 00:14

Cumulative: 00:23

Replay

**Step**

**Video**

**Step**

**Video**

**Step Info**

**Video Tutorial**

**Step**

**Info**

**Video**

**Tutorial**

**Step**

**Video**

**Procedure**

Engine assembly

Duration: 00:09

Cumulative: 00:09

Body assembly

Duration: 00:14
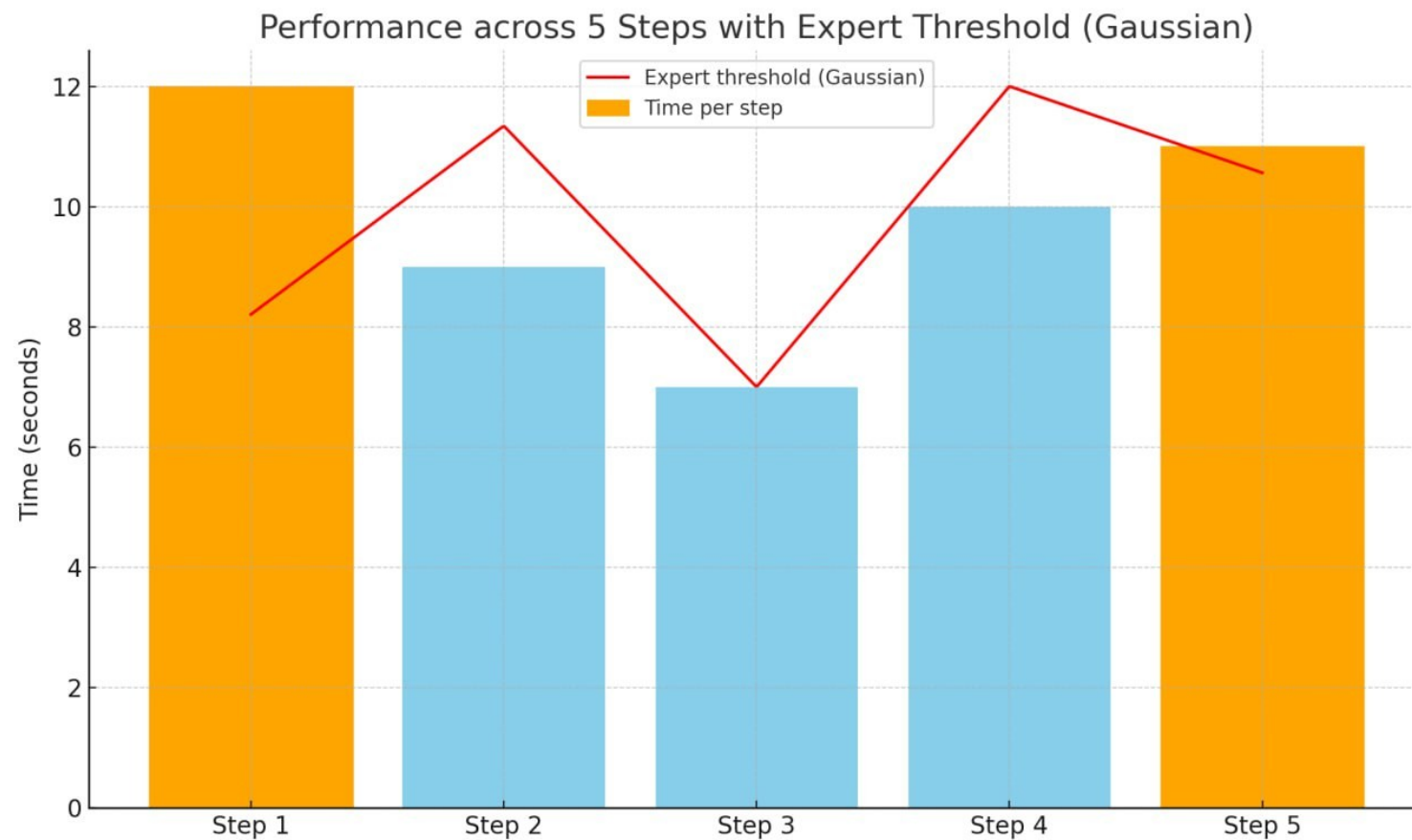
Cumulative: 00:23

Wheel assembly

3/3 STEP

03:06

Replay

Engine assembly

Duration: 00:09

Cumulative: 00:09

Body assembly

Duration: 00:14

Cumulative: 00:23

Wheel assembly

3/3 STEP

03:06

Replay

**Step**

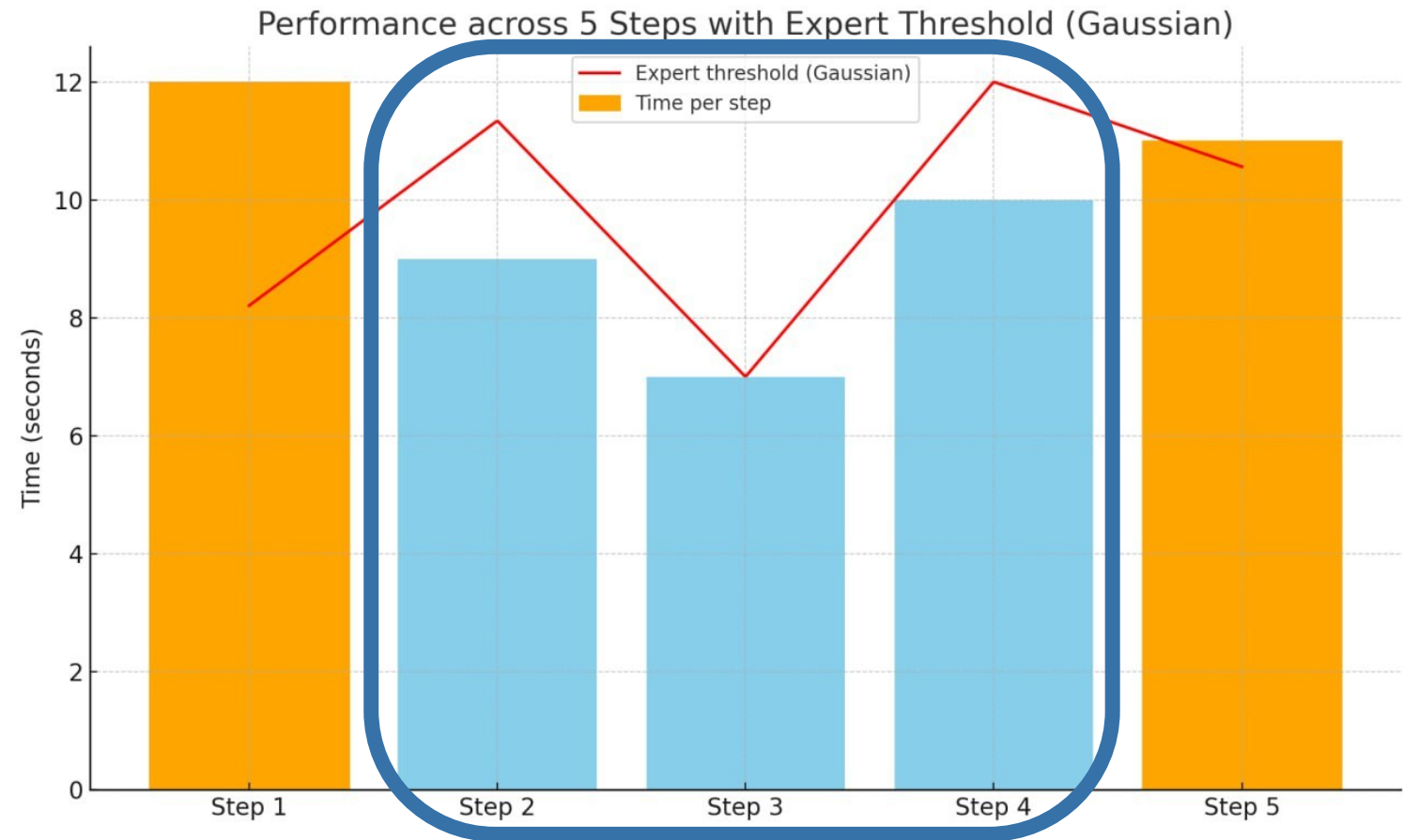**Video**

**Procedure**

Wheel assembly

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

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

3/3
STEP

03:06

Replay

**Step**

**Video**

**Procedure**

**Step**

**Video**

**Procedure**

**Step Info**   **Procedure Info**

**Step Info**  **Procedure Info**

**Step Info**

**Procedure Info**

Performance across 5 Steps with Expert Threshold (Gaussian)

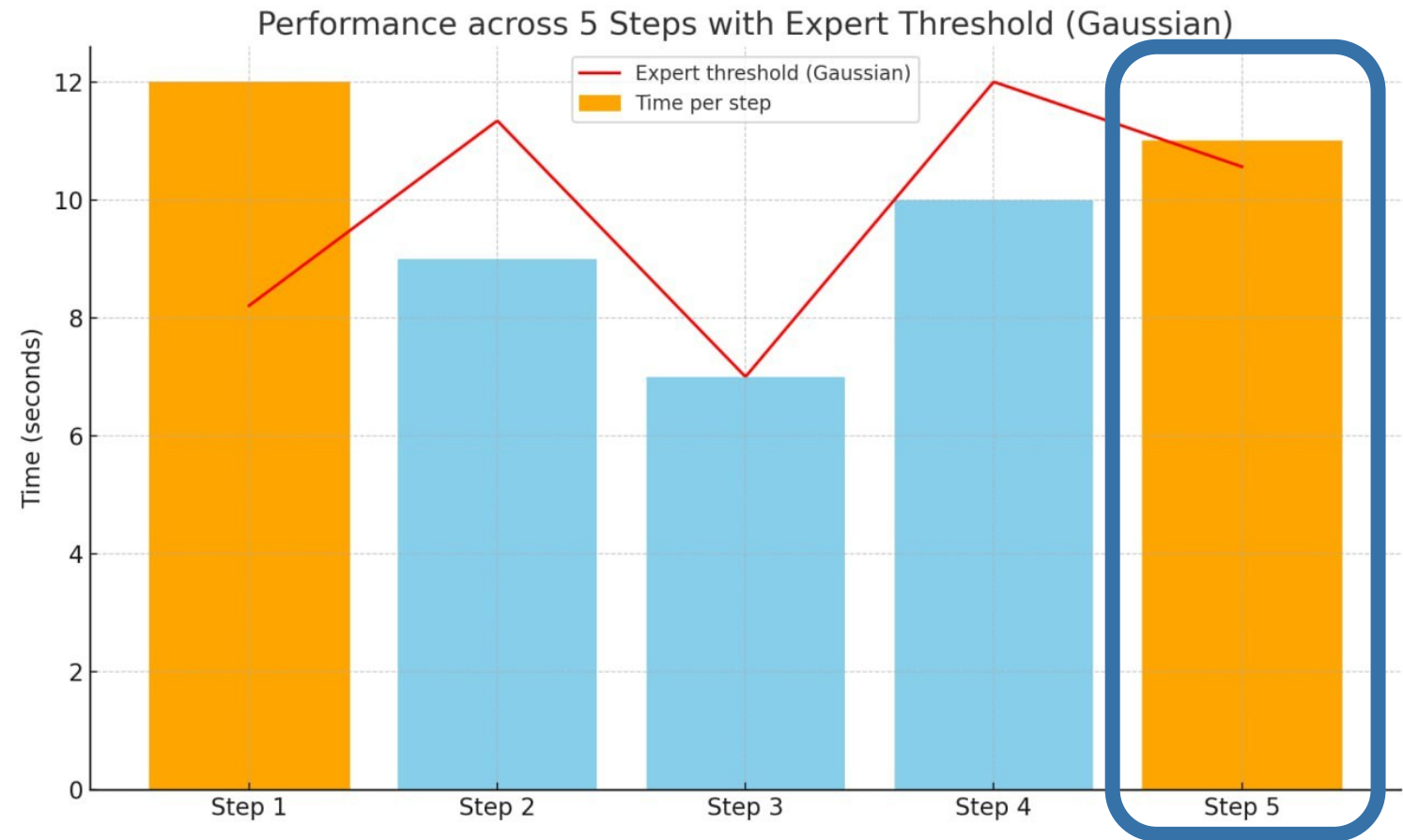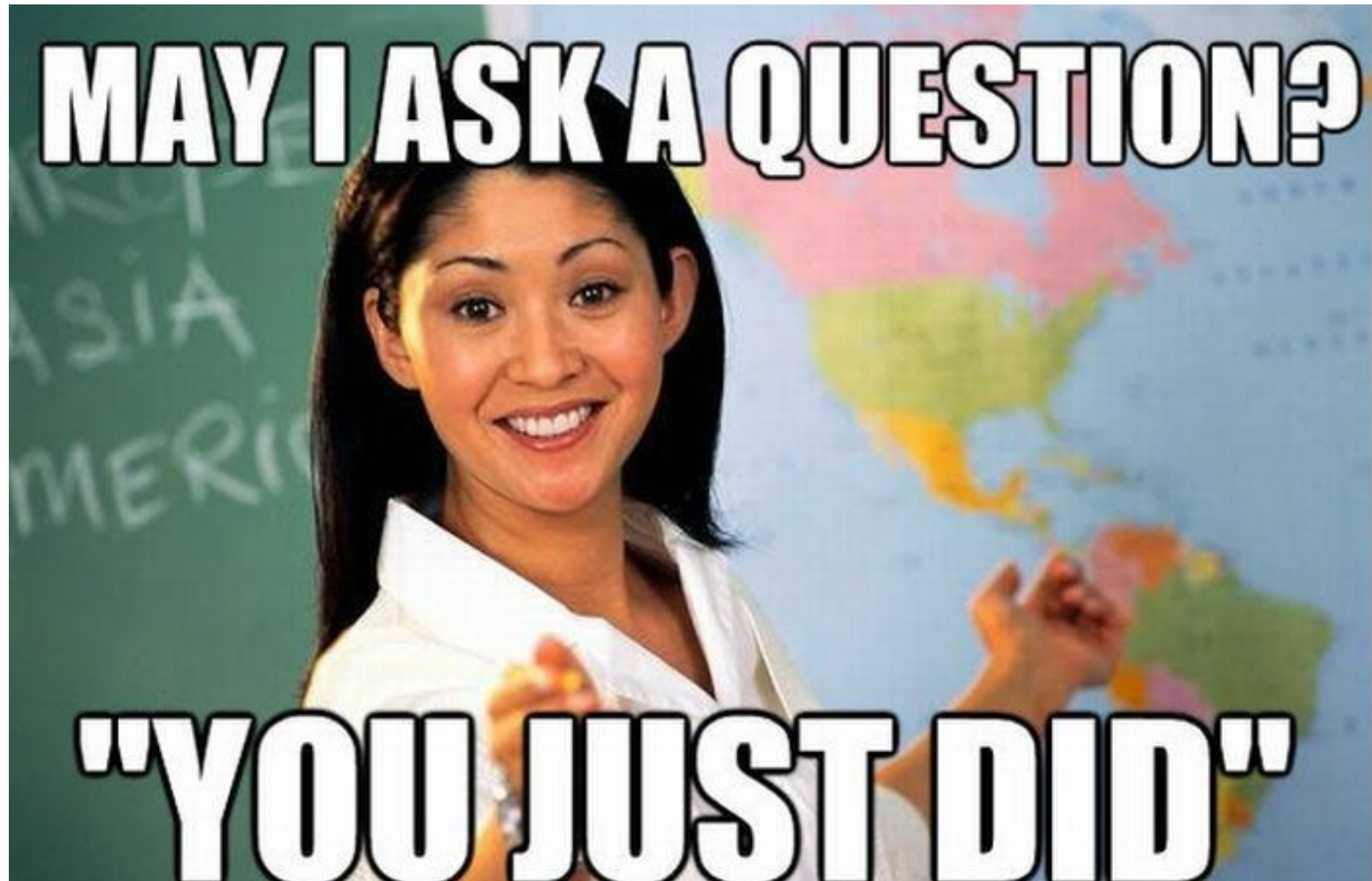**Step Info**

**Procedure Info**

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



Hardware is increasingly available as big tech gests interested.



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

**THANK YOU!**

# Egocentric Vision:
# Emerging Trends and Human-Centric Applications

## 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/

2) **Part II: Hand-Object Interactions in Egocentric Vision [15.50 – 16.50]**

   a) **Introduction to Hand-Object Interactions Detection**

   b) **Datasets and Benchmarks for Hand-Object Interactions in Egocentric Vision**

   c) **Models and Architectures for Hand-Object Interactions Detection**

   d) **Open Challenges**