# 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 - 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.**

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

# Part II

## Fundamental Tasks for First Person Vision Systems

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

Video Quality

Field of View

Wearing Modality

Other Modalities

Università di Catania

- Try to get a high quality camera to get high quality images!
- Egocentric video is subject to motion blur and exposure issues.
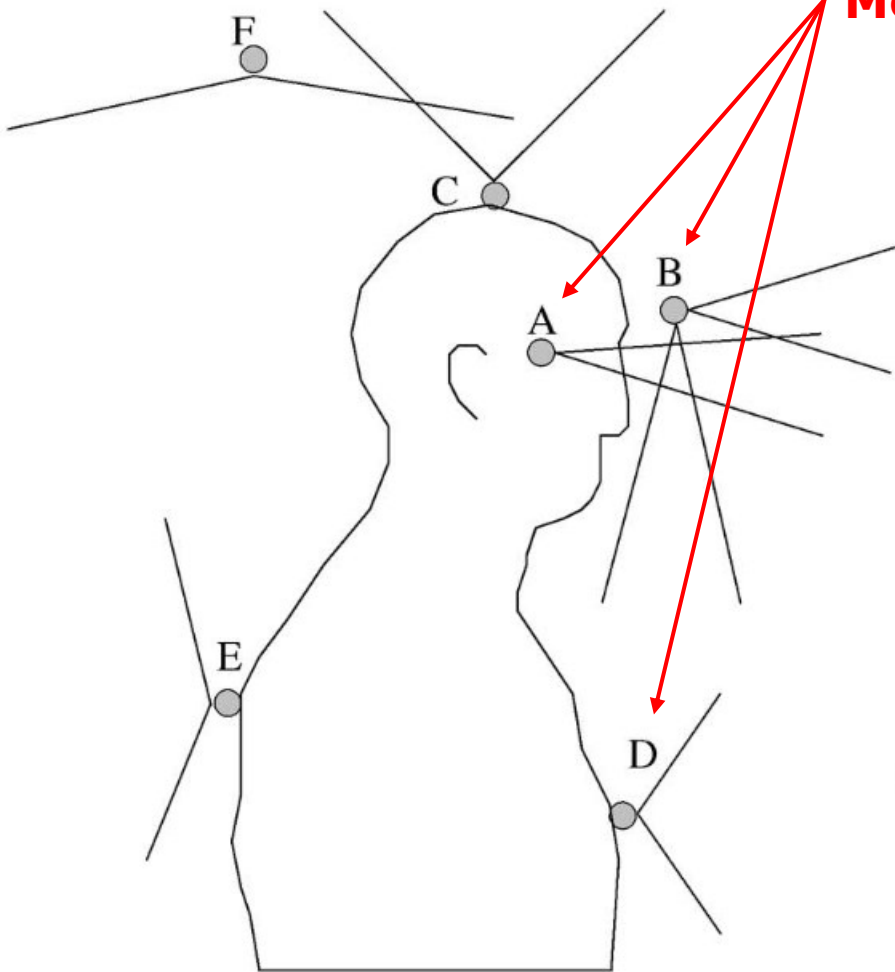
**High Quality Video Obtained with a GoPro**

**Average Quality Video**

**A,B: head mounted, D: chest mounted**

**Most Common Wearing Modalities**



**A**

**B (frontward)**

**B (downward)**

**D**

Mayol-Cuevas, W. W., Tordoff, B. J., & Murray, D. W. (2009). On the choice and placement of wearable vision sensors. *IEEE Transactions on Systems, Man, and Cybernetics-Part A: Systems and Humans*, *39*(2), 414-425.

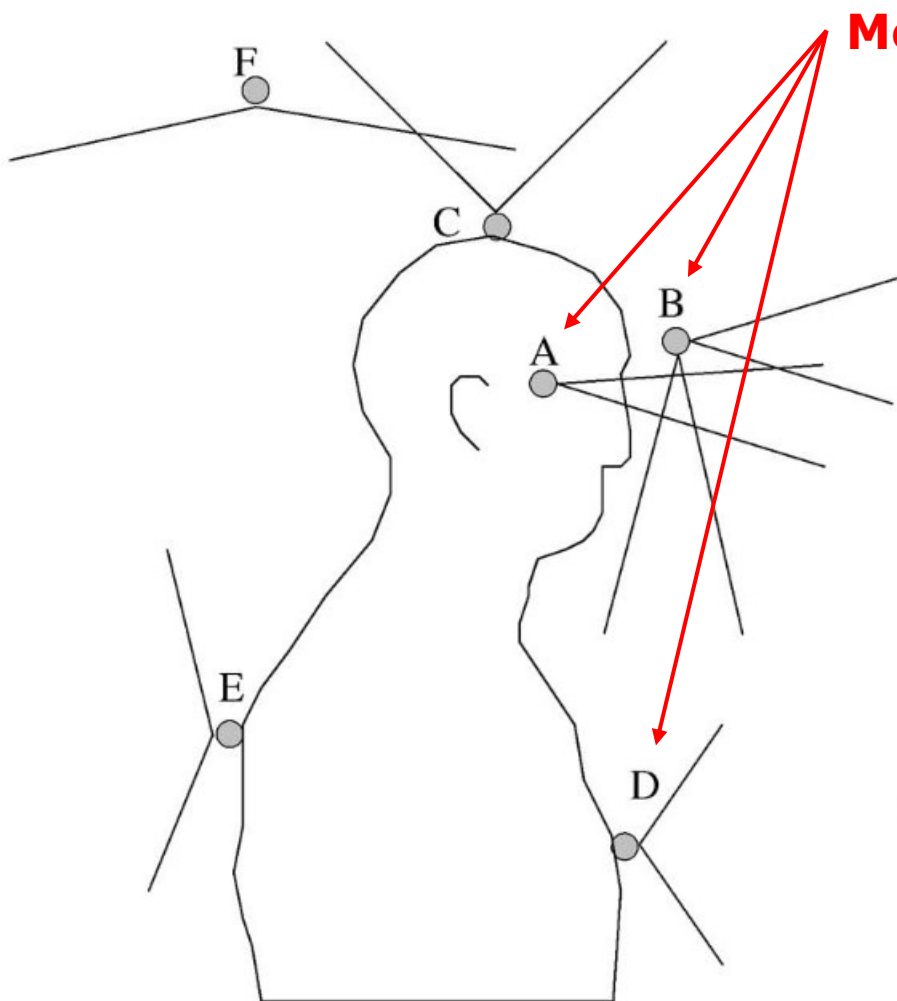**Most Common Wearing Modalities**

- A-B are best to capture objects:

  - A, B (frontward) to capture objects in front of the subjects (e.g., paintings in a museum);

  - B (downward) to capture objects manipulated with hands (e.g., kitchen);

- Chest-mounted cameras (D) are less obtrusive and give stable video, but they may miss details on what the user is looking at;
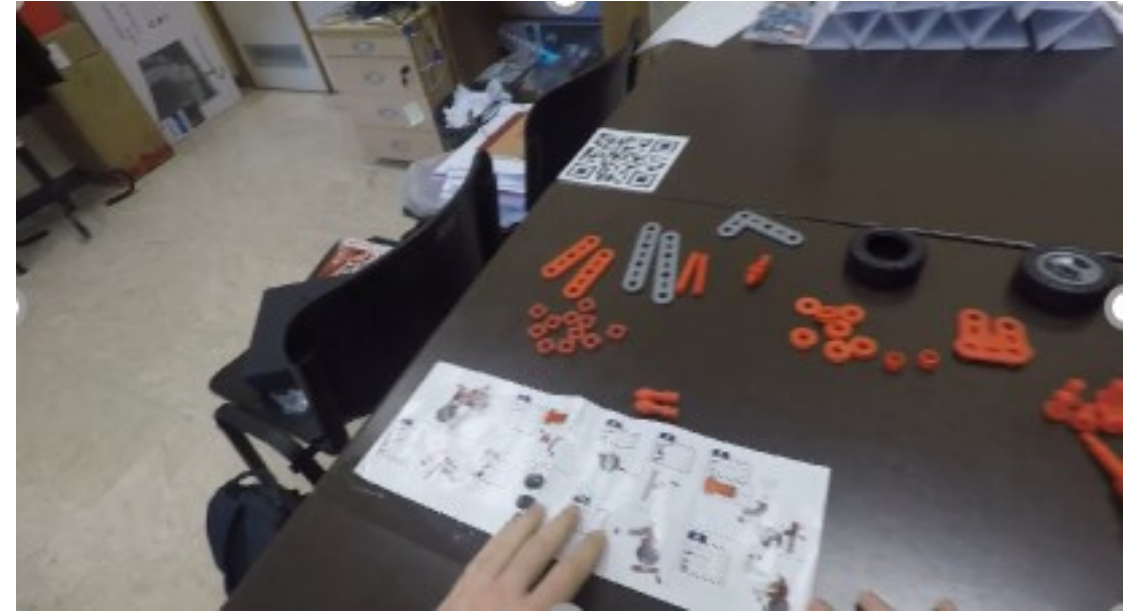
Mayol-Cuevas, W. W., Tordoff, B. J., & Murray, D. W. (2009). On the choice and placement of wearable vision sensors. *IEEE Transactions on Systems, Man, and Cybernetics-Part A: Systems and Humans*, *39*(2), 414-425.

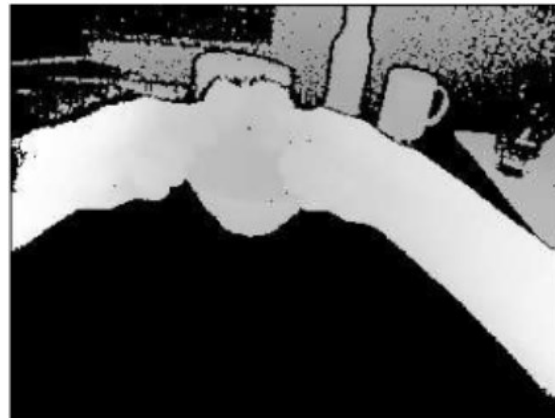A wide FOV allows to capture more scene but it may introduce distortion

**Narrow Angle**                    **Wide Angle**

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



Wan, S., & Aggarwal, J. K. (2015). Mining discriminative states of hands and objects to recognize egocentric actions with a wearable RGBD camera. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops* (pp. 36-43).

https://github.com/microsoft/HoloLensForCV

**Microsoft HoloLens Research Mode**

- Microsoft HoloLens has a «Research Mode» which allows to access:

  - short-range depth

  - long-range depth;

  - IR reflectivity;



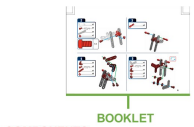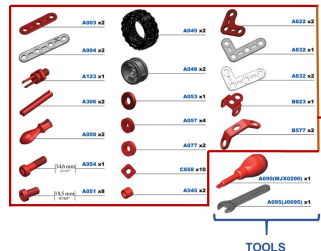https://docs.microsoft.com/en-us/windows/mixed-reality/research-mode

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

However, gaze trackers generally require a calibration process (and some expertise).
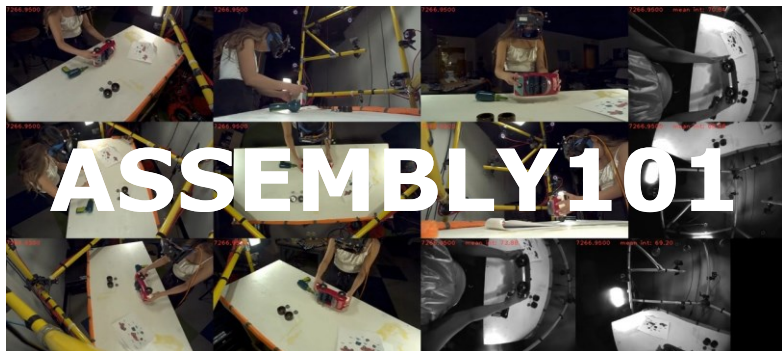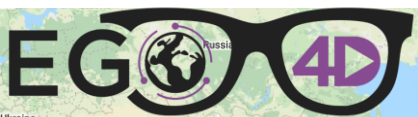
F. Ragusa, A. Furnari, S. Livatino, G. M. Farinella. The MECCANO Dataset: Understanding Human-Object Interactions from Egocentric Videos in an Industrial-like Domain. WACV 2021 (ORAL) (https://arxiv.org/abs/2010.05654).

**Università di Catania**

**MECCANO**

EPIC KITCHENS

HOI4D: A 4D Egocentric Dataset for Category-Level Human-Object Interaction

EGO 4D

ASSEMBLY101

ADL

EGO-CH

EGTEA Gaze+

Charades-Ego

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. *arXiv preprint arXiv:2308.07123*. ⬅ **And More**

Università di Catania

| Dataset | URL | Settings | Annotations | Goal |
|---------|-----|----------|-------------|------|
| EGO-EXO4D | https://ego-exo4d-data.org/ | 839 participants performing procedural and physical activities. | Natural language descriptions, segmentation masks, temporal segments of keysteps, task-graphs, profiency labels, 3D human pose | Keystep Recognition, Proficiency Estimation, Relation, Pose Estimation |
| EGO4D | https://ego4d-data.org/ | 931 participants performing different activities in different domains. | Different temporal and spatial annotations related to 5 benchmarks | Episodic Memory, Hand-Object Interaction, Audio-Visual Diarization, Social Interactions, Forecasting |
| EPIC-KITCHENS-100 | https://epic-kitchens.github.io/2020-100 | Subjects performing unscripted actions in their native kitchens. | Temporal segments | Action recognition, detection, anticipation, retrieval. |
| MECCANO | https://iplab.dmi.unict.it/MECCANO/ | 20 subjects assembling a toy motorbike. | Temporal segments, active objects, human-object interactions | Action recognition, Active object detection, Egocentric Human-Object Interaction Detection |
| ASSEMBLY101 | https://assembly-101.github.io/ | 53 subjects assembling in a cage settings 101 children's toys. | Temporal segments, 3D hand poses | Action recognition, Action Anticipation, Temporal Segmentation |

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

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

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

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

# Datasets (non-exhaustive)

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

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

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

## 12 Egocentric Vision Research Tasks

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

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. *arXiv preprint arXiv:2308.07123*.

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

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. *arXiv preprint arXiv:2308.07123*.

# Localization

## SCENE RECOGNITION



INSIDE CITY

off-the-shelf detectors

## CAMERA POSE-ESTIMATION



ROOM D ROOM C ROOM B

ROOM D

●user

**coordinates or estimated camera pose**

ROOM A

3D reconstruction of the building

**Level of Description**

−

＋

**Amount of Data**

## ROOM-LEVEL RECOGNITION



ROOM D ROOM C ROOM B

**room-level localization**

●user

ROOM D

ROOM A

moderate amount of training data

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

? → Inside city

Street

Highway

Coast

…

## COMPUTATIONALLY INEXPENSIVE ALGORITHMS

### GIST Descriptor

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

### DCT-GIST (runs on the IGP pipeline)

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

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



GT: cafeteria
top-1: cafeteria (0.179)
top-2: restaurant (0.167)
top-3: dining hall (0.091)
top-4: coffee shop (0.086)
top-5: restaurant patio (0.080)

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

*A 10 million Image Database for Scene Recognition* B. Zhou, A. Lapedriza, A. Khosla, A. Oliva, and A. Torralba *IEEE Transactions on Pattern Analysis and Machine Intelligence, 2017*

**SCENE RECOGNITION**



INSIDE CITY

off-the-shelf detectors

**CAMERA POSE-ESTIMATION**

ROOM D    ROOM C    ROOM B

ROOM D

● user

**coordinates or estimated camera pose**

ROOM A

3D reconstruction of the building

**Level of Description**

—

+

**Amount of Data**

**ROOM-LEVEL RECOGNITION**

ROOM D    ROOM C    ROOM B

**room-level localization**

● user

ROOM D

ROOM A

moderate amount of training data

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

(videos acquired in the different contexts)

Training Set

(frames extracted from videos
acquired in the different contexts)

F. Ragusa, A. Furnari, S. Battiato, G. Signorello, G. M. Farinella (2020). EGO-CH: Dataset and Fundamental Tasks for Visitors Behavioral Understanding using Egocentric Vision . Pattern Recognition Letters.

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

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

**Training Set (room-based images)**

CNN

**There are no training negatives!**

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

**1. Discrimination**   estimation of $P(y_i | I_i, y_i \neq 0)$

temporal window

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

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

**2. Negative Rejection**   estimation of $P(y_i | I_i)$

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

**3. Sequential Modelling**   application of HMM

A. Furnari, G. M. Farinella, S. Battiato, Personal-Location-Based Temporal Segmentation of Egocentric Video for Lifelogging Applications, Journal of Visual Communication and Image Representation, 2017.

F. Ragusa, A. Furnari, S. Battiato, G. Signorello, G. M. Farinella. Egocentric Visitors Localization in Cultural Sites. In Journal on Computing and Cultural Heritage (JOCCH), 2019.

G. M. Farinella, G. Signorello, S. Battiato, A. Furnari, F. Ragusa, R. Leonardi, E. Ragusa, E. Scuderi, A. Lopes, L. Santo, M. Samarotto. VEDI: Vision Exploitation for Data Interpretation. In 20th International Conference on Image Analysis and Processing (ICIAP), 2019

**SCENE RECOGNITION**



INSIDE CITY

off-the-shelf detectors

**CAMERA POSE-ESTIMATION**



ROOM D    ROOM C    ROOM B

ROOM D

● user

**coordinates or estimated camera pose**

ROOM A

3D reconstruction of the building

**Level of Description**
─────────
**Amount of Data**

−
+

**ROOM-LEVEL RECOGNITION**



ROOM D    ROOM C    ROOM B

**room-level localization**

● user

ROOM D

ROOM A

moderate amount of training data

**Images**

**3D Model**

**Structure from Motion (SfM)**



P1    P2    P3

(P,Q)

**Attach estimated 6DOF pose to each image**

**Arbitrary Coordinate System (pose/scale)**

PCA

**camera poses**

**rotated poses**

**scaled/aligned poses**

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



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

Use deep metric learning to <u>learn</u> a representation function $\varphi$ which maps close to each other images of nearby locations



**1-NN Search**

(15,21,15°)   (37,144,-12°)

(16,19,13°)

query image

(15,21,15°)

representation space

E. Spera, A. Furnari, S. Battiato, G. M. Farinella, Egocentric Shopping Cart Localization, International Conference on Pattern Recognition (ICPR), 2018
S. A. Orlando, A. Furnari, S. Battiato, G. M. Farinella. Image-Based Localization with Simulated Egocentric Navigations. VISAPP 2019

Large-Scale Visual Localization

Home  Course Description  Organizers

ICCV 2021 Tutorial

Large-Scale Visual Localization

Sunday, October 17th, 2021

## Course Information

- **When:** Sunday, October 17th, 2021

- **Where:** Online at https://youtu.be/RaVPiIGhdWk

- **Time:** half-day tutorial - starts at 2:30 pm CEST (ics)

- **Preliminary Schedule**

  - Part I: Image Retrieval for Coarse Localization (Giorgos, Yannis)

    - Image Retrieval & Visual Representation [50 min] (Giorgos) [slides]

    - Metric learning: knowledge transfer, data augmentation, and attention [20 min] (Yannis) [slides]

https://sites.google.com/view/lsvpr2021/home

# Object Detection/Interaction

Università di Catania

**Objects and Actions are tight!**

**Useful to know what is in the scene**

**Useful to know what actions can be performed**



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

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).

Faster-RCNN
(bounding boxes)

RetinaNet
(bounding boxes - faster)

Mask-RCNN
(boxes + segments)

YOLO
(much faster, but less accurate)

https://github.com/facebookresearch/detectron2

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

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

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

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



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

# Train/Finetune your own object detector

**Homes**

**ADL**

**(2012)**

**20 subjects, 42 classes, 32k images, 137k boxes**

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

**Kitchens**

**EPIC-KITCHENS-55**

**(2018)**

**32 subjects, 323 classes, 221k images, 450k boxes**

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

**Museums**

**EGO-CH**

**(2020)**

**10 subjects, 226 classes, 177k images**

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

**Industial-Like**

**MECCANO**

**(2023)**

**20 subjects, 20 object classes, 300k boxes**

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

- In some scenarios, it could be necessary to fine-tune an object-detector with application-specific data.

- Main egocentric datasets providing bounding box annotations.

- EGO4D is multi- domain annotated with 295K bounding boxes.

NEW EgoObjects!
114K annotated frames
https://github.com/facebookresearch/EgoObjects

# Understanding human-object interactions (HOI)

**Hands**

**Active Objects**

**Passive Objects**

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



An «augmented» detector which recognizes:
- The left hand;
- The right hand;
- The interacted object.

**VISOR DATASET**

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

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

Standard approach:

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

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

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

## Can simulated data help?

**ENIGMA Laboratory**

**19 objects categories**



Rosario Leonardi, Francesco Ragusa, Antonino Furnari, Giovanni Maria Farinella (2022). Egocentric Human-Object Interaction Detection Exploiting Synthetic Data . In International Conference on Image Analysis and Processing (ICIAP)

Rosario Leonardi, Francesco Ragusa, Antonino Furnari, Giovanni Maria Farinella (2022). Egocentric Human-Object Interaction Detection Exploiting Synthetic Data . In International Conference on Image Analysis and Processing (ICIAP)

Rosario Leonardi, Francesco Ragusa, Antonino Furnari, Giovanni Maria Farinella (2022). Egocentric Human-Object Interaction Detection Exploiting Synthetic Data . In International Conference on Image Analysis and Processing (ICIAP)

# Action Recognition

Model

**VERB**        **NOUN**

Open - Box

$v = 3$        $n = 23$

$t_s$        $t_e$

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

As defined in EPIC-KITCHENS-2020

# TAKE SCREWDRIVER



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.

**TAKE SCREWDRIVER**



**Start Action**

**Start Interaction (H-O)**



**Frame of Contact**

**TAKE SCREWDRIVER**



**End Interaction**

**Start Action**

**Start Interaction (H-O)**

**End Action**

**Frame of Contact**

**Frame of Decontact**

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

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

Università di Catania

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



Feichtenhofer, C., Fan, H., Malik, J., & He, K. (2019). Slowfast networks for video recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (pp. 6202-6211).

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



- X-Fast
- X-Temporal
- X-Spatial
- X-Depth
- X-Width
- X-Bottleneck

Feichtenhofer, C. (2020). X3D: Expanding Architectures for Efficient Video Recognition. 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 200-210.

≡ **README.md**

# PySlowFast

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

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

https://github.com/facebookresearch/SlowFast

# Anticipation

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

Next-active-object: **LOCKER**
Next action: **OPEN LOCKER**

Ivan Rodin, Antonino Furnari, Dimitrios Mavroedis, Giovanni Maria Farinella (2021). Predicting the Future from First Person (Egocentric) Vision: A Survey. Computer Vision and Image Understanding, 211, pp. 103252.

(observed video)

EPIC KITCHENS

Model

Take - Plate

(unobserved)

$$t_s - \tau_a - \tau_o$$

$$t_s - \tau_a \quad t_s$$

$$t_e$$

$\tau_o$ arbitrary

$\tau_a = 1s;$

Damen, Dima, et al. "Scaling egocentric vision: The epic-kitchens dataset." *Proceedings of the European Conference on Computer Vision (ECCV)*. 2018.

Dima Damen et al. Rescaling Egocentric Vision . International Journal on Computer Vision (IJCV). 2021

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

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

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



A. Furnari, S. Battiato, K. Grauman, G. M. Farinella, Next-Active-Object Prediction from Egocentric Videos, Journal of Visual Communication and Image Representation, 2017

**prediction**

bbox = [1391,101,531,713]
noun = *wooden block*
verb = *take*
ttc = 0.75s
score = 0.83

**Last observed frame** $(V_t)$

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

frame of contact

Input video: $V_{:t}$

$\delta$

t

t + δ

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



Francesco Ragusa, Giovanni Maria Farinella, Antonino Furnari (2023). StillFast: An End-to-End Approach for Short-Term Object Interaction Anticipation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops.

Next-active-object: **LOCKER**
Next action: **OPEN LOCKER**

What will happen in 1 second?

- The factory is a natural place for a wearable assistant;

- Closed-world assumption;

- Current research has considered different scenarios;

- No datasets in industrial-like scenarios;

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

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

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



**COMPONENTS**

| | | |
|---|---|---|
| A003 x2 | A045 x2 | A622 x2 |
| A004 x2 | | A632 x1 |
| A123 x1 | A046 x2 | A632 x2 |
| A306 x2 | A053 x1 | B823 x1 |
| A050 x2 | A057 x4 | B577 x2 |
| A054 x1 | A077 x2 | |
| A051 x8 | C658 x10 | |
| | A545 x2 | |

**TOOLS**

A090(MJX0200) x1

A095(J0095) x1

**BOOKLET**

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

Francesco Ragusa, Antonino Furnari, Salvatore Livatino, Giovanni Maria Farinella (2021). The MECCANO Dataset: Understanding Human-Object Interactions from Egocentric Videos in an Industrial-like Domain. In IEEE Winter Conference on Application of Computer Vision (WACV).

# The ENIGMA-51 Dataset



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).

# Industrial Applications

Intelligent Navigation

Image-based Localization

Augmented Reality

Multi-platform



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

https://drive.google.com/file/d/1lle4yF6b1kLp9P3ywqKOi77koTvn5OuE/view?usp=share_link

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

- **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

# Doing Research in Egocentric Vision: Where to start?

Data nowadays carries a lot of privacy/social/economic implications, so modern datasets are usually licensed.

**! pay attention to which uses are permitted!**



**Disclaimer**

EPIC-KITCHENS-55 and EPIC-KITCHENS-100 were collected as a tool for research in computer vision. The dataset may have unintended biases (including those of a societal, gender or racial nature).

**Copyright**

All datasets and benchmarks on this page are copyright by us and published under the Creative Commons Attribution-NonCommercial 4.0 International License. This means that you must give appropriate credit, provide a link to the license, and indicate if changes were made. You may do so in any reasonable manner, but not in any way that suggests the licensor endorses you or your use. You may not use the material for commercial purposes.

For commercial licenses of EPIC-KITCHENS and any of its annotations, email us at uob-epic-kitchens@bristol.ac.uk

**EGO4D License Agreement**

Obtaining the dataset or any annotations requires you first review our license agreement and accept the terms. Go here (ego4ddataset.com) to review and execute this agreement, and you will be emailed a set of AWS access credentials when your license agreement is approved, which will take ~48hrs. In the meantime, you can check out data overview & sample notebooks here to get familiar with the dataset, and can download the CLI & dataloaders to get setup in advance.

Note that licenses have the option to execute our license agreements as either an individual or on behalf of your institution. You will likely sign the license as an individual. Typically, only institutional signatories at a director or executive level can agree to license terms on behalf of an entire organization.

Also note that once approved your access credentials will expire in 14 days - you're expected to download the data locally, not to consume it from AWS. You can easily renew your license once it expires though: license renewal FAQ

**Ego4D Dataset**

This information you enter below will be used to generate a data usage agreement. You will receive an email from HelloSign which will step you through the process of signing all the agreements. You can review the data usage agreement at —

http://ego4d.github.io/pdfs/Ego4D-Licenses-Draft.pdf

Note: Only official signatories can sign as organisation

○ Individual    ○ Organization

First name    Last name

Email

Home Address

City    State / Province / County    Country

Submit

**Università di Catania**


EPIC KITCHENS

## Modern datasets are HUGE!
- EPIC-KITCHENS ~ 796 GB
- EGO4D ~ 30+ TB

### Download only certain data types

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

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

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

If you want to download only videos, then:

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

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

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

Will download both RGB and optical flow frames.

### Specifying participants

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

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

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

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

Will download only videos from `P01`, `P02` and `P03`.

https://github.com/epic-kitchens/epic-kitchens-download-scripts

### Data download

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

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

v2.0 annotations can be downloaded with:

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

### Detailed Flags

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

### Datasets
The following datasets are available (not exhaustive):

| Dataset | Description |
|---|---|
| annotations | The full set of annotations for the majority of benchmarks. |
| full_scale | The full scale version of all videos. (Provide `benchmarks` or `video_uids` filters to reduce the 5TB download size.) |
| clips | Clips available for benchmark training tasks. (Provide `benchmarks` or `video_uids` filters to reduce the download size.) |
| video_540ss | The downscaled version of all videos - rescaled to 540px on the short side. (Provide `benchmarks` or `video_uids` filters to reduce the 5TB download size.) |
| annotations_540ss | The annotations corresponding to the downscaled `video_540ss` videos - primarily differing only in spatial annotations (e.g. bounding boxes). |
| 3d | Annotations for the 3D VQ benchmark. |
| 3d_scans | 3D location scans for the 3D VQ benchmark. |
| 3d_scan_keypoints | 3D location scan keypoints for the 3D VQ benchmark. |
| imu | IMU data for the subset of videos available |
| slowfast8x8_r101_k400 | Precomputed action features for the Slowfast 8x8 (R101) model |
| omnivore_video_swinl | Precomputed action features for the Omnivore Video model |
| omnivore_image_swinl | Precomputed action features for the Omnivore Image model |
| fut_loc | Images and annotations for the future locomotion benchmark. |
| av_models | Model checkpoints for the AV/Social benchmark. |
| lta_models | Model checkpoints for the Long Term Anticipation benchmark. |
| moments_models | Model checkpoints for the Moments benchmark. |
| nlq_models | Model checkpoints for the NLQ benchmark. |
| sta_models | Model checkpoints for the Short Term Anticipation benchmark. |
| vq2d_models | Model checkpoints for the 2D VQ benchmark. |

https://github.com/facebookresearch/Ego4d/tree/main/ego4d/cli

EGO4D

## Command Line Interfaces Provided to Simplify Download

Università di Catania

**EPIC KITCHENS**

ABOUT  STATS  DOWNLOADS  CHALLENGES  TEAM

# EPIC-KITCHENS-100 2023 CHALLENGES

Challenge Details with links to ★NEW★ Codalab Leaderboards

**New** leaderboards are now open for the **challenge phase from Mon Jan 2023**. Check the results of the 2022 chalenge results below

**In 2023, we have 9 open challenges. These are**

- **New** Semi-Supervised Video Object Segmentation Challenge
- **New** Hand-Object Segmentation Challenge
- **New** TREK-150 Object Tracking Challenge
- **New** EPIC-SOUNDS Audio-Based Interaction Recognition
- Action Recognition
- Action Detection
- Action Anticipation
- UDA for Action Recognition
- Multi-Instance Retrieval

## EPIC-Kitchens 2023 Challenges

| | |
|---|---|
| Jan 23rd 2023, | All leaderboards are open (note new challenges for 2023) |
| June 1st 2023, | Server Submission Deadline at 23:00:00 UTC |
| June 6th 2023, | Deadline for Submission of Technical Reports on CMT |
| Mon June 19 2023, | Results announced at 11th EPIC@CVPR2023 workshop in Vancouver 11th EPIC@CVPR2023 workshop in Vancouver |

## Challenges Guidelines

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

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

https://epic-kitchens.github.io/2023#challenges

# Ego4D Challenge 2023

**Episodic memory:**

- Visual queries with 2D localization (VQ2D) and Visual Queries 3D localization (VQ3D): Given an egocentric video clip and an image crop depicting the query object, return the most recent occurrence of the object in the input video, in terms of contiguous bounding boxes (2D + temporal localization) or the 3D displacement vector from the camera to the object in the environment.
  - Quickstart: [Open in Colab]
- Natural language queries (NLQ): Given a video clip and a query expressed in natural language, localize the temporal window within all the video history where the answer to the question is evident.
  - Quickstart: [Open in Colab]
- Moments queries (MQ): Given an egocentric video and an activity name (e.g., a "moment"), localize all instances of that activity in the past video
- EgoTracks: Given an egocentric video and a visual template of an object, localize the bounding box containing the object in each frame of the video along with a confidence score representing the presence of the object. **[NEW for 2023]**
- PACO Zero-Shot: Retrieve the bounding box of a specific object instance from a dataset, based on a textual query describing the instance. Query is composed using object and part attributes describing the object of interest. **[NEW for 2023]**

**Hands and Objects:**

- Temporal localization: Given an egocentric video clip, localize temporally the key frames that indicate an object state change.
- Object state change classification: Given an egocentric video clip, indicate the presence or absence of an object state change.

**Audio-Visual Diarization:**

- Audio-visual speaker diarization: Given an egocentric video clip, identify which person spoke and when they spoke.
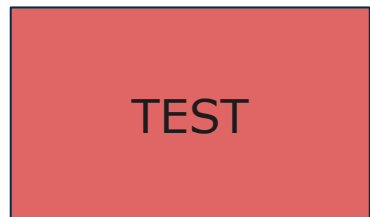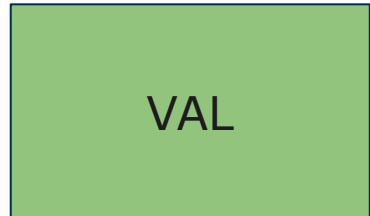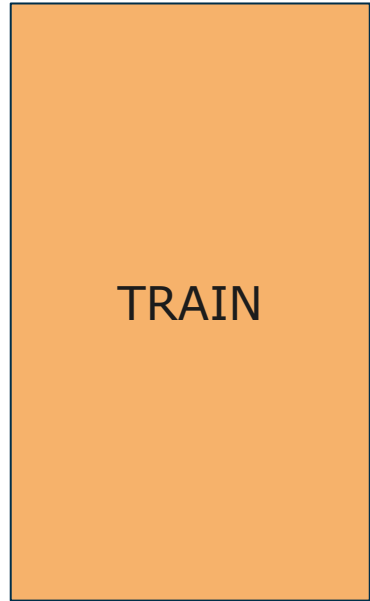- Speech transcription: Given an egocentric video clip, transcribe the speech of each person.

**Social Understanding:**

- Talking to me: Given an egocentric video clip, identify whether someone in the scene is talking to the camera wearer.
- Looking at me: Given an egocentric video clip, identify whether someone in the scene is looking at the camera wearer.
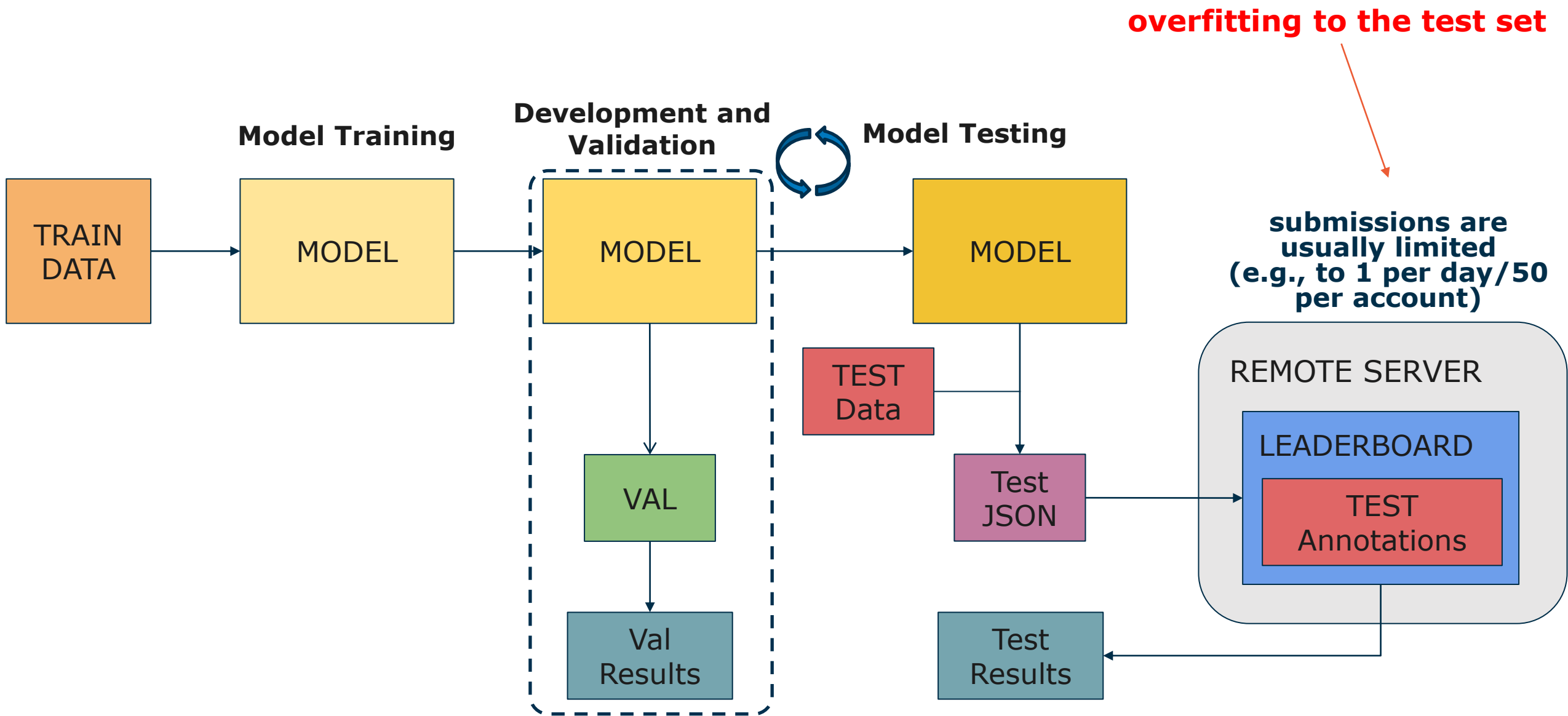
**Forecasting:**

- Short-term hand object prediction: Given a video clip, predict the next active objects, and, for each of them, predict the next action, and the time to contact.
  - Quickstart: [Open in Colab]
- Long-term activity prediction: Given a video clip, the goal is to predict what sequence of activities will happen in the future. For example, after kneading dough, list the actions that the baker will do next.

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

**TRAIN**

**VAL**

**TEST**

- Datasets are usually divided into train/val/test splits;
- All videos are publicly released;

- <u>Train</u> annotations are publicly released and meant for training models for the different challenges;

- <u>Val</u> annotations are publicly released and meant for model development and hyperparameter search;

- <u>Test</u> annotations are <u>private</u> and meant for assessing the performance of models <u>avoiding bias</u> in model design and optimization;

- Hence, the <u>only way</u> to obtain results on the test set is to send model predictions to an evaluation server.

**Model Training**

**Development and Validation**

**Model Testing**

**overfitting to the test set**

**submissions are usually limited (e.g., to 1 per day/50 per account)**

TRAIN DATA

MODEL

MODEL

MODEL

VAL

Val Results

TEST Data

Test JSON

Test Results

REMOTE SERVER

LEADERBOARD

TEST Annotations

- First Person Vision paves the way to a variety of user-centric applications;

- However, we are still missing solid building blocks related to fundamental problems of First Person Vision such as action recognition, object detection, action anticipation and human-object interaction detection;

- Consumer devices are starting to appear, but the near future of First Person Vision is in focused applications such as the ones in industrial scenarios.

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

# 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 - https://francescoragusa.github.io/