# Opportunities in Egocentric Video Understanding

University of BRISTOL
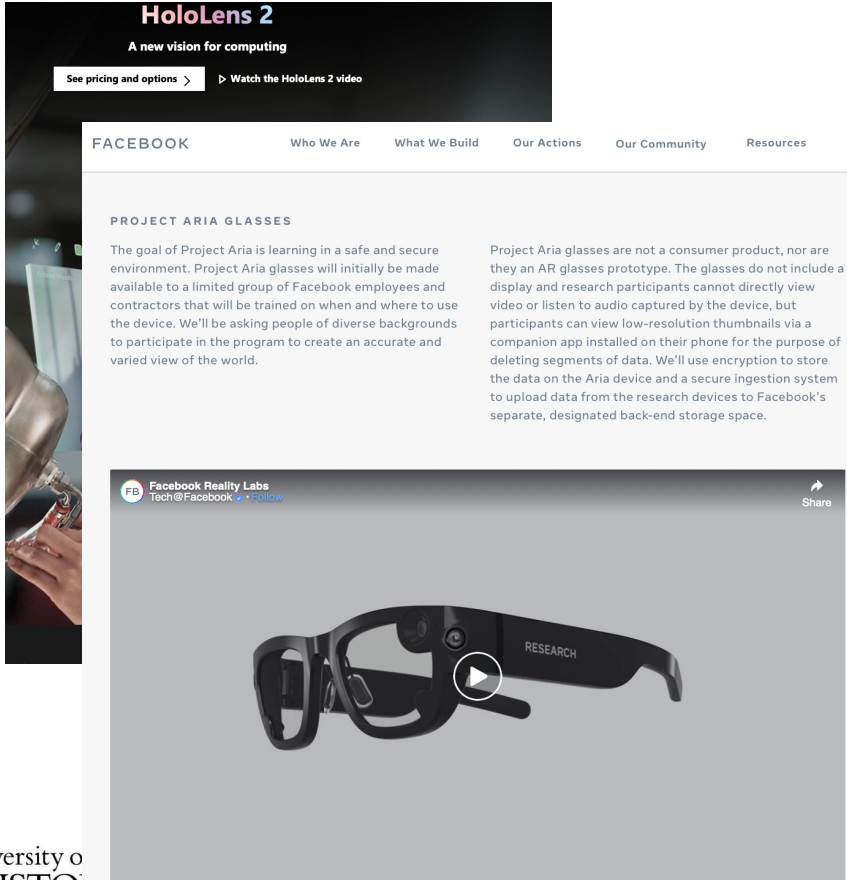
Dima Damen
WACV2024 – Waikoloa, Hawaii

Photo *Illustration* by Pelle Cass

BRISTOL

# The future…



HoloLens 2
A new vision for computing

See pricing and options

▷ Watch the HoloLens 2 video

## Samsung patent application reveals augmented reality headset design

*It comes as the Gear VR slowly fades away*

By Jon Porter

FACEBOOK

Who We Are | What We Build | Our Actions | Our Community | Resources

### PROJECT ARIA GLASSES

The goal of Project Aria is learning in a safe and secure environment. Project Aria glasses will initially be made available to a limited group of Facebook employees and contractors that will be trained on when and where to use the device. We'll be asking people of diverse backgrounds to participate in the program to create an accurate and varied view of the world.

Project Aria glasses are not a consumer product, nor are they an AR glasses prototype. The glasses do not include a display and research participants cannot directly view video or listen to audio captured by the device, but participants can view low-resolution thumbnails via a companion app installed on their phone for the purpose of deleting segments of data. We'll use encryption to store the data on the Aria device and a secure ingestion system to upload data from the research devices to Facebook's separate, designated back-end storage space.

Facebook Reality Labs
Tech@Facebook ✔ • Follow

RESEARCH

3

WACV2024 – Waikoloa, Hawaii

University of BRISTOL

# Surveillance vs Sousveillance

## Surveillance



## Sousveillance

GEORGE FLOYD

**Teen with 'cell phone and sheer guts' credited for Derek Chauvin's murder conviction**

CNNWire By Holly Yan, CNN
Wednesday, April 21, 2021 6:07PM

**Video shows Charlotte officer repeatedly hitting pinned woman during arrest: 'Not easy to watch'**

WTVD-AP
Friday, November 17, 2023

**'They could've killed him': Jacksonville family wants justice after video of arrest goes viral**

**Cyclist's GoPro footage captures**

University of BRISTOL

# Egocentric cameras are coming

What can we do with such footage?

# Data Collection Exercises

**EPIC KITCHENS**

2017 - now

100 hours
45 kitchens
4 countries
Long-term recording
Kitchen-based activities

**EGO 4D**

2020 - now

6730 hours
923 participants
74 locations
9 countries
Short-term recording
All daily activities

# Data Collection Exercises



**EGO-EXO4D**

2022 - now

Released Dec 2023
1422 hours
8 skilled activities
839 camera wearers
Ego-Exo recordings

2024 – [coming]

[new recordings]

University of
BRISTOL

**Object Recognition**

baboon

**Dataset**

mammal → dog → primate → gorilla → baboon

# Kinetics Dataset



**Action Recognition**

YouTube ᴵᵀ

🔍 absailing

University of BRISTOL

Object Recognition

**Let's collect Data!**

EPIC KITCHENS-100

University of BRISTOL

Pascal VOC
ImageNet
Kinetics
Something-Something

**Labels**

**Data**

EPIC-KITCHENS
Ego4D
…
KITTI

# EPIC-KITCHENS

University of BRISTOL

Store

Lay-down    **PUT**

Place    Rest-on

stove    burner

**hob**

stovetop    gas

*open vocab*    lay-down    stovetop

*closed vocab*    put    hob

*category*    leave    appliance

**Data** 

**Labels** 

| Data | Labels |
|---|---|
| Naturally unbalanced | Unnaturally balanced (or nearly) |
| Harder to label (exposes ambiguity) | Easier to label (hides ambiguity) |
| Closer to application | Can be expanded |
| Many research opportunities… | Single task |

**Data first brings out many opportunities**

# Opportunities in Egocentric Video Understanding

University of BRISTOL

# Tasks are harder

Detection, 3D Mapping, Tracking, VOS, Hand-Object, Generative, …

Solutions prove more rewarding

Weak supervision, Domain Adap/Gen., Audio-Visual, long-term understanding

# Tasks are harder

Detection, 3D Mapping, Tracking, VOS, Hand-Object, Generative, …

# Solutions prove more rewarding

Weak supervision, Domain Adap/Gen., Audio-Visual, long-term understanding

# Tasks are harder

Detection, 3D Mapping, Tracking, VOS, Hand-Object, Generative, …

# Solutions prove more rewarding

Weak supervision, Domain Adap/Gen., Audio-Visual, long-term understanding

# Action Detection

with: Hanyuan Wang
Majid Mirmehdi
Toby Perrett

| Task | Method | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 | Avg |
|---|---|---|---|---|---|---|---|
| Verb | BMN [18,36] | 10.8 | 9.8 | 8.4 | 7.1 | 5.6 | 8.4 |
| | G-TAD [76] | 12.1 | 11.0 | 9.4 | 8.1 | 6.5 | 9.4 |
| | Ours | **26.6** | **25.6** | **24.4** | **22.4** | **18.3** | **23.4** |
| Noun | BMN [18,36] | 10.3 | 8.3 | 6.2 | 4.5 | 3.4 | 6.5 |
| | G-TAD [76] | 11.0 | 10.0 | 8.6 | 7.0 | 5.4 | 8.4 |
| | Ours | **25.5** | **24.3** | **22.6** | **20.3** | **16.6** | **21.9** |

Table 2. **Results on EPIC-Kitchens 100 validation set**.

Zhang et al (2022). ActionFormer: Localizing Moments of Actions with Transformers. ECCV



Figure 3. Qualitative results on the EPIC-KITCHENS-100 validation set. Ground truth and predictions are shown with colour-coded class

University of BRISTOL

H Wang et al (2024). Refining Action Boundaries for One-stage Detection. *BMVCW*

Dima Damen          28
WACV2024 – Waikoloa, Hawaii

# Tasks are harder

Detection, 3D Mapping, Tracking, VOS, Hand-Object, Generative, …

# Solutions prove more rewarding

Weak supervision, Domain Adap/Gen., Audio-Visual, long-term understanding

with: Chiara Plizzari
Toby Perrett



Plizzari et al (2023). What can a cook in Italy teach a mechanic in India? Action Recognition Generalisation Over Scenarios and Locations. IEEE/CVF International Conference on Computer Vision (ICCV).

University of BRISTOL

with: Chiara Plizzari
Toby Perrett



Plizzari et al (2023). What can a cook in Italy teach a mechanic in India? Action Recognition Generalisation Over Scenarios and Locations. IEEE/CVF International Conference on Computer Vision (ICCV).

with: Chiara Plizzari
Toby Perrett

Plizzari et al (2023). What can a cook in Italy teach a mechanic in India? Action Recognition Generalisation Over Scenarios and Locations. IEEE/CVF International Conference on Computer Vision (ICCV).

University of BRISTOL

Dima Damen
WACV2024 – Waikoloa, Hawaii

32

with: Chiara Plizzari
Toby Perrett

Plizzari et al (2023). What can a cook in Italy teach a mechanic in India? Action Recognition Generalisation Over Scenarios and Locations. IEEE/CVF International Conference on Computer Vision (ICCV).

Dima Damen
WACV2024 – Waikoloa, Hawaii

33

with: Chiara Plizzari
Toby Perrett

Plizzari et al (2023). What can a cook in Italy teach a mechanic in India? Action Recognition Generalisation Over Scenarios and Locations. IEEE/CVF International Conference on Computer Vision (ICCV).

# Generalisation across Scenarios and Locations

with: Chiara Plizzari
Toby Perrett



Plizzari et al (2023). What can a cook in Italy teach a mechanic in India? Action Recognition Generalisation Over Scenarios and Locations. IEEE/CVF International Conference on Computer Vision (ICCV).

University of BRISTOL

Dima Damen
WACV2024 – Waikoloa, Hawaii

with: Chiara Plizzari
Toby Perrett

Plizzari et al (2023). What c... ... mechanic in India? Action Recognition Generalisation Over Scenarios and Locations. IEE... ...national Conference on Computer Vision (ICCV).

Dima Damen                                                                    36
WACV2024 – Waikoloa, Hawaii

with: Chiara Plizzari
Toby Perrett

- We introduce **ARGO1M**, the first dataset to perform **Action Recognition Generalisation** Over Scenarios and Locations



EGO 4D

13 locations

10 scenarios

Plizzari et al (2023). What can a cook in Italy teach a mechanic in India? Action Recognition Generalisation Over Scenarios and Locations. IEEE/CVF International Conference on Computer Vision (ICCV).

University of BRISTOL

Dima Damen
WACV2024 – Waikoloa, Hawaii

# Dataset: ARGO1M

- We introduce **ARGO1M**, the first dataset to perform **Action Recognition Generalisation** Over Scenarios and Locations



13 locations

10 scenarios

60 action classes

University of BRISTOL

Dima Damen
WACV2024 – Waikoloa, Hawaii

38

- We introduce **ARGO1M**, the first dataset to perform **Action Recognition Generalisation** Over Scenarios and Locations

**NEW**

## 1.1M samples



EGO4D

13 locations

10 scenarios

60 action classes

University of BRISTOL

Dima Damen
WACV2024 – Waikoloa, Hawaii

# Generalisation across Scenarios and Locations

with: Chiara Plizzari
Toby Perrett

ARGO1M: 1.05M action clips from 60 action classes recorded in 13 locations within 10 scenarios

University of BRISTOL

with: Chiara Plizzari
Toby Perrett

**ARGO1M**

University of BRISTOL

Dima Damen
WACV2024 – Waikoloa, Hawaii

with: Chiara Plizzari
Toby Perrett



**ARGO1M**

Cooking

**Japan**

University of BRISTOL

Dima Damen
WACV2024 – Waikoloa, Hawaii

42

with: Chiara Plizzari
Toby Perrett

**ARGO1M**

Cooking in **Japan**

**Cooking**

**Japan**

Plizzari et al (2023). What can a cook in Italy teach a mechanic in India? Action Recognition Generalisation Over Scenarios and Locations. IEEE/CVF International Conference on Computer Vision (ICCV).

Dima Damen
WACV2024 – Waikoloa, Hawaii

University of BRISTOL

with: Chiara Plizzari
Toby Perrett



**Training set**
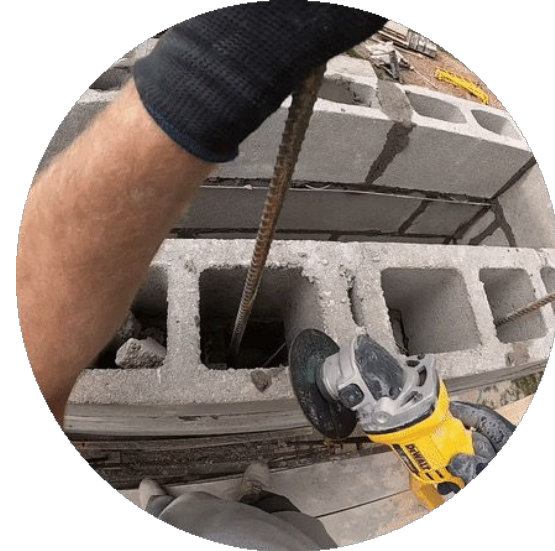
**Test set**

**Cooking** in **Japan**

Plizzari et al (2023). What can a cook in Italy teach a mechanic in India? Action Recognition Generalisation Over Scenarios and Locations. IEEE/CVF International Conference on Computer Vision (ICCV).

Dima Damen
WACV2024 – Waikoloa, Hawaii

44

with: Chiara Plizzari
Toby Perrett

**Training set**

**Test set**

**Cooking** in **Japan**

10 test sets

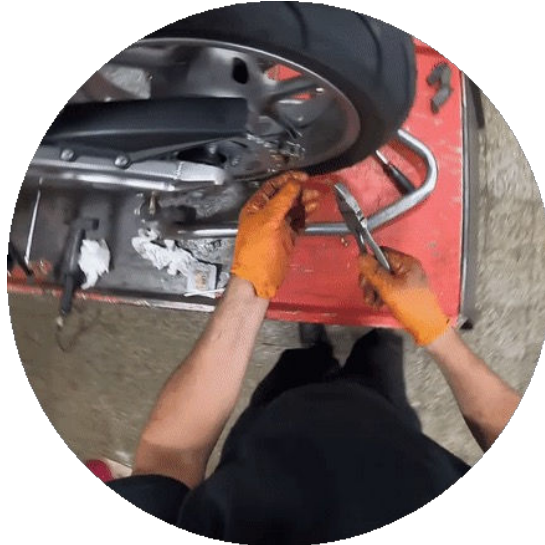| | | | | |
|---|---|---|---|---|
| **Ga** US-PNA | **Cl** US-MN | **Kn** IND | **Sh** IND | **Bu** US-PNA |
| **Me** SAU | **Sp** COL | **Co** JPN | **Ar** ITA | **Pl** US-IN |

University of BRISTOL

Plizzari et al (2023). What can a cook in Italy teach a mechanic in India? Action Recognition Generalisation Over Scenarios and Locations. IEEE/CVF International Conference on Computer Vision (ICCV).

Dima Damen
WACV2024 – Waikoloa, Hawaii

45

with: Chiara Plizzari
Toby Perrett



*He cuts the lemon strand*

University of BRISTOL

with: Chiara Plizzari
Toby Perrett



action classification

$f(v)$

$f$

$h \rightarrow L_C$

She picks tomatoes

$g$

$g(t)$

Plizzari et al (2023). What can a cook in Italy teach a mechanic in India? Action Recognition Generalisation Over Scenarios and Locations. IEEE/CVF International Conference on Computer Vision (ICCV).

University of BRISTOL

with: Chiara Plizzari
Toby Perrett



support set

She picks tomatoes

$f(v)$

$g(t)$

$L_c$

University of BRISTOL

Dima Damen
WACV2024 – Waikoloa, Hawaii

# Proposed method: CIR

with: Chiara Plizzari
Toby Perrett



**first** cross-instance reconstruction

support set

She picks tomatoes

$f(v)$

$g(t)$

$L_c$

$\oplus v$

Plizzari et al (2023). What can a cook in Italy teach a mechanic in India? Action Recognition Generalisation Over Scenarios and Locations. IEEE/CVF International Conference on Computer Vision (ICCV).

University of BRISTOL

# Proposed method: CIR

with: Chiara Plizzari
Toby Perrett



**first** cross-instance reconstruction

support set

video-text association

$L_{rt}$

$f(v)$

$g(t)$

She picks tomatoes

Plizzari et al (2023). What can a cook in Italy teach a mechanic in India? Action Recognition Generalisation Over Scenarios and Locations. IEEE/CVF International Conference on Computer Vision (ICCV).

University of BRISTOL

Dima Damen
WACV2024 – Waikoloa, Hawaii

50

with: Chiara Plizzari
Toby Perrett



first cross-instance reconstruction

support set

She picks tomatoes

$f(v)$

$g(t)$

video-text association

$L_{rt}$

Plizzari et al (2023). What can a cook in Italy teach a mechanic in India? Action Recognition Generalisation Over Scenarios and Locations. IEEE/CVF International Conference on Computer Vision (ICCV).

University of BRISTOL

Dima Damen
WACV2024 – Waikoloa, Hawaii

51

with: Chiara Plizzari
Toby Perrett



support set

second cross-instance reconstruction

$f(v)$

$g(t)$

She picks tomatoes

$\oplus v'$

$\oplus v$

$L_c$

$L_{rt}$

University of BRISTOL

Dima Damen
WACV2024 – Waikoloa, Hawaii

# Proposed method: CIR

with: Chiara Plizzari
Toby Perrett



action classification

$f(v)$

$f$

$h \rightarrow L_c$

support set

second cross-instance reconstruction

$\oplus v'$

$h \rightarrow L_{rc}$

$g(t)$

She picks tomatoes

$g$

$\oplus v$

$L_{rt}$

Plizzari et al (2023). What can a cook in Italy teach a mechanic in India? Action Recognition Generalisation Over Scenarios and Locations. IEEE/CVF International Conference on Computer Vision (ICCV).

University of BRISTOL

with: Chiara Plizzari
Toby Perrett

action classification

$f(v)$

$f$

$h \rightarrow L_c$

----------→ inference

Plizzari et al (2023). What can a cook in Italy teach a mechanic in India? Action Recognition Generalisation Over Scenarios and Locations. IEEE/CVF International Conference on Computer Vision (ICCV).

Dima Damen
WACV2024 – Waikoloa, Hawaii

Chiara Plizzari
Toby Perrett
Dima Damen

#C C drops the cut vegetables



**query**



**support 1**     **support 2**     **support 3**     **support 4**     **support 5**
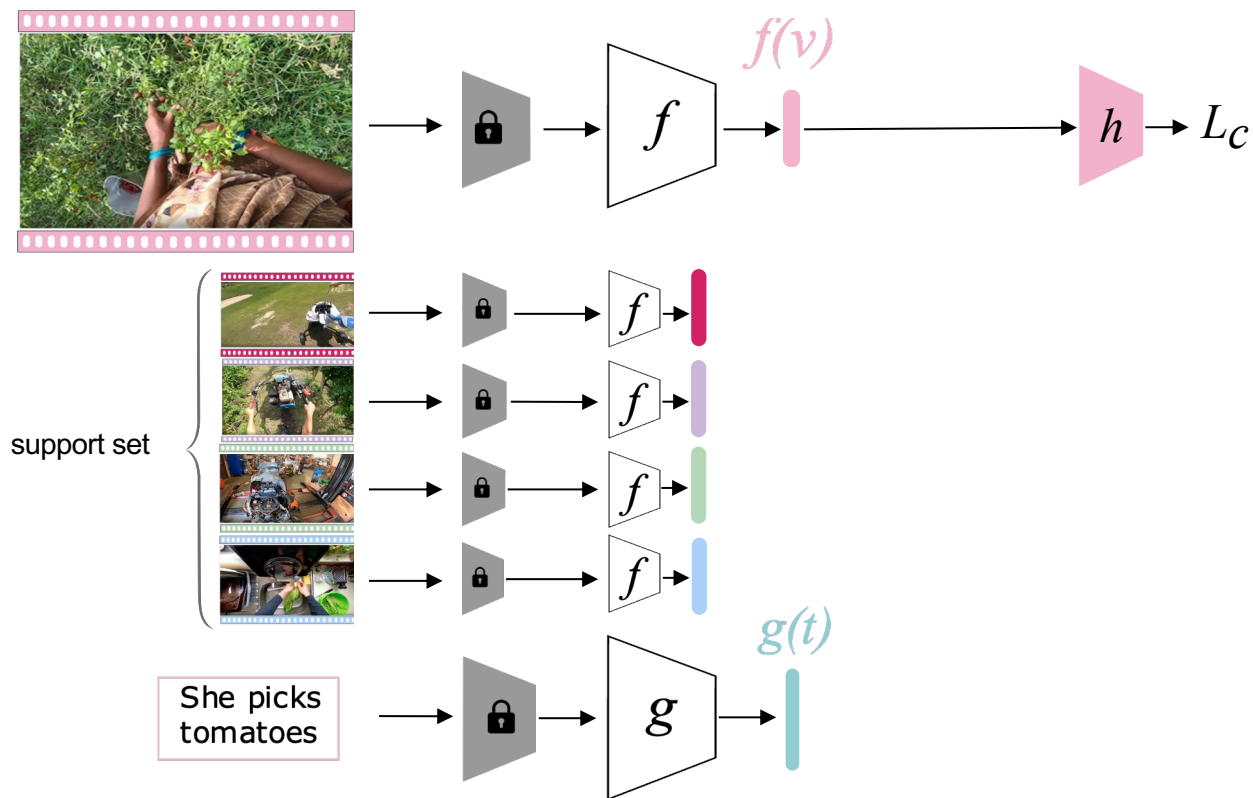
55
Hawaii

Plizzari et al (2023). What can a cook in Italy teach a mechanic in India? Action Recognition Generalisation Over Scenarios and Locations. IEEE/CVF International Conference on Computer Vision (ICCV).

**What can a cook in Italy teach a mechanic in India?**
**Action Recognition Generalisation Over Scenarios and Locations**

Chiara Plizzari    Toby Perrett    Barbara Caputo    Dima Damen

Politecnico di Torino, Italy    University of Bristol, United Kingdom

**Abstract**

We propose and address a new generalisation problem: can a model trained for action recognition successfully classify actions when they are performed within a previously unseen scenario and in a previously unseen location? To answer this question, we introduce the Action Recognition Generalisation Over scenarios and locations dataset (ARGO1M), which contains 1.1M video clips from the large-scale Ego4D dataset, across 10 scenarios and 13 locations. We demonstrate recognition models struggle to generalise over 10 proposed test splits, each of an unseen scenario in an unseen location. We thus propose CIR, a method to represent each video as a Cross-Instance Reconstruction of videos from other domains. Reconstructions are paired with text narrations to guide the learning of a domain generalisable representation. We provide extensive analysis and ablations on ARGO1M that show CIR outperforms prior domain generalisation works on all test splits. Code and data: https://chiaraplizz.github.io/what-can-a-cook/.

Figure 1: Problem statement and samples from the ARGO1M dataset. The same action, e.g. "cut", is performed differently based on the scenario and the location in which it is carried out. We aim to generalise so as to recognise the same action within a new scenario, unseen during training, and in an unseen location, e.g., Mechanic in India.

ARGO1M Dataset
CIR Method
Code and Models

RELEASED

Plizzari et al (2023). What can a cook in Italy teach a mechanic in India? Action Recognition Generalisation Over Scenarios and Locations. IEEE/CVF International Conference on Computer Vision (ICCV).
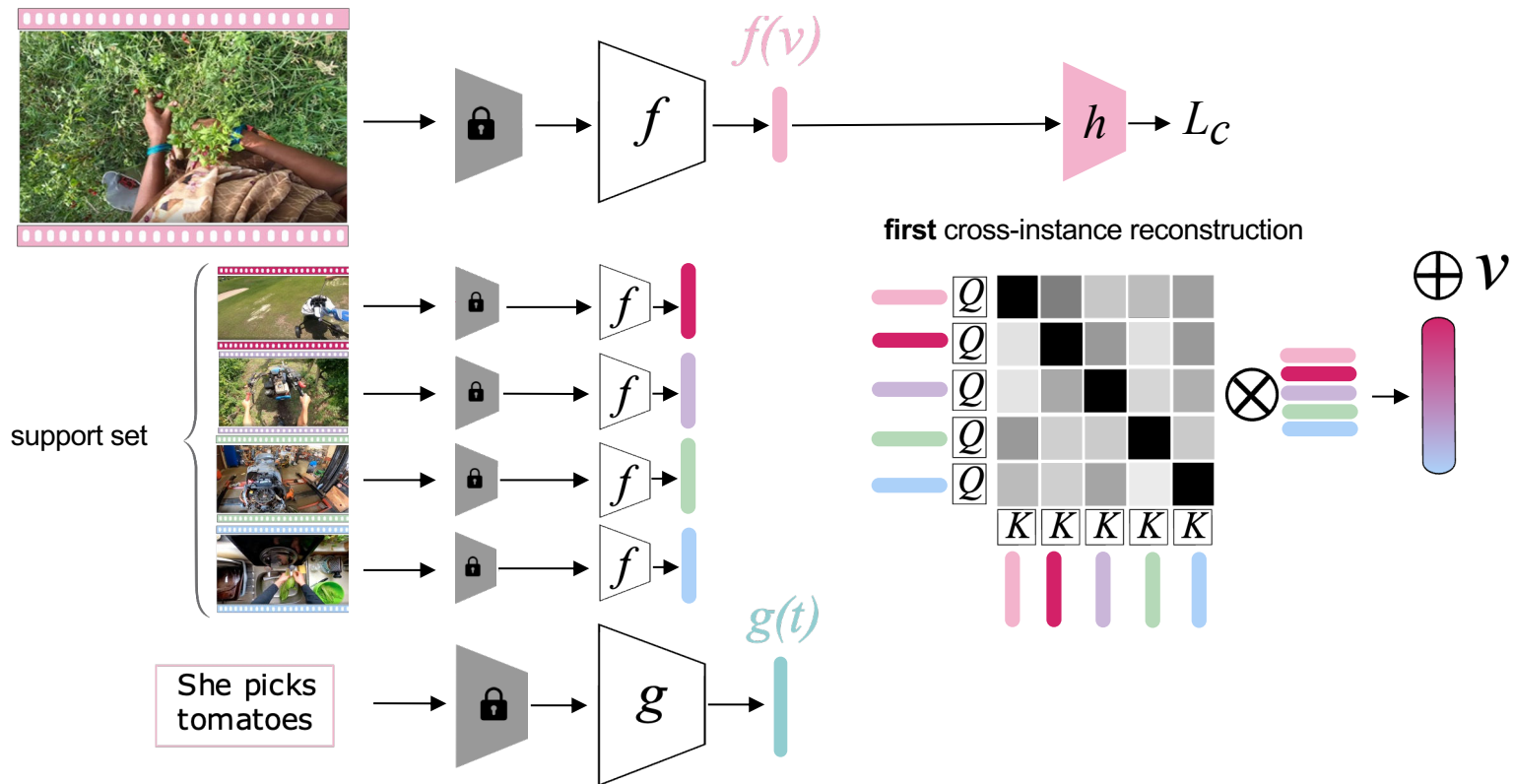
Dima Damen
WACV2024 – Waikoloa, Hawaii

57

University of BRISTOL

## Tasks are harder

Detection, 3D Mapping, Tracking, VOS, Hand-Object, Generative, …

## Solutions prove more rewarding

Weak supervision, Domain Adap/Gen., Audio-Visual, long-term understanding

- Hands transform objects….



♠ = avocado

**Input**     **peeled♠on chopping board**     **♠in a blender**     **♠ smoothie in a blender**

University of BRISTOL

T Soucek et al (2023). GenHowTo: Learning to Generate Actions and State Transformations from Instructional Videos. ArXiv

Dima Damen          59
WACV2024 – Waikoloa, Hawaii

with: Tomas Soucek    Michael Wray
Ivan Laptev    Josef Sivic



Input        GenHowTo        EF-DDPM        InstructPix2Pix

Prompt: a frosted cake with strawberries around the top

Prompt: a person kneading dough on a cutting board

Prompt: a person cutting a fish on a cutting board

University of
BRISTOL

- Two contributions…. Dataset & Method

- Two contributions…. **Dataset** & Method

**Instructional (HowTo) Video**



Initial State

Self-supervised temporal detection model 🔒

Action State

Final State

Image captioning model 🔒

action prompt $P_{ac}$

"a person cutting an avocado"

"avocado halves on a chopping board"

state prompt $P_{st}$

Tomas Soucek, Jean-Baptiste Alayrac, Antoine Miech, Ivan Laptev, and Josef Sivic (2022). Multi-task learning of object state changes from uncurated videos.

- Two contributions…. Dataset & <mark>Method</mark>



Input frame

VAE Enc. $A_E$

Control Net $U'_E$

$z_t$

Input Prompt P

"avocado halves on chopping board"

Transf.

Enc. $U_E$

Dec. $U_D$

VAE Dec. $A_D$

Stable Diffusion

Target frame

$t = T…1$

frozen
fine-tuned

→ text conditioning
↻ iterative generation

T Soucek et al (2023). GenHowTo: Learning to Generate Actions and State Transformations from Instructional Videos. ArXiv

with: Tomas Soucek   Michael Wray
Ivan Laptev   Josef Sivic

Input          *less noise* ⟵―――――――――――――――――⟶ *more noise*

- Qualitative Evaluation…

  ○ Initial vs Final State

  ○ Binary Classifier

| Method | $Acc_{ac}$ ↑ | $Acc_{st}$ ↑ |
|---|---|---|
| *test set categories unseen during training* | | |
| (a) Stable Diffusion | 0.51 | 0.50 |
| (b) Edit Friendly DDPM | 0.60 | 0.61 |
| (c) InstructPix2Pix | 0.55 | 0.63 |
| (d) *CLIP (manual prompts)* | 0.52 | 0.62 |
| (e) **GenHowTo** | **0.66** | **0.74** |
| *test set categories seen during training* | | |
| (f) Edit Friendly DDPM[†] | 0.69 | 0.80 |
| (g) **GenHowTo**[†] | **0.77** | **0.88** |
| (h) *Real images* | 0.96 | 0.97 |

[†] Models trained also on the test set *categories*.

with: Tomas Soucek    Michael Wray
        Ivan Laptev    Josef Sivic

a person is wrapping a tortilla on a plate

REAL IMAGE ———— GENERATED

a plate with two burritos on it

REAL IMAGE ———— GENERATED

a man pouring beer into a glass

REAL IMAGE ———— GENERATED

a man sitting at a table holding a glass of beer

REAL IMAGE ———— GENERATED

University of BRISTOL

# Tasks are harder

Detection, 3D Mapping, Tracking, VOS, Hand-Object, Generative, …

# Solutions prove more rewarding

Weak supervision, Domain Adap/Gen., Audio-Visual, long-term understanding

University of BRISTOL

Videos are multi-modal

# Multi-modal learning…

with: Vangelis Kazakos    Jaesung Huh
Arsha Nagrani.    Jacob Chalk
Andrew Zisserman

- The magic of audio-visual understanding…

- Object-Object interactions

with: Vangelis Kazakos    Jaesung Huh
Arsha Nagrani.    Jacob Chalk
Andrew Zisserman

- The magic of audio-visual understanding…

- Object-Object interactions

- Material sounds

University of BRISTOL

with: Vangelis Kazakos   Jaesung Huh
Arsha Nagrani.   Jacob Chalk
Andrew Zisserman

- The magic of audio-visual understanding…

- Object-Object interactions

- Material sounds

- Sound-emitting objects

with: Vangelis Kazakos
Arsha Nagrani
Andrew Zisserman

Harmonic Sounds

Percussive Sounds

EPIC-KITCHENS

E Kazakos, A Nagrani, A Zisserman, D Damen (2021). Slow-Fast Auditory Streams For Audio Recognition. ICASSP

with: Vangelis Kazakos
Arsha Nagrani
Andrew Zisserman

Harmonic Sounds

Percussive Sounds

VGG-Sound

Dima Damen                                    73
WACV2024 – Waikoloa, Hawaii

with: Vangelis Kazakos
Arsha Nagrani
Andrew Zisserman

# Auditory Slow-Fast

Outstanding Paper Award – ICASSP 2021

University of BRISTOL

# Audio Slow-Fast

with: Vangelis Kazakos
Arsha Nagrani
Andrew Zisserman

- Slow has low temporal precision and large amount of channels
- Fast has fewer channels but high temporal resolution
- Multi-level lateral connections
- Separable convolutions

E Kazakos, A Nagrani, A Zisserman, D Damen (2021). Slow-Fast Auditory Streams For Audio Recognition. ICASSP

University of BRISTOL

# Audio Slow-Fast

with: Vangelis Kazakos
Arsha Nagrani
Andrew Zisserman

| Slow stream | | Fast stream | |
|---|---|---|---|
| **Animals** | baltimore oriole calling | **Percussive sounds** | footsteps on snow |
| | cheetah chirrup | | snake rattling |
| | zebra braying | | tap dancing |
| | dinosaurs bellowing | | car engine knocking |
| | horse neighing | | woodpecker pecking tree |
| | black capped chickadee calling | | chopping wood |
| | cat hissing | | people clapping |
| | cuckoo bird calling | | lawn mowing |
| | mosquito buzzing | | typing on typewriter |
| | bull bellowing | | opening or closing car doors |
| | whale calling | | playing tennis |
| **Scenes** | volcano explosion | | railroad car |
| | playing lacrosse | | playing tympani |
| | hair dryer drying | | playing drum kit |
| | sea waves | | playing  vibraphone |
| | playing tympani | | popping pop corn |
| | blowtorch igniting | **Voices** | singing choir |
| | opening/closing electric car | | people cheering |
| | windows | | people crowd |
| | thunder | | child speech |
| | electric blender running | | baby laughter |
| | playing shofar | **Others** | cat purring |
| | airplane flyby | | dog barking |
| | playing trumpet | | race car |
| | wind chime | | singing bowl |
| | striking bowling | | vacuum cleaner cleaning floors |
| | | | toilet flushing |
| | | | dog growling |
| | | | splashing water |

University of BRISTOL

E Kazakos, A Nagrani, A Zisserman, D Damen (2021). Slow-Fast Auditory Streams For Audio Recognition. ICASSP

# Audio Slow-Fast

with: Vangelis Kazakos
Arsha Nagrani
Andrew Zisserman

| Slow stream | | Fast stream | |
|---|---|---|---|
| **Animals** | baltimore oriole calling<br>cheetah chirrup<br>zebra braying<br>dinosaurs bellowing<br>horse neighing<br>black capped chickadee calling<br>cat hissing<br>cuckoo bird calling<br>mosquito buzzing<br>bull bellowing<br>whale calling | **Percussive sounds** | footsteps on snow<br>snake rattling<br>tap dancing<br>car engine knocking<br>woodpecker pecking tree<br>chopping wood<br>people clapping<br>lawn mowing<br>typing on typewriter<br>opening or closing car doors<br>playing tennis<br>railroad car<br>playing tympani<br>playing drum kit<br>playing vibraphone<br>popping pop corn |
| **Scenes** | volcano explosion<br>playing lacrosse<br>hair dryer drying<br>sea waves<br>playing tympani<br>blowtorch igniting<br>opening/closing electric car windows<br>thunder<br>electric blender running<br>playing shofar<br>airplane flyby<br>playing trumpet<br>wind chime<br>striking bowling | **Voices** | singing choir<br>people cheering<br>people crowd<br>child speech<br>baby laughter |
| | | **Others** | cat purring<br>dog barking<br>race car<br>singing bowl<br>vacuum cleaner cleaning floors<br>toilet flushing<br>dog growling<br>splashing water |

University of BRISTOL

E Kazakos, A Nagrani, A Zisserman, D Damen (2021). Slow-Fast Auditory Streams For Audio Recognition. ICASSP

**TOWARDS LEARNING UNIVERSAL AUDIO REPRESENTATIONS**

*Luyu Wang, Pauline Luc, Yan Wu, Adrià Recasens, Lucas Smaira, Andrew Brock, Andrew Jaegle,*

**Table 2**: **Evaluating frameworks and architectures on HARES.** We compare the impact of architecture choice under the classification and SimCLR objective. We also show the performance of several other recent strongly performing frameworks. Average scores are reported for tasks in each domain separately, and all three combined. All models are trained on AudioSet except for bidirectional CPC and Wav2Vec2.0, for which we also show results when they are trained on LibriSpeech (LS).

| Architecture | #Params | Input format | Used in | Env. | Speech | Music | HARES | AudioSet (mAP) |
|---|---|---|---|---|---|---|---|---|
| | | | | *Classification/SimCLR* | | | | |
| BYOL-A CNN | 5.3m | Spectrogram | [9] | 69.4/69.9 | 61.4/69.8 | 57.6/63.1 | 63.1/68.2 | 32.2/32.2 |
| EfficientNet-B0 | 4.0m | Spectrogram | [8] | 71.1/63.8 | 43.5/40.7 | 48.0/44.0 | 53.8/49.2 | 34.5/26.2 |
| CNN14 | 71m | Spectrogram | [11, 13] | 74.6/66.4 | 56.0/37.3 | 56.4/44.8 | 62.3/48.9 | 37.8/28.8 |
| ViT-Base | 86m | Spectrogram | [12] | 73.3/74.6 | 50.4/56.5 | 60.3/64.2 | 60.5/64.5 | 36.8/36.8 |
| ResNet50 | 23m | Spectrogram | [19] | 74.8/74.4 | 51.7/65.0 | 59.6/63.7 | 61.4/67.8 | 38.4/36.2 |
| SF ResNet50 | 26m | Spectrogram | [17] | 74.0/74.3 | 56.9/73.4 | 59.6/65.2 | 63.3/71.7 | 37.2/36.6 |
| NFNet-F0 | 68m | Spectrogram | Ours | **76.1**/76.0 | 59.0/65.9 | 61.8/65.5 | 65.4/69.2 | **39.3**/37.6 |
| SF NFNet-F0 | 63m | Spectrogram | Ours | 75.2/75.8 | 65.6/**77.2** | 64.5/**68.6** | 68.5/**74.6** | 38.2/37.8 |

achieve state-of-the-art performance across all domains.

**Index Terms—** audio representations, representation evaluation, speech, music, acoustic scenes

111.12

supervised contrastive learning [8, 15], and comparing them across a large set of model architectures. We find that models trained with contrastive learning tend to generalize better in the speech and music domain, while performing comparably to supervised pretraining for environment sounds. We

University of BRISTOL

Dima Damen                    78
WACV2024 – Waikoloa, Hawaii

with: Jaesung Huh*   & Jacob Chalk*
Vangelis Kazakos    Andrew Zisserman

# EPIC-Sounds: A Large-scale Dataset of Actions That Sound

Jaesung Huh*, Jacob Chalk*, Evangelos Kazakos, Dima Damen, Andrew Zisserman
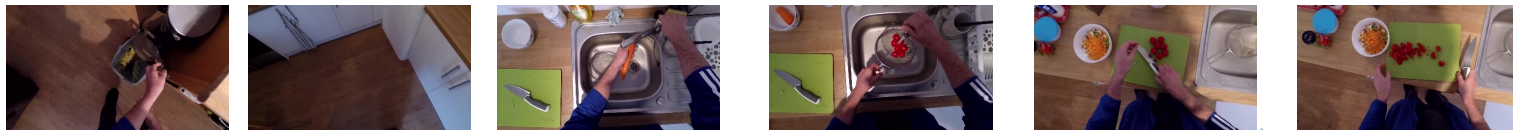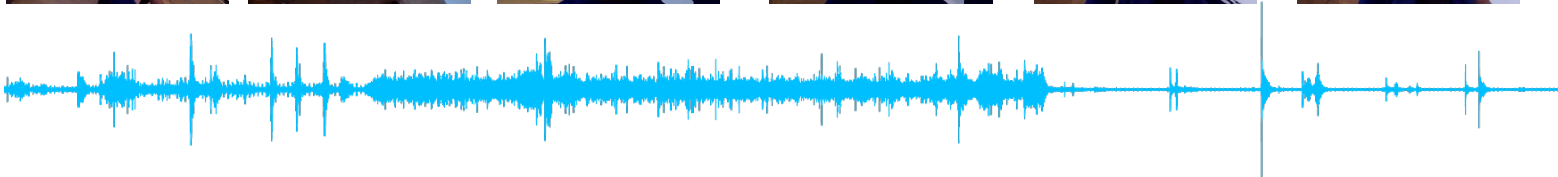* : Equal contribution

WACV2024 – Waikoloa, Hawaii

with: Jaesung Huh*   & Jacob Chalk*
Vangelis Kazakos   Andrew Zisserman



Video

Audio

J Huh*, J Chalk* E Kazakos, D Damen, A Zisserman (2023). EPIC-Sounds: A Large-scale Dataset of Actions That Sound. ICASSP

with: Jaesung Huh*   & Jacob Chalk*
Vangelis Kazakos    Andrew Zisserman

Video

Close bin   Close bag          Wash carrot          Wash tomato        Take knife      Cut tomato

Audio

Dima Damen                                    81
WACV2024 – Waikoloa, Hawaii

with: Jaesung Huh*   & Jacob Chalk*
Vangelis Kazakos   Andrew Zisserman



Video

Close bin   Close bag          Wash          Wash tomato       Take knife    Cut tomato

**Incorrect assumption**

Audio

J Huh*, J Chalk* E Kazakos, D Damen, A Zisserman (2023). EPIC-Sounds: A Large-scale Dataset of Actions That Sound. ICASSP

Dima Damen                   82
WACV2024 – Waikoloa, Hawaii

# Motivation

with: Jaesung Huh*   & Jacob Chalk*
Vangelis Kazakos   Andrew Zisserman



Video

Close bin   Close bag        Wash carrot        Wash tomato        Take knife   Cut tomato

Audio

Dima Damen
WACV2024 – Waikoloa, Hawaii

83

with: Jaesung Huh*   & Jacob Chalk*
Vangelis Kazakos    Andrew Zisserman



Video

Close bin    Close bag                Wash carrot                Wash tomato            Take knife        Cut tomato

Audio

J Huh*, J Chalk* E Kazakos, D Damen, A Zisserman (2023). EPIC-Sounds: A Large-scale Dataset of Actions That Sound. ICASSP

University of BRISTOL

with: Jaesung Huh*   & Jacob Chalk*
Vangelis Kazakos   Andrew Zisserman

Video

Close bin   Close bag          Wash carrot          Wash tomato       Take knife     Cut tomato

Audio

University of BRISTOL

with: Jaesung Huh*   & Jacob Chalk*
Vangelis Kazakos    Andrew Zisserman

Video

Audio

Audio labels

Open/close   Sniffing  Footsteps          Metal              Metal

Rustle                          Tap running

Cut food

Dima Damen                                    86
WACV2024 – Waikoloa, Hawaii

with: Jaesung Huh*   & Jacob Chalk*
Vangelis Kazakos    Andrew Zisserman

Dima Damen
WACV2024 – Waikoloa, Hawaii

# EPIC-SOUNDS

with: Jaesung Huh*   & Jacob Chalk*
Vangelis Kazakos    Andrew Zisserman

**EPIC-KITCHENS VIDEOS**

100 hours

45 kitchens

**Visual Action Annotations**

90K visual actions

97 verb classes

300 noun classes

**EPIC-Sounds**

Audio-Based Annotations

79K categorised audio events
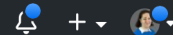
44 sound categories

39K uncategorised events

J Huh*, J Chalk* E Kazakos, D Damen, A Zisserman (2023). EPIC-Sounds: A Large-scale Dataset of Actions That Sound. ICASSP

University of BRISTOL

Dima Damen
WACV2024 – Waikoloa, Hawaii

88

spray

89

# Annotations Pipeline

with: Jaesung Huh*  & Jacob Chalk*
Vangelis Kazakos   Andrew Zisserman

- We annotate all the distinctive sound events which consist of temporal intervals using free-form sound descriptions.
- Using VGG VIA annotation tool



J Huh*, J Chalk* E Kazakos, D Damen, A Zisserman (2023). EPIC-Sounds: A Large-scale Dataset of Actions That Sound. ICASSP

University of BRISTOL

Dima Damen
WACV2024 – Waikoloa, Hawaii

- From free-form descriptions to categories

University of BRISTOL

J Huh*, J Chalk* E Kazakos, D Damen, A Zisserman (2023). EPIC-Sounds: A Large-scale Dataset of Actions That Sound. ICASSP

Dima Damen
WACV2024 – Waikoloa, Hawaii

92

with: Jaesung Huh*   & Jacob Chalk*
Vangelis Kazakos    Andrew Zisserman

- For collision sounds, we annotate the **materials** of the objects that colliding.
- Materials example



Ceramic             Cloth             Metal             Plastic             Glass

University of BRISTOL

with: Jaesung Huh*   & Jacob Chalk*
Vangelis Kazakos   Andrew Zisserman

- Manual check on validation / test set

- We use the overlaps between audio and visual segments for reviewing train set.

J Huh*, J Chalk* E Kazakos, D Damen, A Zisserman (2023). EPIC-Sounds: A Large-scale Dataset of Actions That Sound. ICASSP

# EPIC-SOUNDS

with: Jaesung Huh*  & Jacob Chalk*
Vangelis Kazakos    Andrew Zisserman

**EPIC-KITCHENS VIDEOS**
100 hours
45 kitchens

**Visual Action Annotations**
90K visual actions
97 verb classes
300 noun classes

**EPIC-Sounds**
Audio-Based Annotations
79K categorised audio events
44 sound categories
39K uncategorised events

University of BRISTOL

epic-kitchens / **epic-sounds-annotations**   Public

Edit Pins ⌄    Unwatch 5 ⌄    Fork 3 ⌄    Starred 47 ⌄

<> **Code**    Issues 1    Pull requests    Actions    Projects    Wiki    Security    Insights    Settings

111 lines (91 sloc) | 10.3 KB

<>   Raw   Blame

# EPIC-SOUNDS Dataset

We introduce EPIC-SOUNDS, a large scale dataset of audio annotations capturing temporal extents and class labels within the audio stream of the egocentric videos from EPIC-KITCHENS-100. EPIC-SOUNDS includes 78.4k categorised and 39.2k non-categorised segments of audible events and actions, distributed across 44 classes. In this repository, we provide labelled temporal timestamps for the train / val split, and just the timestamps for the recognition test split. We also provided the temporal timestamps for annotations that could not be clustered into one of our 44 classes, along with the free-form description used during the initial annotation. We train and evaluate two state-of-the-art audio recognition models on our dataset, which we also provide the code and pretrained models for.

## Download the Data

A download script is provided for the videos here. You will have to extract the untrimmed audios from these videos. Instructions on how to extract and format the audio into a HDF5 dataset can be found on the Auditory SlowFast GitHub repo. Alternatively, you can email uob-epic-kitchens@bristol.ac.uk for access to an existing HDF5 file.

Contact: uob-epic-kitchens@bristol.ac.uk

## Citing

When using the dataset, kindly reference our ICASSP 2023 Paper:

# Tasks are harder

Detection, 3D Mapping, Tracking, ==VOS==, Hand-Object, Generative, …

# Solutions prove more rewarding

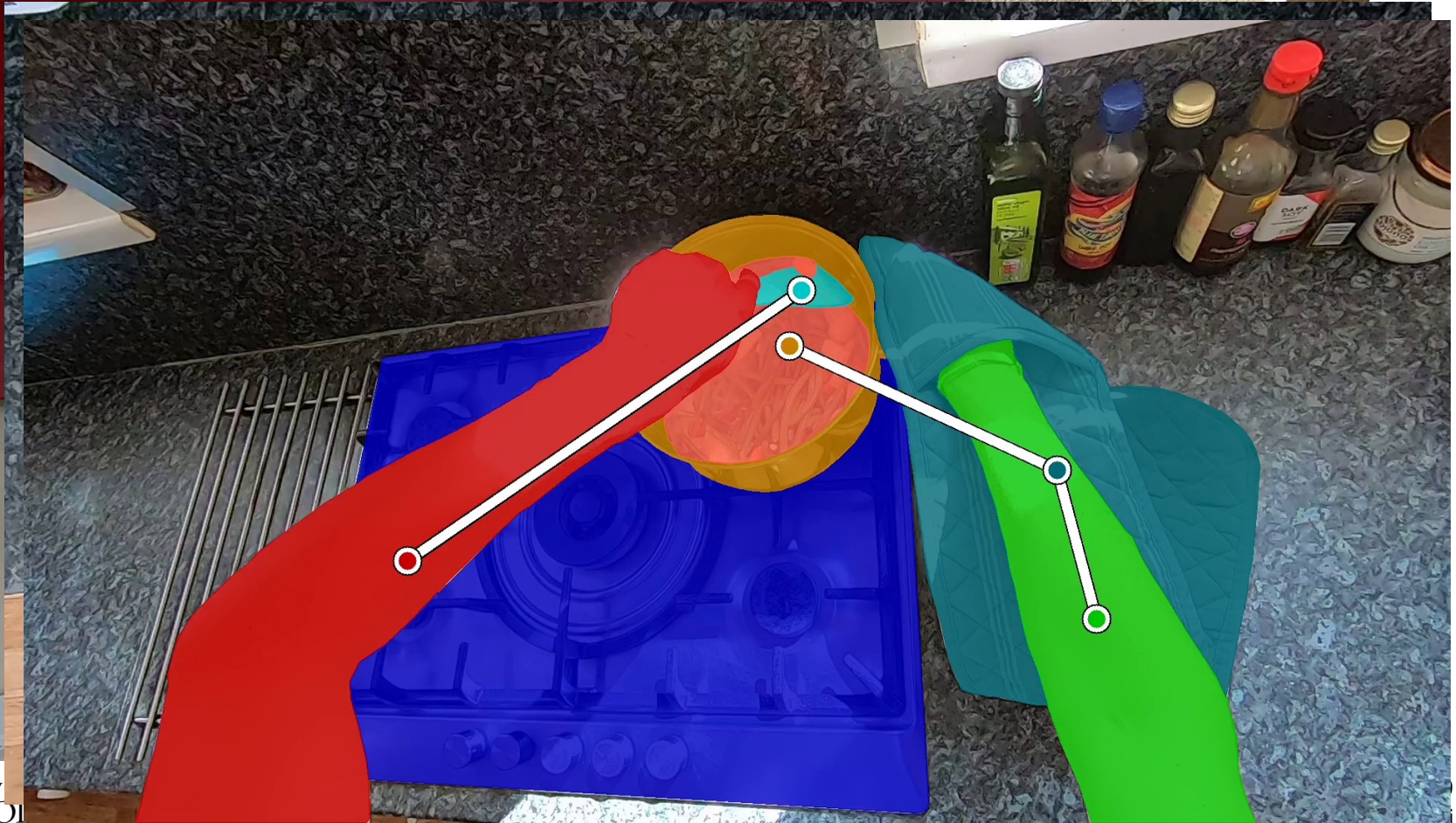Weak supervision, Domain Adap/Gen., Audio-Visual, long-term understanding
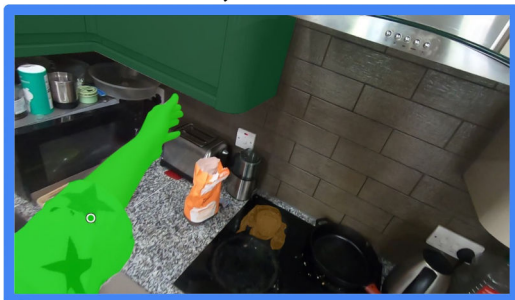
with: Ahmad Darkhalil, Dandan Shan, Bin Zhu, Jian Ma,
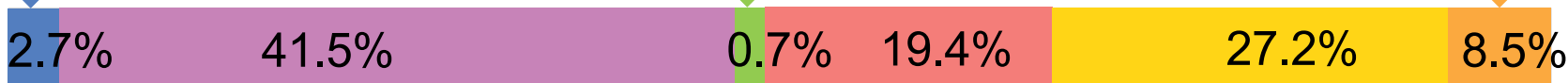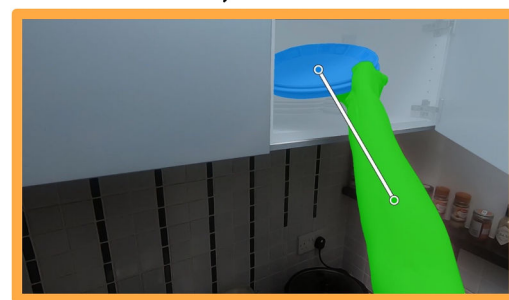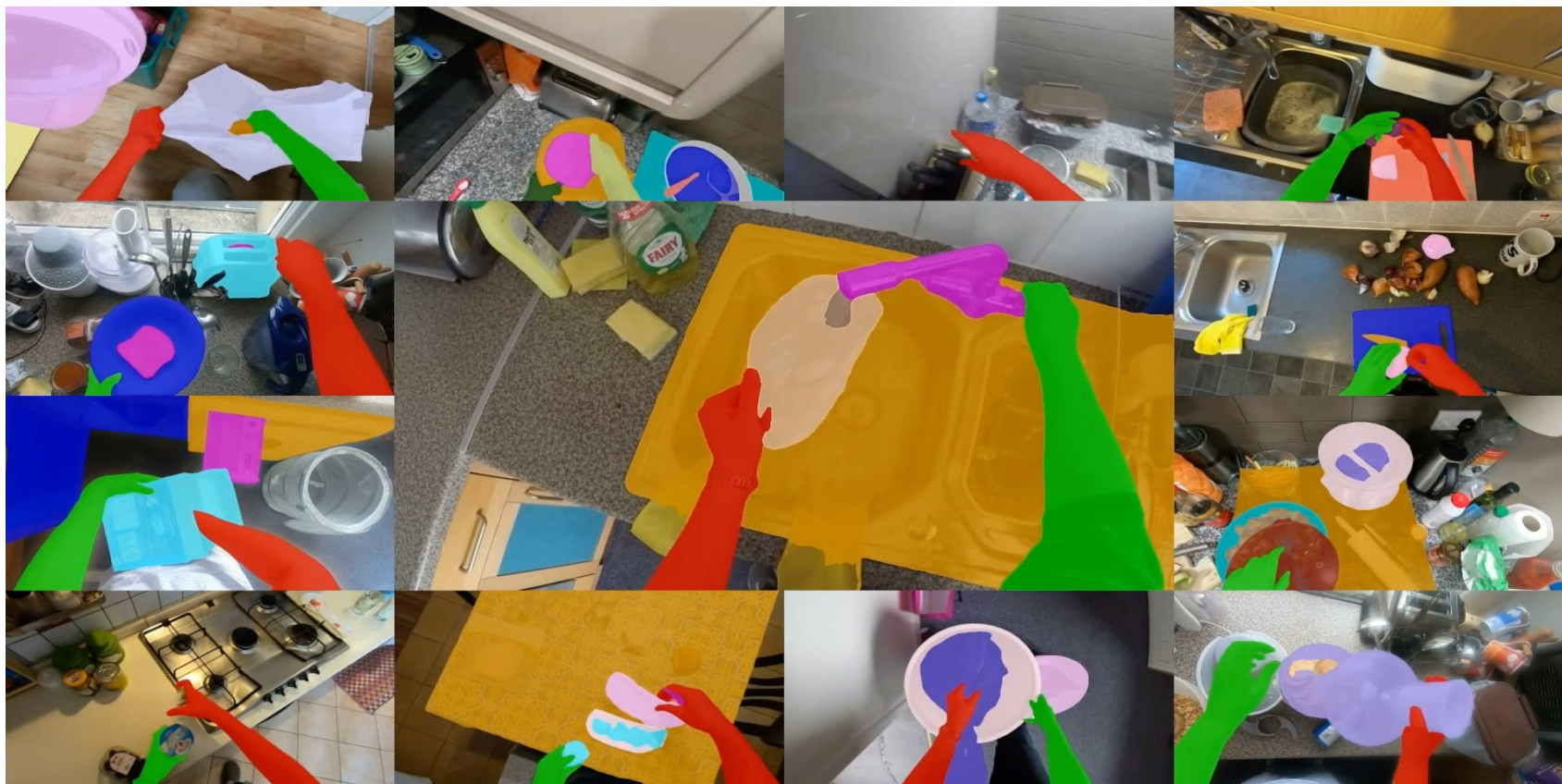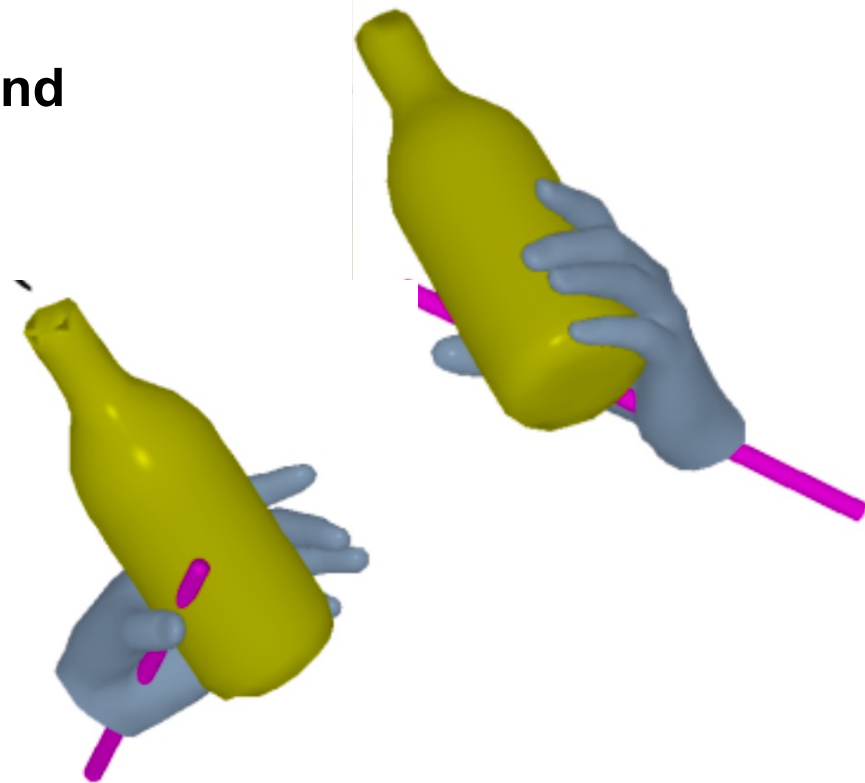Amlan Kar, Richard Higgins, David Fouhey, Sanja Fidler

VISOR annotates videos from EPIC-KITCHENS

University of BRISTOL

with: Ahmad Darkhalil, Dandan Shan, Bin Zhu, Jian Ma, Amlan Kar, Richard Higgins, David Fouhey, Sanja Fidler, Dima Damen

University BRISTOL

# Object relation stats

with: Ahmad Darkhalil, Dandan Shan, Bin Zhu, Jian Ma, Amlan Kar, Richard Higgins, David Fouhey, Sanja Fidler, Dima Damen

1 Hand, No Contact

2 Hands, No Contact

1 Hand, In Contact

2.7%    41.5%    0.7%    19.4%    27.2%    8.5%

2 Hands, 2 Obj Contacts

2 Hands, Same Contact

2 Hands, 1 In Contact

A Darkhalil et al (2022). EPIC-KITCHENS VISOR Benchmark: VIdeo Segmentations and Object Relations. *NeurIPS*

University of BRISTOL

Dima Damen
WACV2024 – Waikoloa, Hawaii

100

with: Ahmad Darkhalil, Dandan Shan, Bin Zhu, Jian Ma, Amlan Kar, Richard Higgins, David Fouhey, Sanja Fidler

with: Ahmad Darkhalil, Dandan Shan, Bin Zhu, Jian Ma, Amlan Kar, Richard Higgins, David Fouhey, Sanja Fidler

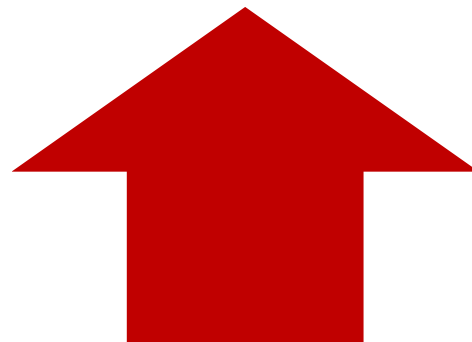University of BRISTOL

Dima Damen
WACV2024 – Waikoloa, Hawaii

# Tasks are harder

Detection, 3D Mapping, Tracking, VOS, ==Hand-Object==, Generative, …
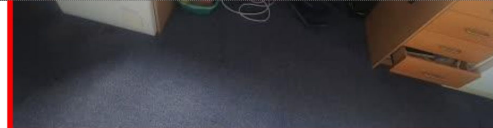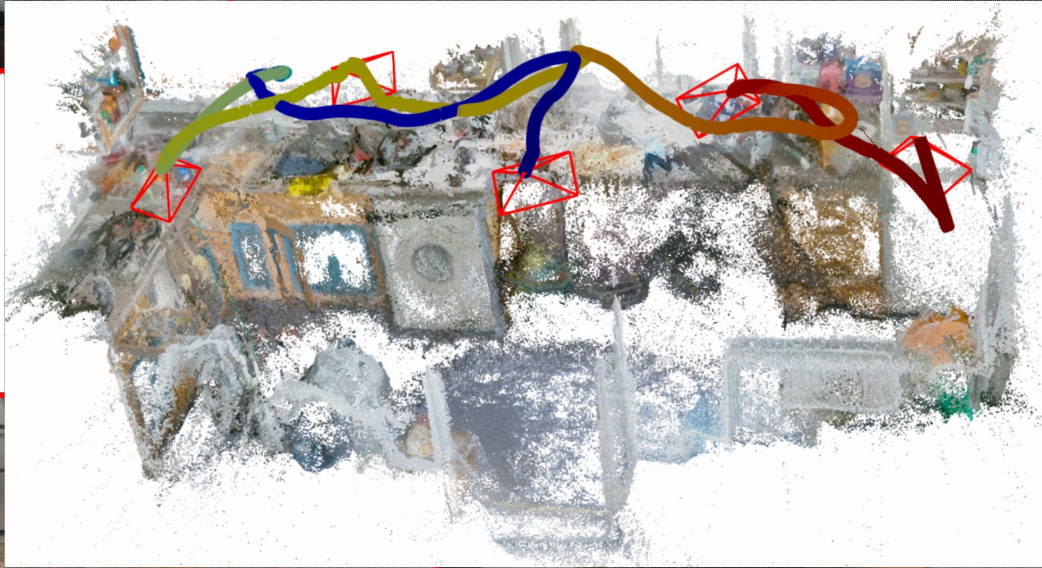
# Solutions prove more rewarding

Weak supervision, Domain Adap/Gen., Audio-Visual, long-term understanding

**left hand**

**bottle**

University of BRISTOL

# Get a Grip

with: Zhifan Zhu

**Non-Ego Views**

**Ego Views**

Invisible Fingers

Z Zhu and D Damen (2023). Get a Grip: Reconstructing Hand-Object Stable Grasps in Egocentric Videos. *ArXiv*

University of BRISTOL

Z Zhu and D Damen (2023). Get a Grip: Reconstructing Hand-Object Stable Grasps in Egocentric Videos. *ArXiv*

with: Zhifan Zhu



contact    grasp          end of grasp     release

Stable Grasp

free hand
standing object

grasping

University of
BRISTOL

Dima Damen
WACV2024 – Waikoloa, Hawaii

with: Zhifan Zhu



Stable Contact Area — IOU=83%    IOU=29%

Finger Pose

(actual v.s. residual) — actual, residual — act.=10° act.=40° res.=3° res.=27°

Objects within the stable grasp move within 1 DoF

Z Fan, O Taheri, D ...las, M Kocabas, M Kaufmann, M J Black, and O Hilliges (2023).
ARCTIC: A dataset for dexterous bimanual hand- object manipulation. CVPR

(left hand) Outside Grasp

with: Zhifan Zhu

**Input**



$t$

University of BRISTOL

Z Zhu and D Damen (2023). Get a Grip: Reconstructing Hand-Object Stable Grasps in Egocentric Videos. *ArXiv*

University of BRISTOL

## Tasks are harder

Detection, 3D Mapping, Tracking, VOS, Hand-Object, Generative, …

## Solutions prove more rewarding

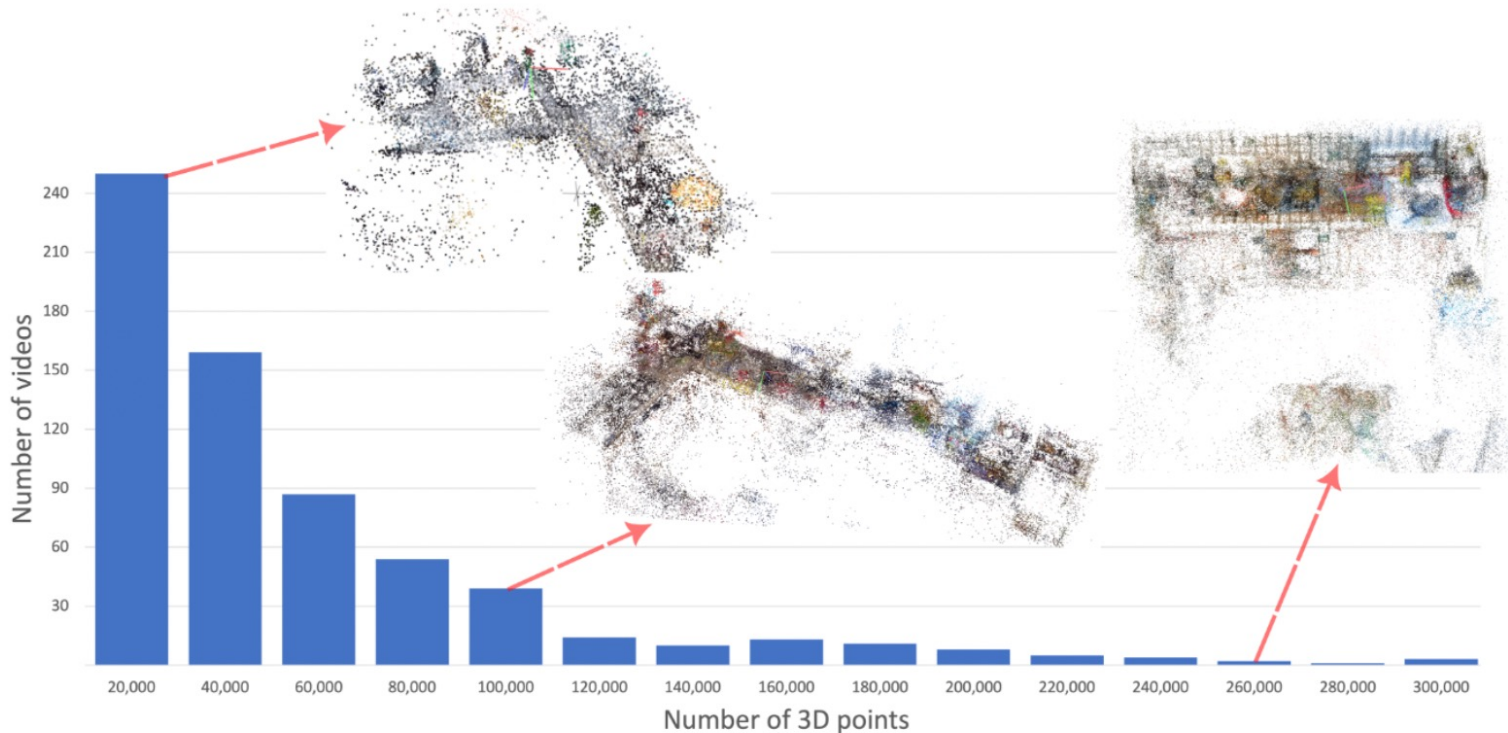Weak supervision, Domain Adap/Gen., Audio-Visual, long-term understanding

with: V Tschernezki*, A Darkhalil*, Z Zhu*,
D Fouhey, I Laina, D Larlus, A Vedaldi

EPIC-KITCHENS

Figure 4: **Number of 3D points histogram.** The majority of our reconstructions generate less than 40,000 points that are enough to represent the kitchen. However, some reconstructions have more than 100,000, we include the point clouds for each points range showing the fine details covered by having more points

Table 1: Comparison of datasets commonly used in dynamic new-view synthesis.

| Dataset | #Scenes | Seq. Length | Monocular | Semantics |
|---|---|---|---|---|
| Nerfies [37] | 4 | 8–15 sec | - | - |
| D-NeRF [41] | 8 | 1–3 sec | - | - |
| Plenoptic Video [22] | 6 | 10-60 sec | - | - |
| NVIDIA Dynamic Scene Dataset [65] | 12 | 1–5 sec | 4 / 12 | - |
| HyperNeRF [38] | 16 | 8–15 sec | 13 / 16 | - |
| iPhone [13] | 14 | 8–15 sec | 7 / 14 | - |
| SAFF [25] | 8 | 1–5sec | - | ✓ |
| **EPIC Fields** (ours) | 50 | 6–37 min (Avg 22) | 50 / 50 | ✓ |

with: Kristen Grauman
+102 authors

Dima Damen
WACV2024 – Waikoloa, Hawaii

with: Kristen Grauman
+102 authors



Ego-Exo Relation

K Grauman et al (2023). Ego-Exo4D: Understanding Skilled Human Activity from First- and Third-Person Perspectives.. *ArXiv*

University of BRISTOL

Dima Damen
WACV2024 – Waikoloa, Hawaii

with: Kristen Grauman
+102 authors



Ego Pose

Aria RGB

Dima Damen
WACV2024 – Waikoloa, Hawaii

cam01

cam04

# EGO-EXO4D

A **diverse**, large-scale **multi-modal**, **multi-view**, video dataset and benchmark collected across 13 cities worldwide by 839 camera wearers, capturing **1422 hours** of video of skilled human activities.

Learn More ↓    Watch Video ↗    Start Here ↗

# 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 degli STUDI di CATANIA

University of BRISTOL

Dima Damen
WACV2024 – Waikoloa, Hawaii

with: Chiara Plizzari, Gabriele Goletto, Antonio Furnari, Siddhant Bansal, Francesco Ragusa, Giovanni Maria Farinella, Tatiana Tommasi

# **Envisioning** an Ambitious Future and **Analysing** the Current Status of Egocentric Vision

How did we do this?

University of BRISTOL

Dima Damen
WACV2024 – Waikoloa, Hawaii

with: Chiara Plizzari, Gabriele Goletto, Antonio Furnari, Siddhant Bansal, Francesco Ragusa, Giovanni Maria Farinella, Tatiana Tommasi

# We imagined a device – *EgoAI* and envisioned its utility in multiple scenarios



**EGO-Designer**

**EGO-Worker**

**EGO-Tourist**

**EGO-Home**

**Ego-Police**

University of BRISTOL

Dima Damen
WACV2024 – Waikoloa, Hawaii

# EGO-Home

with: Chiara Plizzari, Gabriele Goletto, Antonio Furnari, Siddhant Bansal, Francesco Ragusa, Giovanni Maria Farinella, Tatiana Tommasi

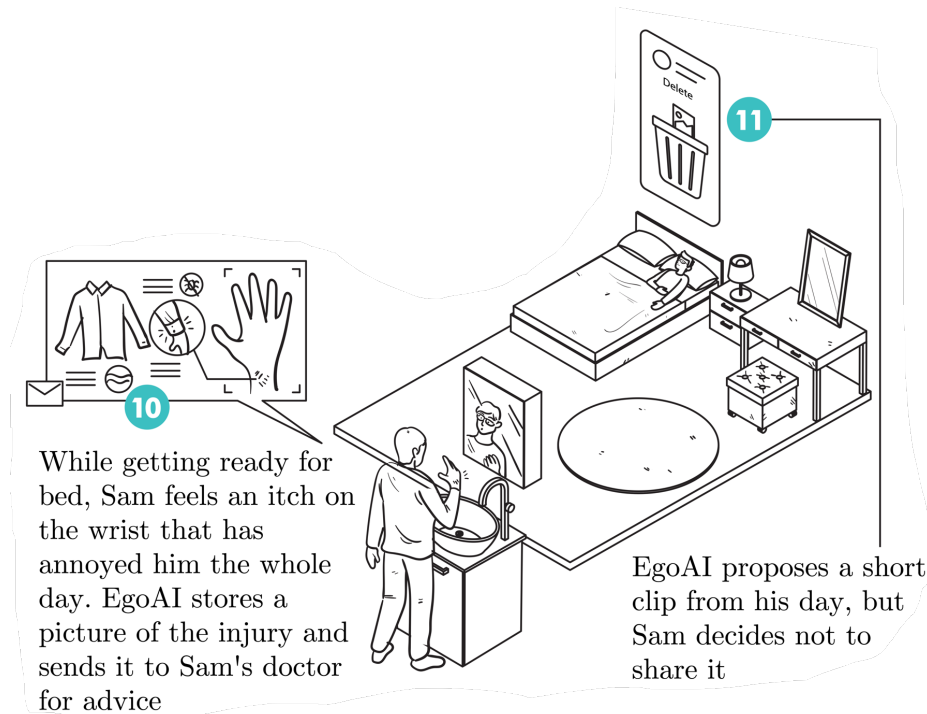**8** Waves hitting the shore look and sound natural

**7** Transferred to a beach he visited last summer

**9** After dinner, Sam enjoys a group card game with his friends, who are connected through their own EgoAI
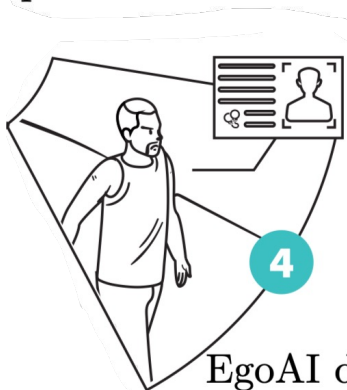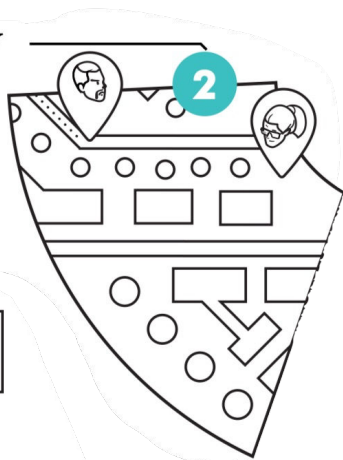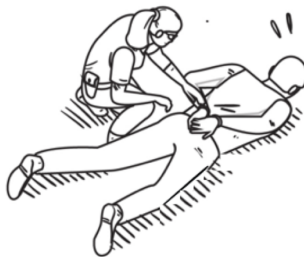
EGO-Home

with: Chiara Plizzari, Gabriele Goletto, Antonio Furnari, Siddhant Bansal, Francesco Ragusa, Giovanni Maria Farinella, Tatiana Tommasi

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

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

# EGO-Home

with: Chiara Plizzari, Gabriele Goletto, Antonio Furnari, Siddhant Bansal, Francesco Ragusa, Giovanni Maria Farinella, Tatiana Tommasi

EgoAI helps Judy navigate through the shortest safe path to target places

**2**

EgoAI detected and re-identi-fied the man before he passed Judy

**4**

**EGO-Police**

Localisation and Navigation — **1** **2**

Messaging — **1** **3** **11**

Action Recognition — **2** **13**

Person Re-ID — **2** **4**

Object Detection and Retrieval — **7**

Measuring System — **8** **9**

Decision Making — **9**

3D Scene Understanding — **10**

Hand-Object Interaction — **12**

Summarisation — **13**

Privacy — **14**

University of BRISTOL

# EGO-Home

with: Chiara Plizzari, Gabriele Goletto, Antonio Furnari, Siddhant Bansal, Francesco Ragusa, Giovanni Maria Farinella, Tatiana Tommasi
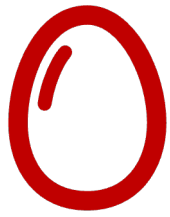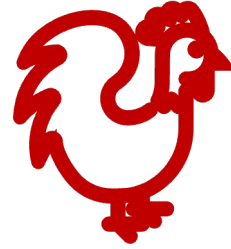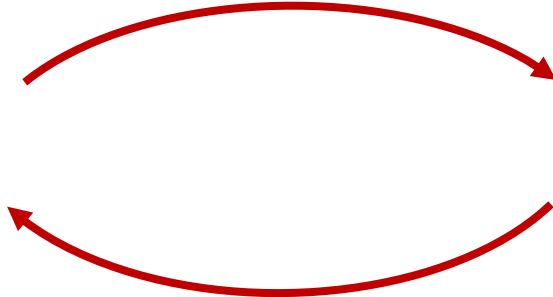
- 12 tasks

  - Seminal Works

  - SOTA methods

  - Datasets

  - Future Perspective
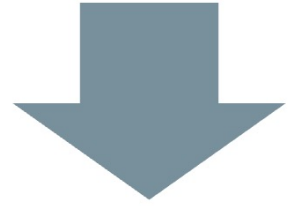
- 44 pages

- 385 references
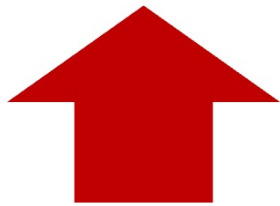
**Labels**

**Data**

EGO-EXO**4D**

Tasks are harder

Detection, 3D Mapping, Tracking, VOS, Hand-Object, Generative, ...

Solutions prove more rewarding

Weak supervision, Domain Adap/Gen., Audio-Visual, long-term understanding

University of BRISTOL

Dima Damen
WACV2024 – Waikoloa, Hawaii

# Visit us…

## 2024 Summer of Research@Bristol

### Enjoy a research-oriented and fun summer!

Get supervision and support to develop cutting-edge research projects with a three-month internship with one of the leading researchers in Machine Learning and Computer Vision at the University of Bristol.

Open to PhD students in any related discipline.



### Click to apply online by:
### Fri 19 Jan 2024

# Thank you

For further info, datasets, code, publications…

http://dimadamen.github.io

@dimadamen

http://www.linkedin.com/in/dimadamen

# Q&A

University of
BRISTOL