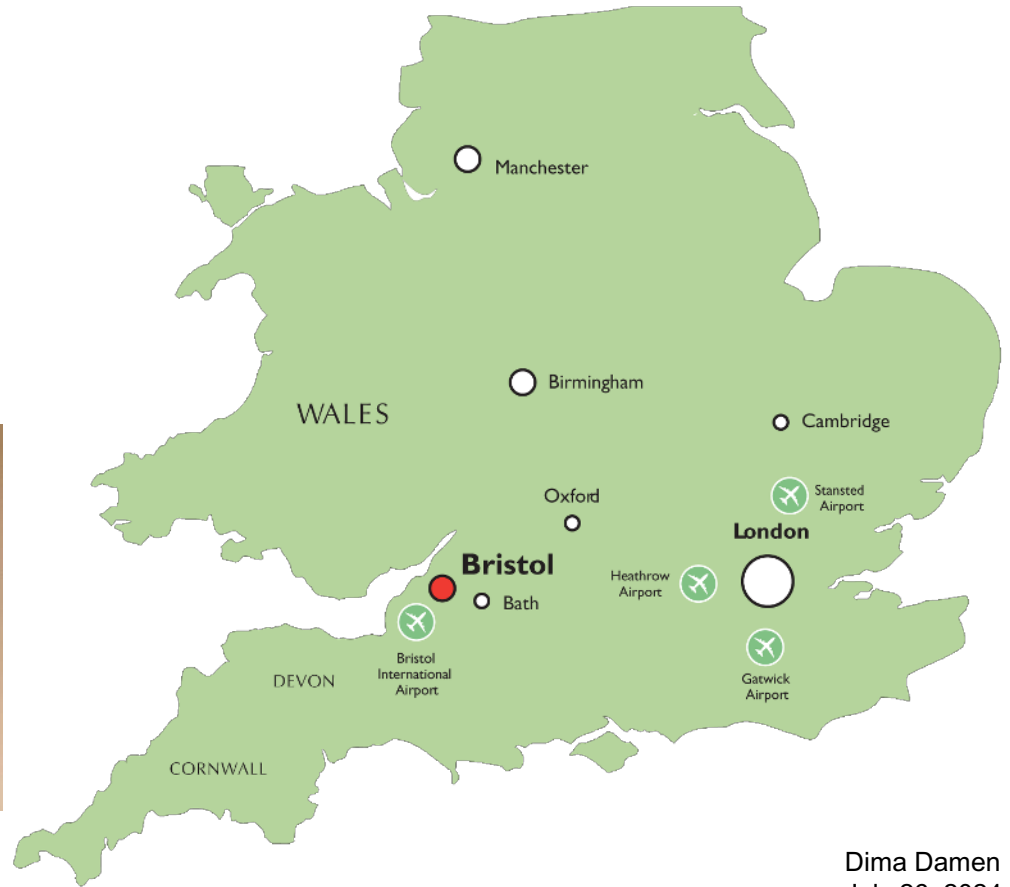




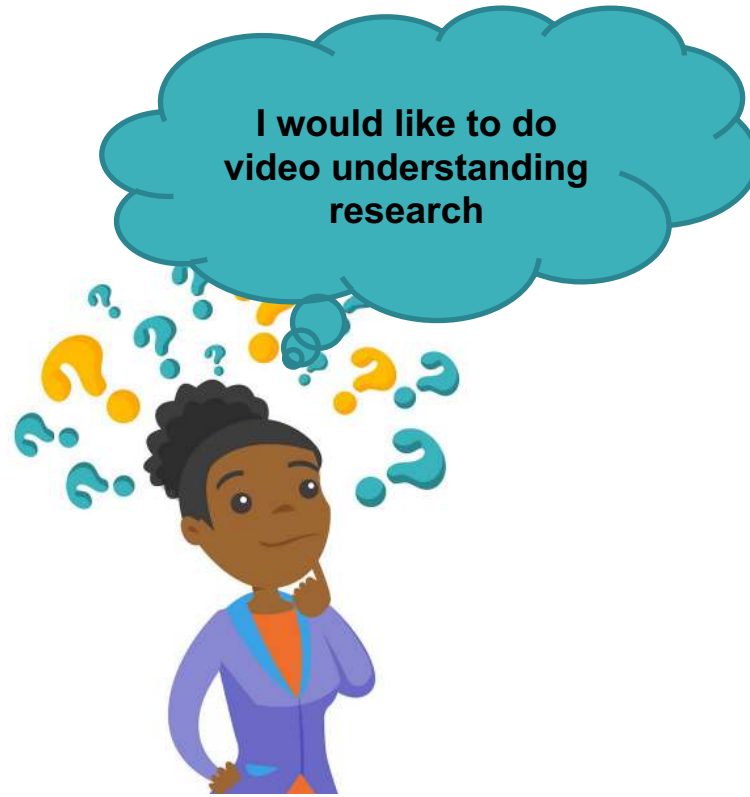
Video Understanding

An Egocentric Perspective

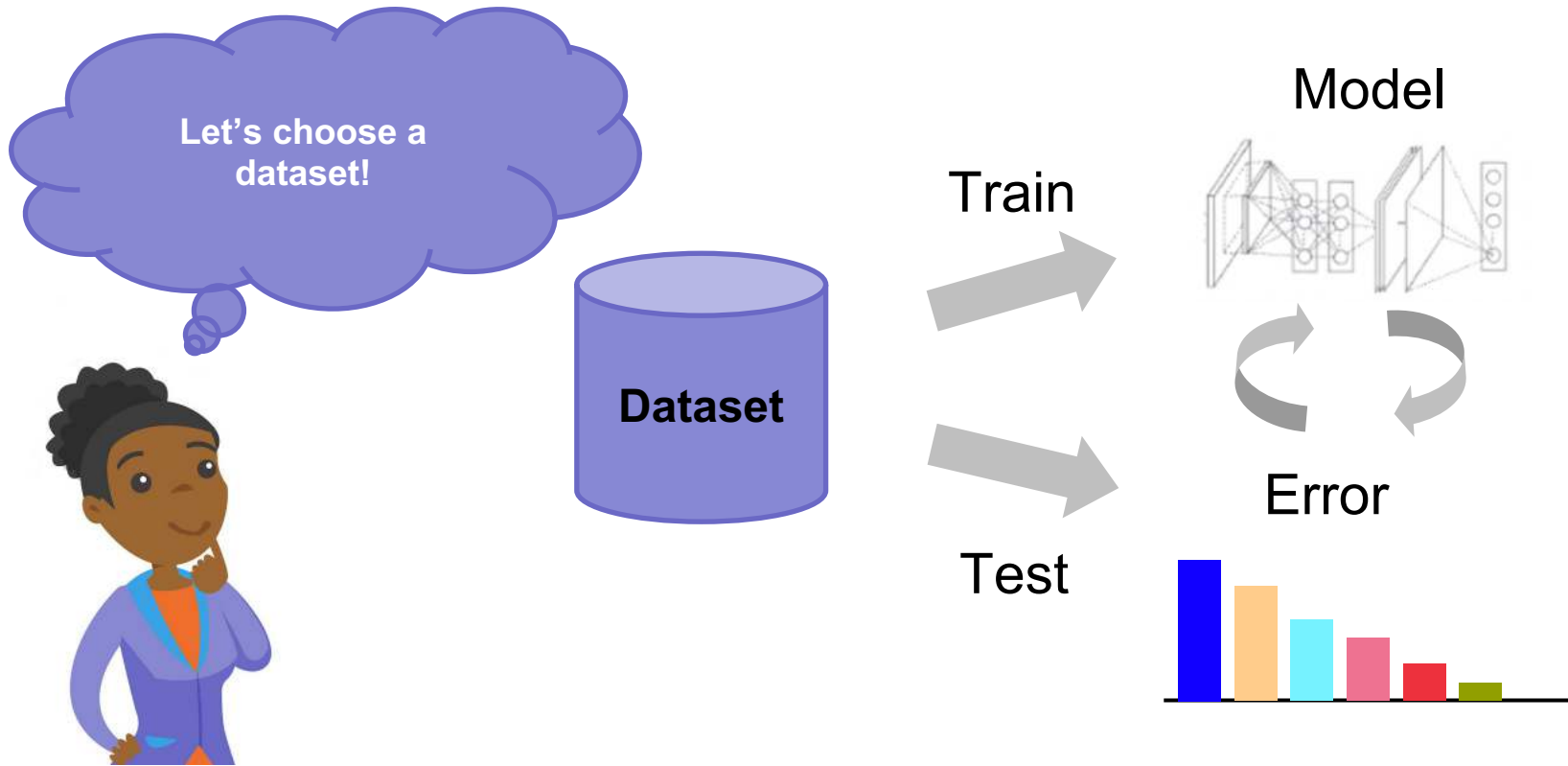
Introduction...



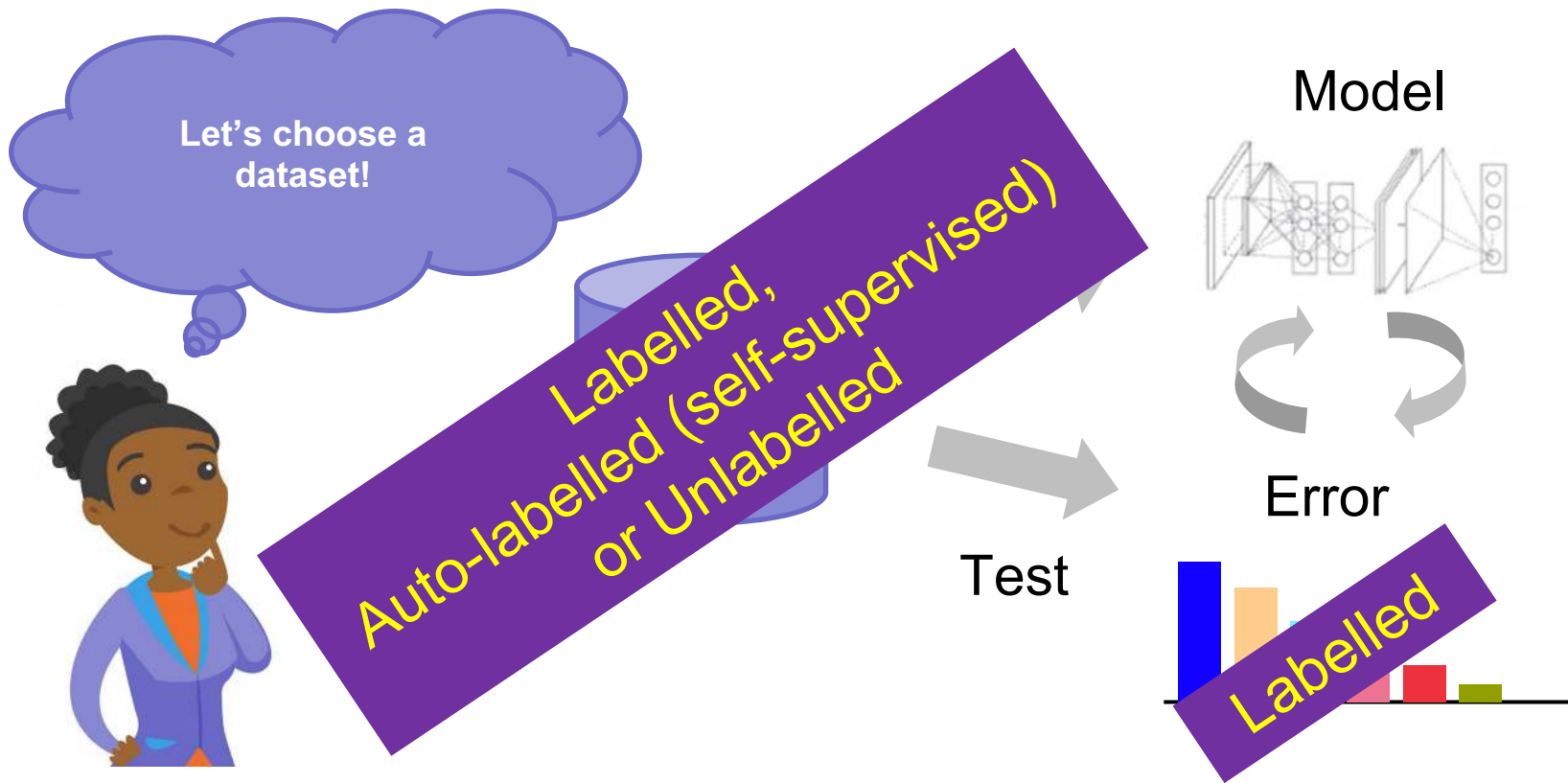
The current paradigm of Computer Vision Research



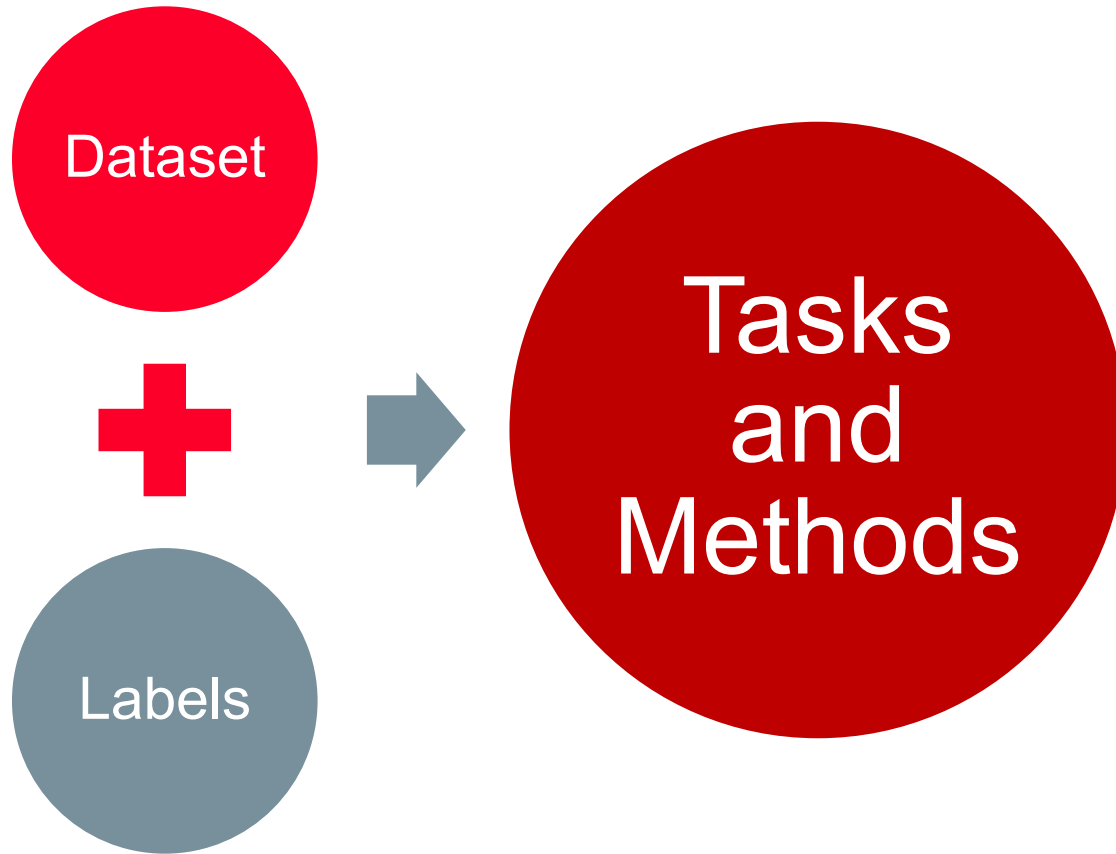
The current paradigm of Computer Vision Research

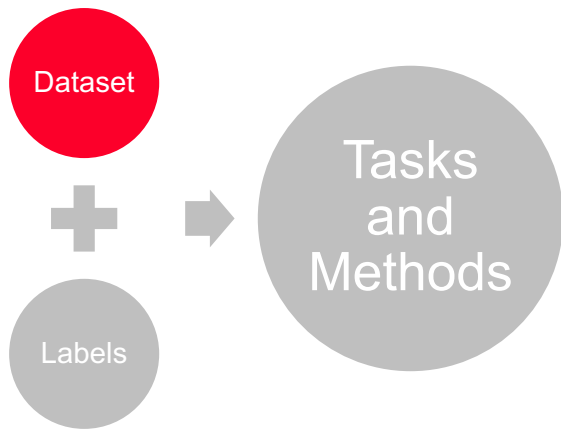


The current paradigm of Computer Vision Research



In this talk... on Video Understanding



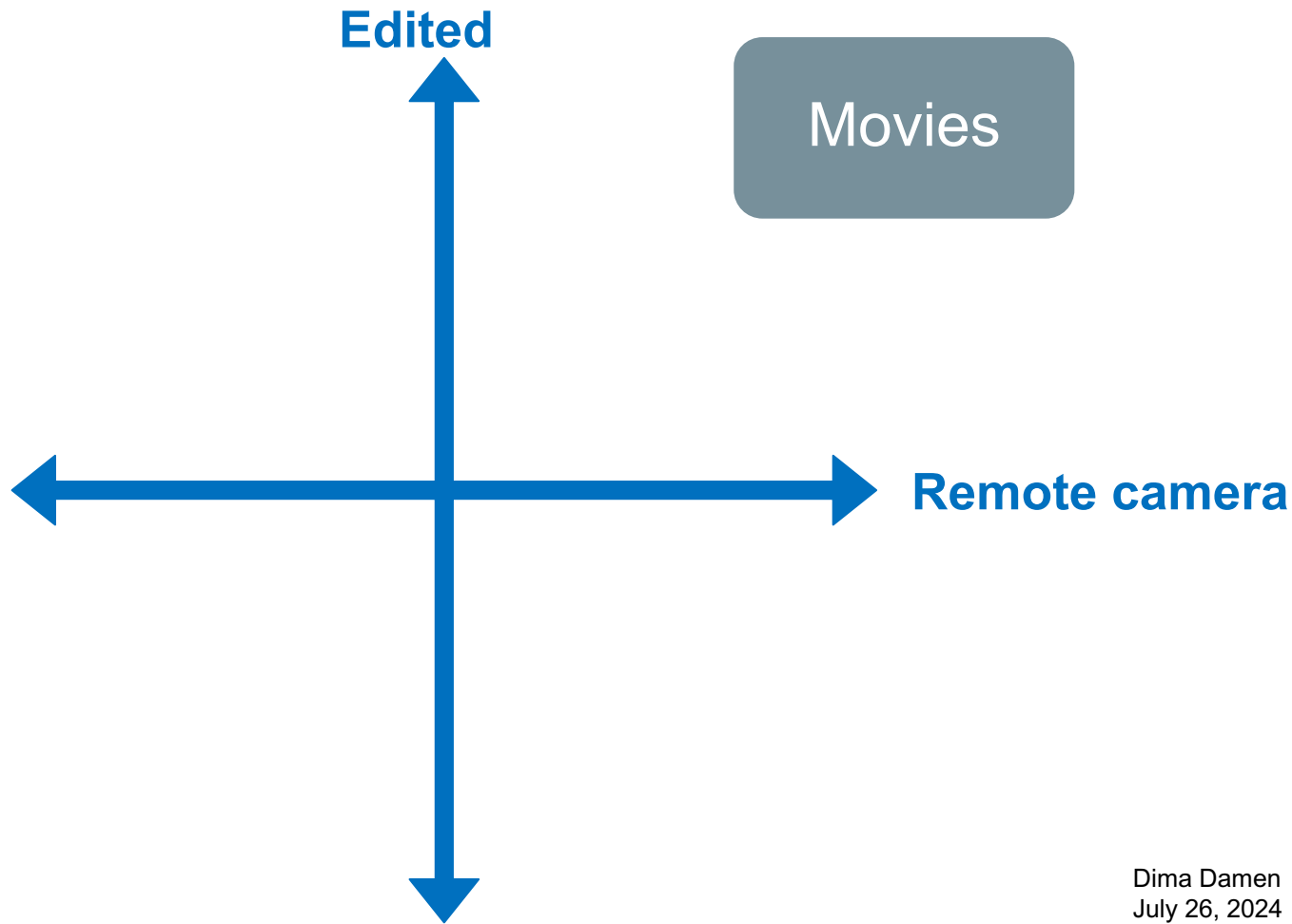


Part I: Collecting a Dataset



Our dataset is made up of... *videos*

The history of *VIDEO Understanding*



The history of VIDEO Understanding



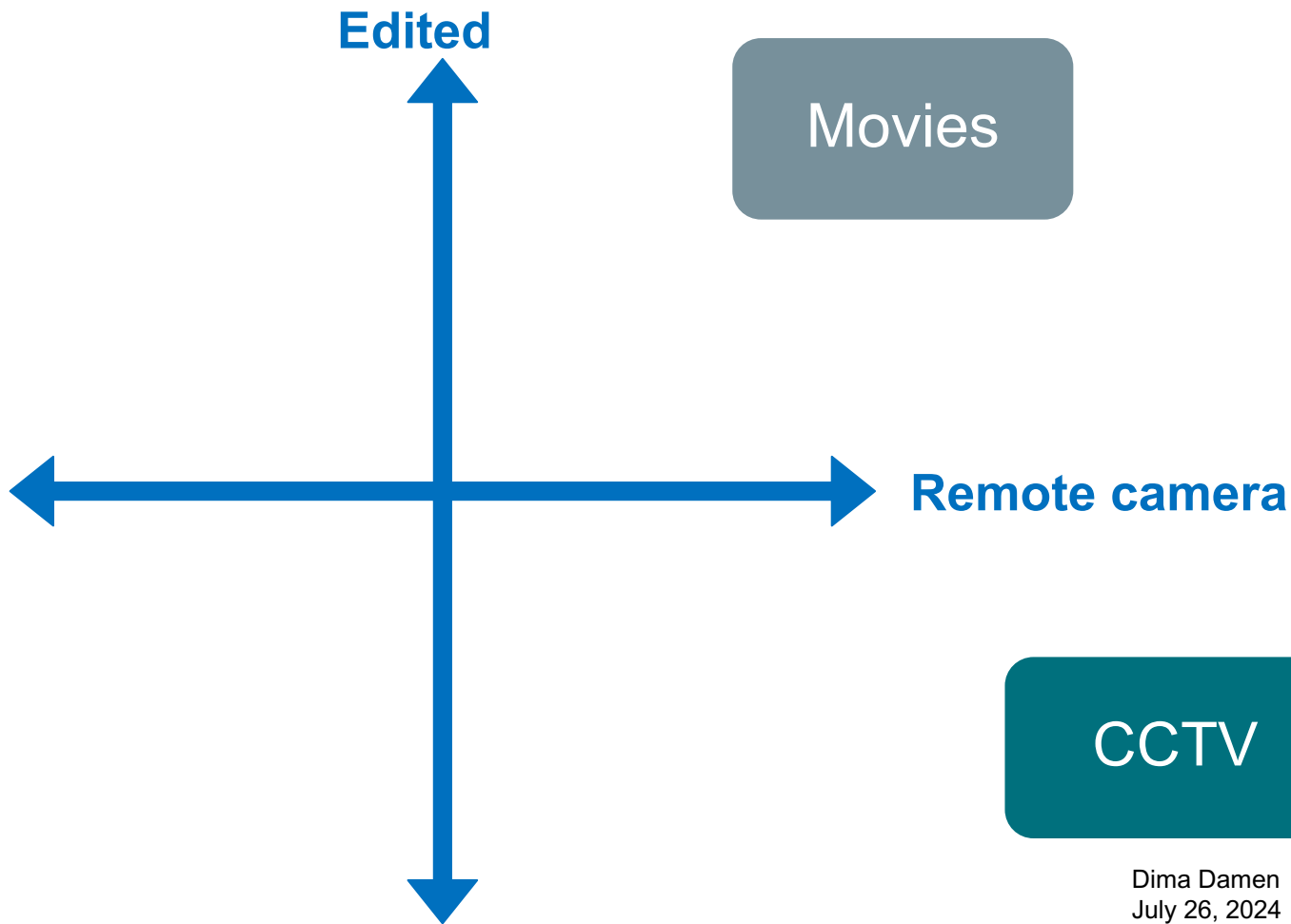
Figure 1. Examples of two action classes (drinking and smoking) from the movie “Coffee and Cigarettes”. Note the high within-

Laptev and Perez (2007)

The history of VIDEO Understanding



The history of *VIDEO Understanding*



The history of VIDEO Understanding



Damen and Hogg (2009). Recognizing linked events: Searching the space of feasible explanations. CVPR

The history of VIDEO understanding



How

**Templated,
Multilingual Domain
Queries:**

“Morning routine”,
“realistic ditl 2015”,
“mijn realistische
routine”, “Ma routine
d'apres-midi”, ...

216K Video Candidates (2.5 Years)
Low *Video-level* Purity



two stitches on two
and we'll slip stitch



by skipping the first
three stitches



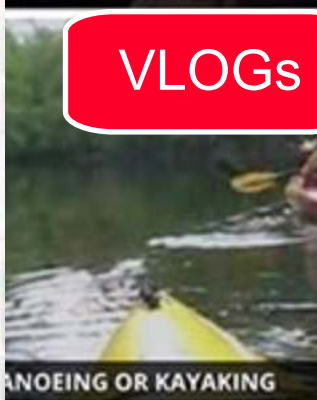
two stitches on two
and we'll slip stitch

**Egocentric
unscripted**

such and just going
to Mariel all the way

TAUCTIONING sta.com

VLOGs



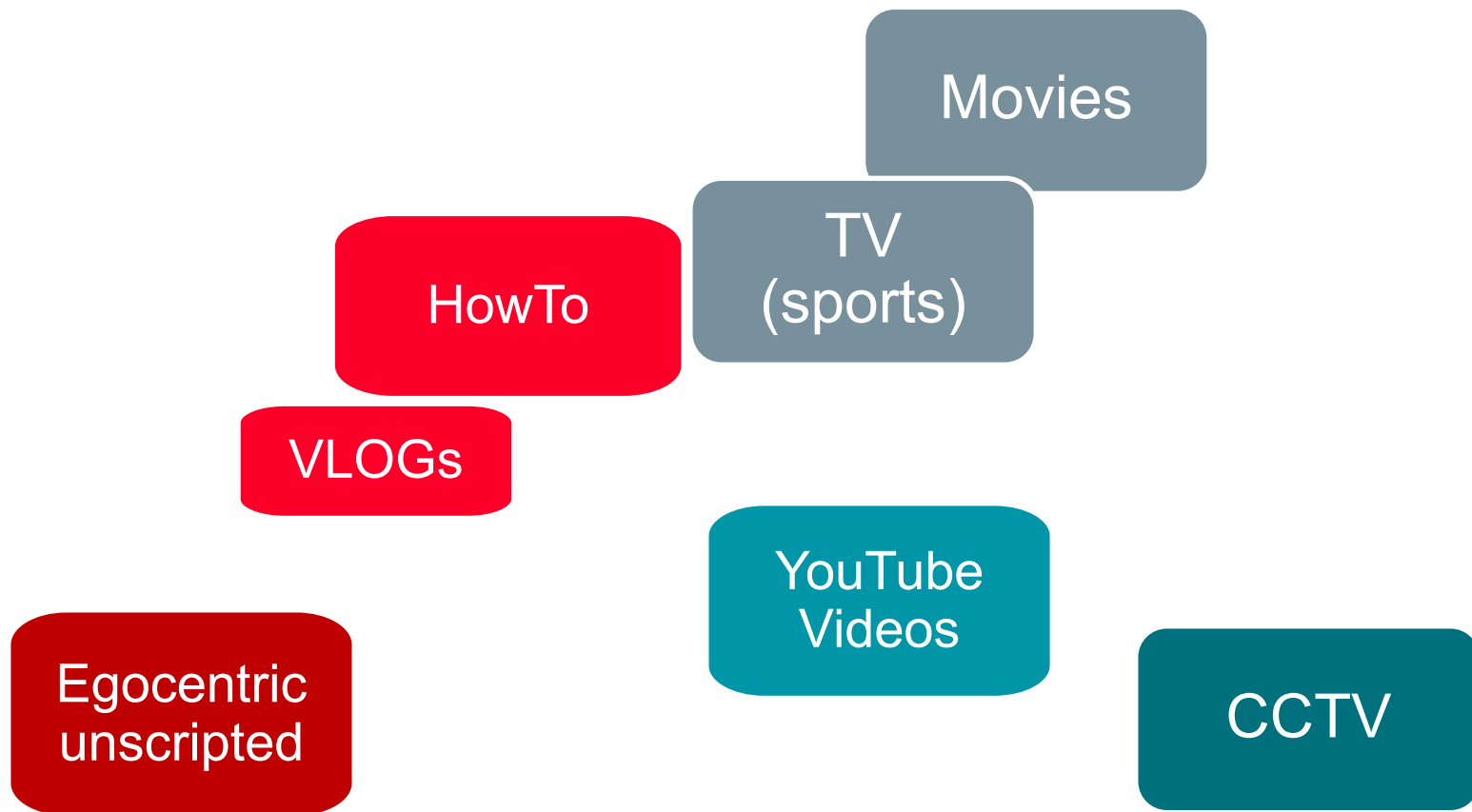
**YouTube
Videos**

Remote camera

CCTV

realted

The history of Video Understanding



Video Understanding

Speech/Plot

Movies

HowTo

VLOGs

Edits/Shots

Movies

HowTo

Audio-Visual

Movies

YouTube

Egocentric

Hand-Obj

HowTo

Egocentric

**Guidance/
Assistance**

HowTo

Egocentric

The Egocentric Perspective

with: Kristen Grauman
+83 authors



Egocentric Videos?



Data Collection Exercises



EPIC
KITCHENS

2017 - now

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



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

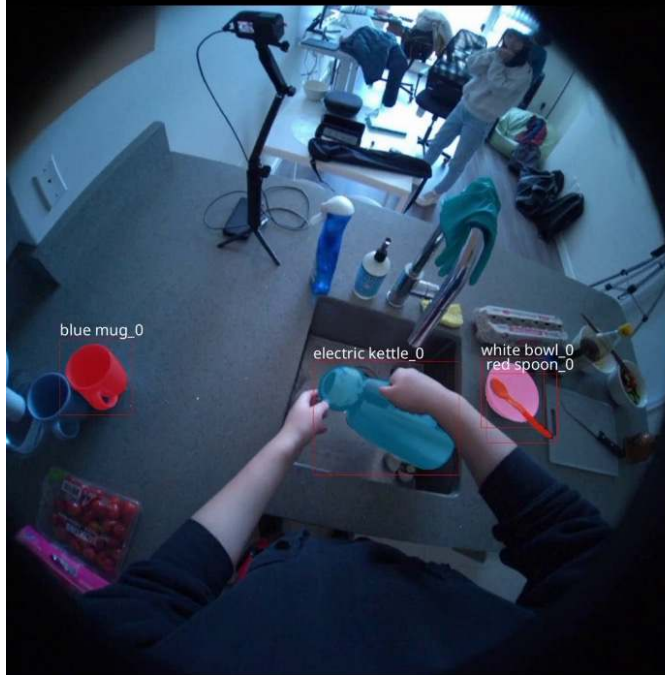


2024 – [coming]

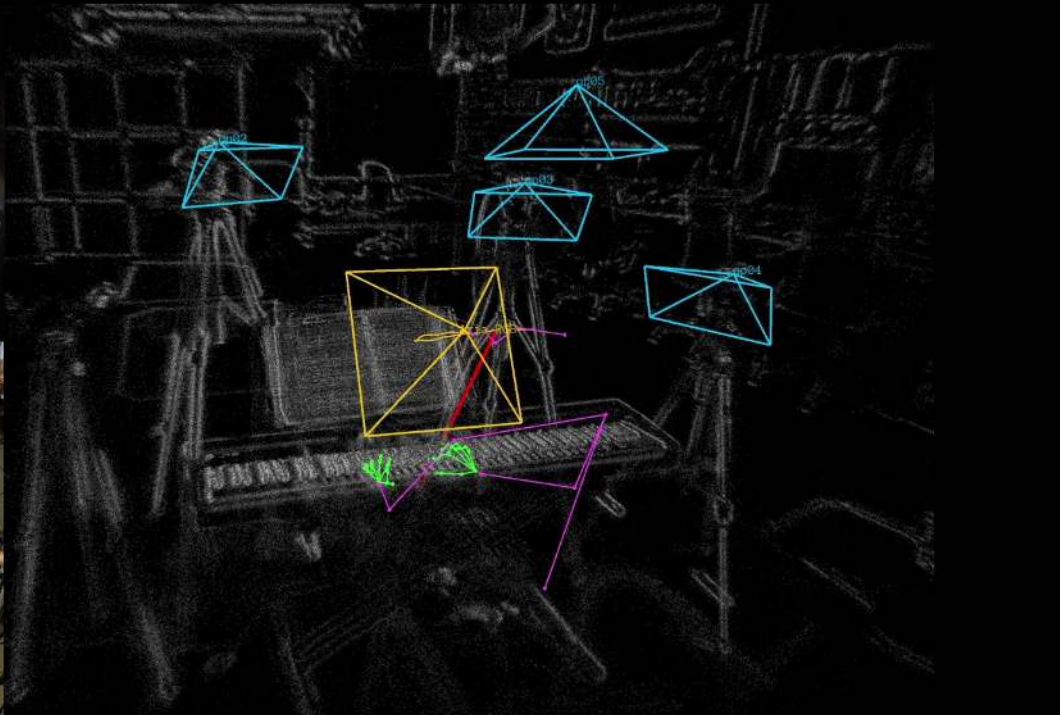
[new recordings]

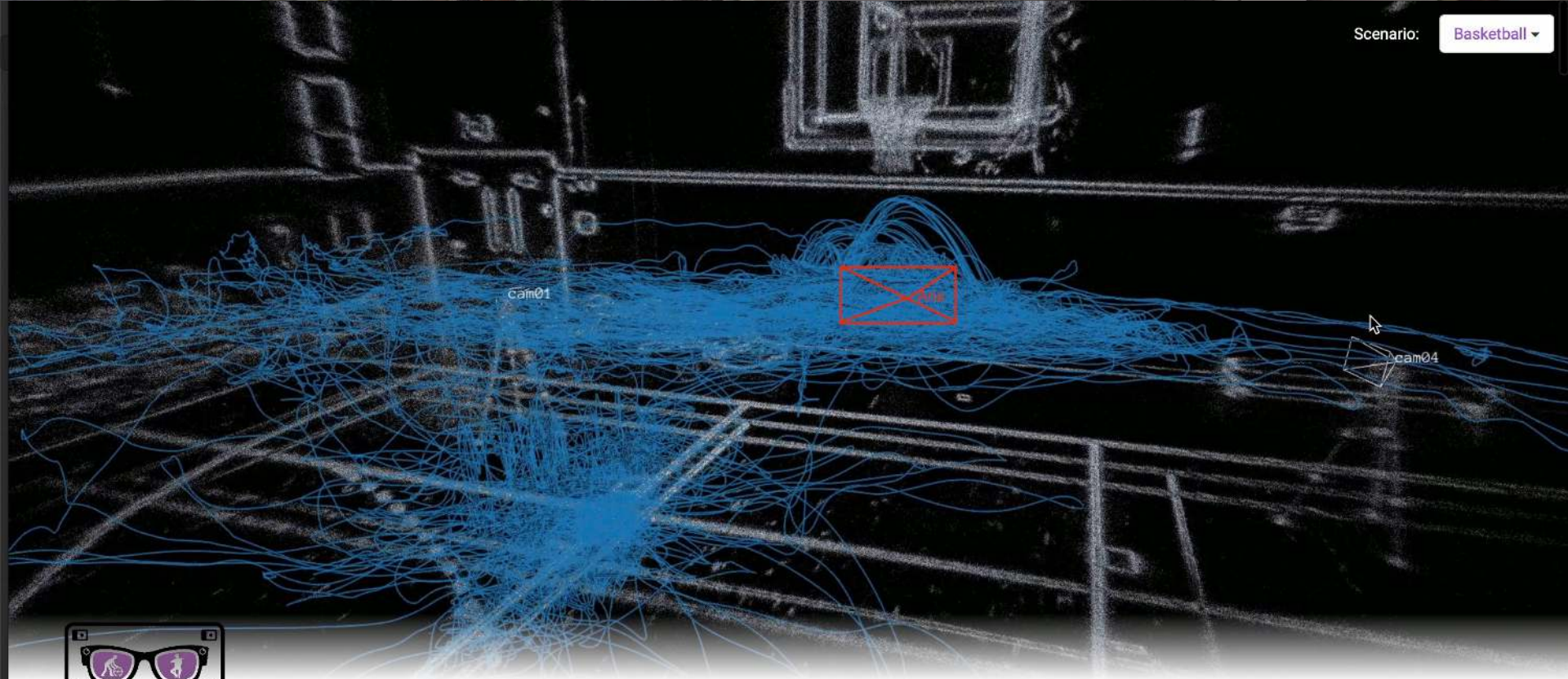


Ego-Exo Relation



Ego Pose





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.

Hover your mouse over scene cameras above to see a sample video for the chosen scenario.

Learn More ↓

Watch Video ↗

Start Here ↗

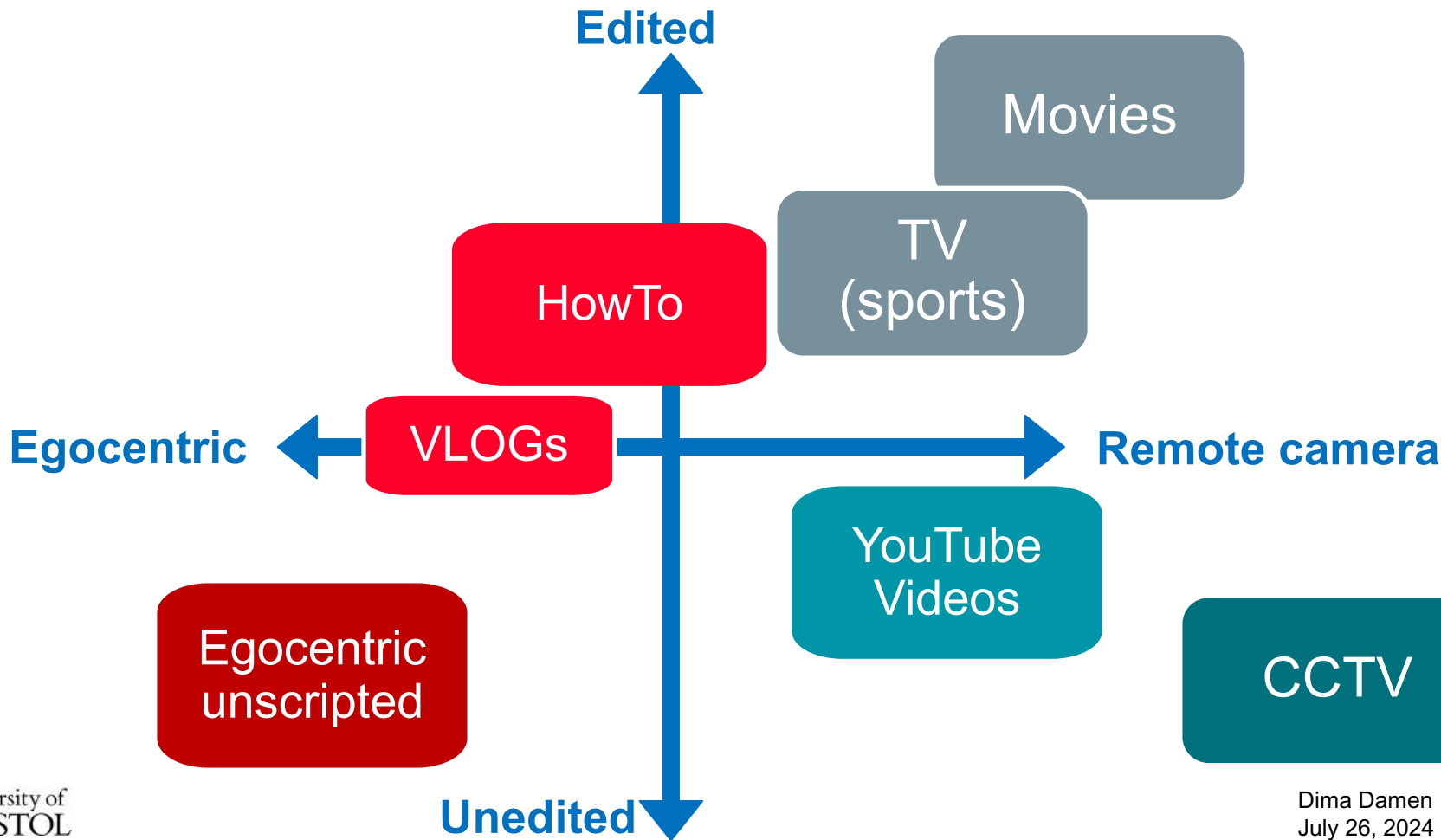


WHAT
DO YOU
THINK?



If you were to do research in **video understanding**, which video type(s) would you explore? Why?

The history of *VIDEO* understanding





sli.do

Joining as a participant?

#3639 120



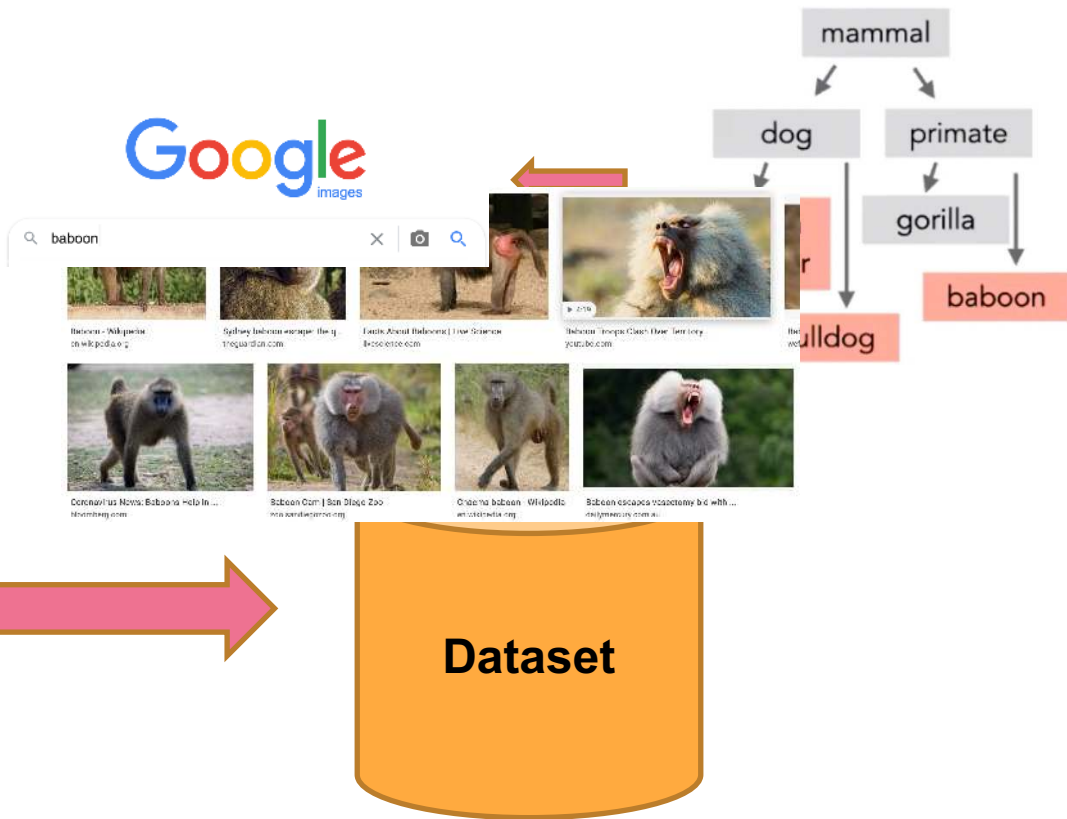


**How to collect a
dataset of videos?**



ImageNet Dataset

Object Recognition



Kinetics Dataset

A. List of Kinetics Human Action Classes

This is the list of classes included in the human action video dataset. The number of clips for each action class is given by the number in brackets following each class name.

1. abseiling (1146)
2. air drumming (1132)
3. answering questions (478)
4. applauding (411)
5. applying cream (478)
6. archery (1147)



YouTube

abseiling

Showing results for 'abseiling'. Search instead for 'abseiling'

- How to abseil! 140K views · 5 years ago
- How to set up an abseil - easy! 5.7K views · 1 year ago
- Abseiling 21K views · 6 years ago
- Abseiling down Northampton Ice tower 2.6K views · 1 year ago
- How to set up an Abseil | Climbing Daily Ep. 10-5 13K views · 2 years ago

D

A. List of Kinetics Human Action Classes

This is the list of classes included in the human action video dataset. The number of clips for each action class is given by the number in brackets following each class name.

1. abseiling (1146)
2. air drumming (1132)
3. answering questions (478)
4. applauding (411)
5. applying cream (478)
6. archery (1147)
7. arm wrestling (1123)
8. arranging flowers (583)
9. assembling computer (1147)
10. auctioning (478)
11. baby waking up (611)
12. baking cookies (927)
13. balloon blowing (826)
14. bandaging (569)
15. barbecuing (1070)

Statistics: The dataset has 400 human action classes, with 400–1150 clips for each action, each from a unique video. Each clip lasts around 10s. The current version has 306,245 videos, and is divided into three splits, one for training having 250–1000 videos per class, one for validation with 50 videos per class and one for testing with 100 videos per class. The statistics are given in table 2.

One Exception



Machine Learning in Practice

- Autonomous Driving...

Welcome to the KITTI Vision Benchmark Suite!

We take advantage of our [autonomous driving platform Annieway](#) to develop novel challenging real-world computer vision benchmarks. Our tasks of interest are: stereo, optical flow, visual odometry, 3D object detection and 3D tracking. For this purpose, we equipped a standard station wagon with two high-resolution color and grayscale video cameras. Accurate ground truth is provided by a Velodyne laser scanner and a GPS localization system. Our datasets are captured by driving around the mid-size city of [Karlsruhe](#), in rural areas and on highways. Up to 15 cars and 30 pedestrians are visible per image. Besides providing all data in raw format, we extract benchmarks for each task. For each of our benchmarks, we also provide an evaluation metric and this evaluation website. Preliminary experiments show that methods ranking high on established benchmarks such as [Middlebury](#) perform below average when being moved outside the laboratory to the real world. Our goal is to reduce this bias and complement existing benchmarks by providing real-world benchmarks with novel difficulties to the community.



To get started, grab a cup of your favorite beverage and watch our video trailer (5 minutes):

stereo flow sceneflow depth odometry object tracking road semantics raw data

Machine Learning in Practice

Object Recognition

Let's collect Data!



EPIC
KITCHENS-100



Scaling and Rescaling Egocentric Vision: The **EPIC-KITCHENS** Dataset



Dima Damen



Hazel Doughty



Giovanni M. Farinella



Sanja Fidler



Antonino Furnari



Evangelos Kazakos



Jian Ma



Davide Moltisanti



Jonathan Munro



Toby Perrett



Will Price



Michael Wray

EPIC-KITCHENS



Scaling and Rescaling Egocentric Vision

- Head-Mounted Go-Pro, adjustable mounting
- Recording starts immediately before entering the kitchen
- Only stopped before leaving the kitchen



EPIC-KITCHENS





open oven



put spoon on counter



put on glove



pick up fork



put down glass
pick up glass



put down plate

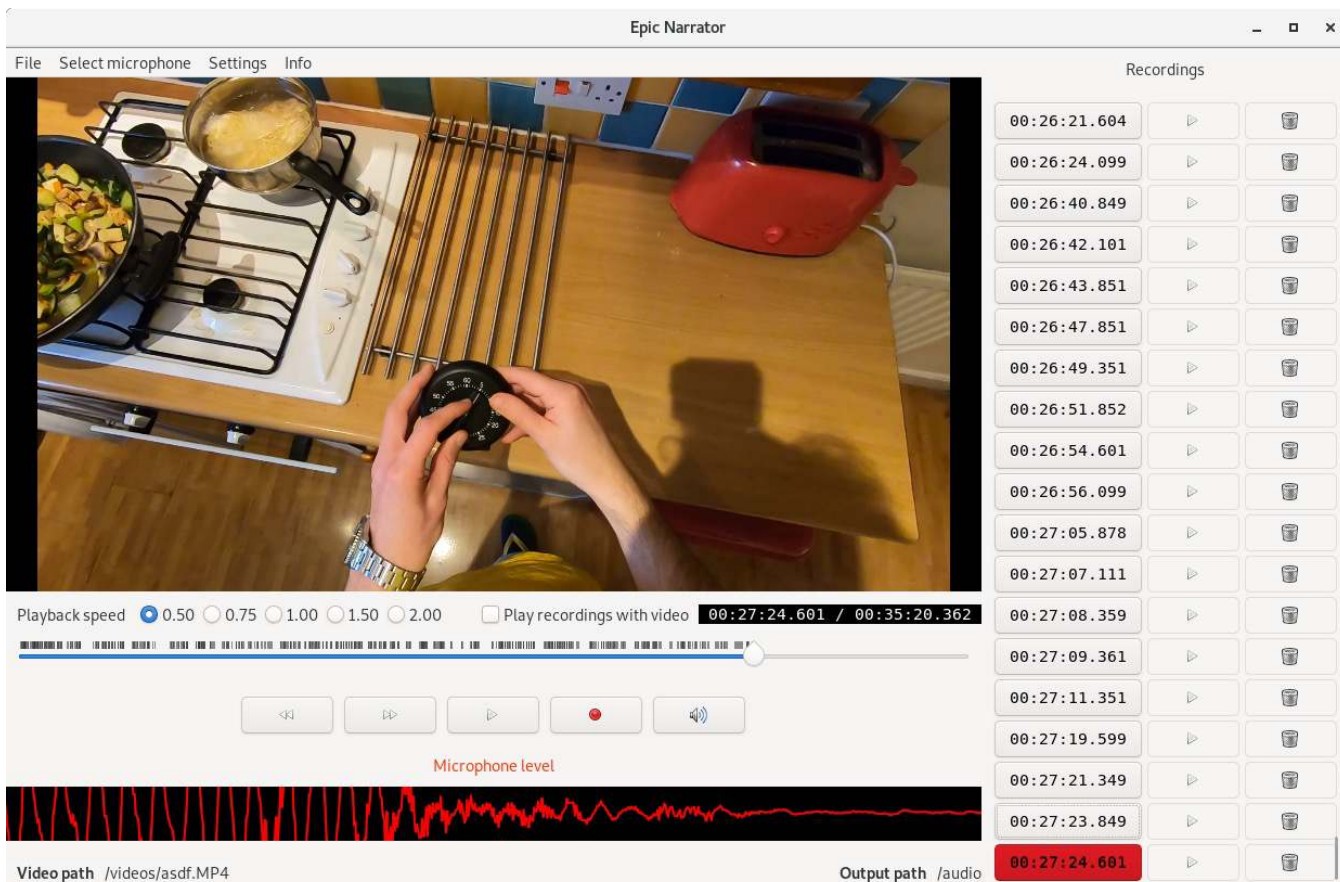


put down spoon
pick up aeropress filter

EPIC-KITCHENS

Epic Narrator

File Select microphone Settings Info



Recordings

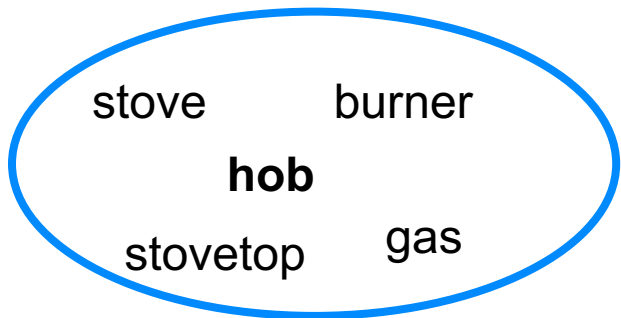
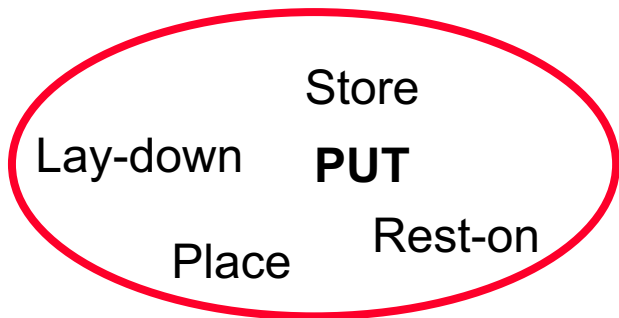
00:26:21.604	▶	🗑️
00:26:24.099	▶	🗑️
00:26:40.849	▶	🗑️
00:26:42.101	▶	🗑️
00:26:43.851	▶	🗑️
00:26:47.851	▶	🗑️
00:26:49.351	▶	🗑️
00:26:51.852	▶	🗑️
00:26:54.601	▶	🗑️
00:26:56.099	▶	🗑️
00:27:05.878	▶	🗑️
00:27:07.111	▶	🗑️
00:27:08.359	▶	🗑️
00:27:09.361	▶	🗑️
00:27:11.351	▶	🗑️
00:27:19.599	▶	🗑️
00:27:21.349	▶	🗑️
00:27:23.849	▶	🗑️
00:27:24.601	▶	🗑️

Playback speed 0.50 0.75 1.00 1.50 2.00 Play recordings with video 00:27:24.601 / 00:35:20.362

Microphone level

Video path /videos/asdf.MP4 Output path /audio

EPIC-KITCHENS and Ego4D



open vocab

lay-down

stovetop



closed vocab

put

hob

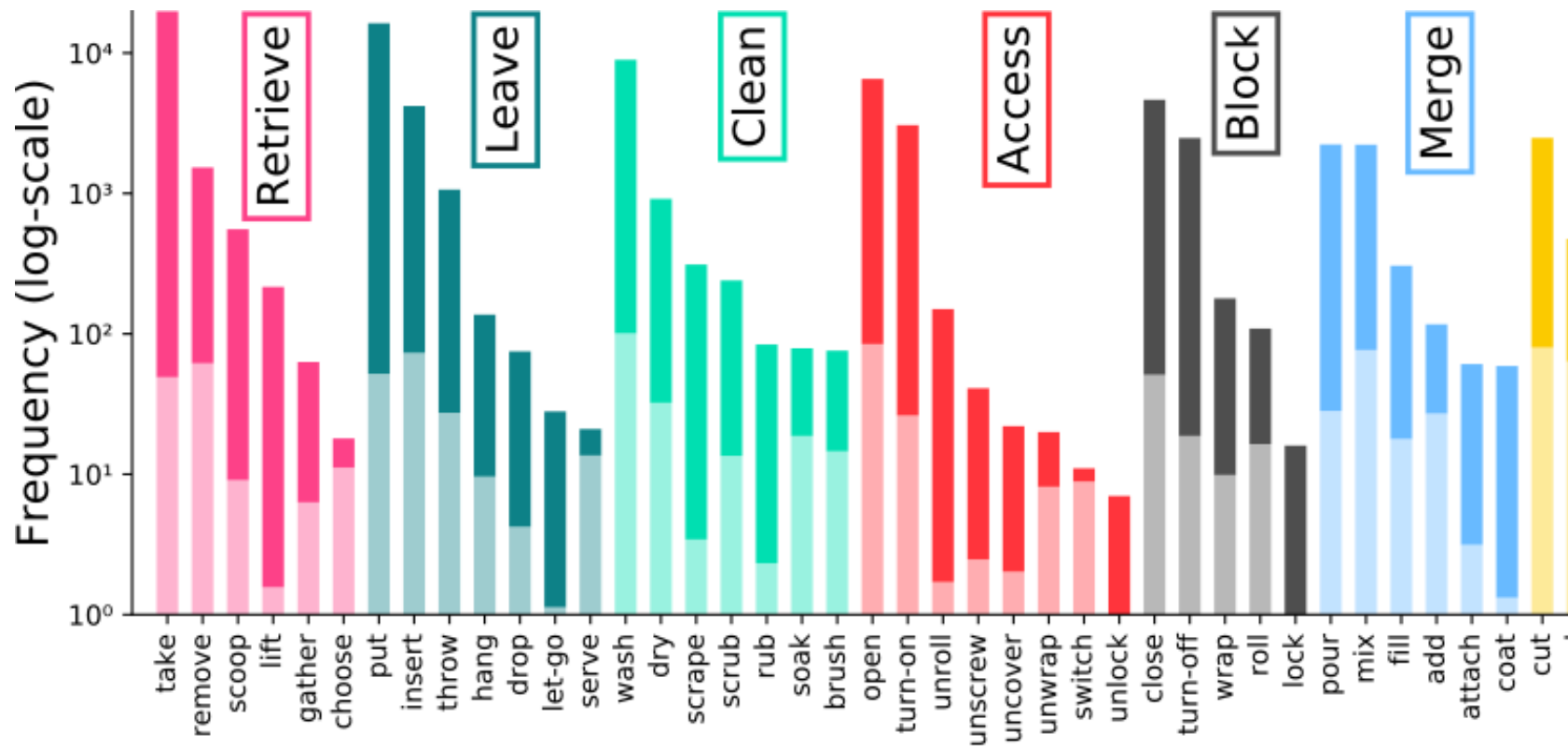


category

leave

appliance

EPIC-KITCHENS-100 Statistics



Data Collection Exercise



Labels

Pascal VOC
ImageNet
Kinetics
Something-Something



Data

EPIC-KITCHENS
Ego4D
...
KITTI

The chicken or the egg...

Data



Naturally unbalanced

Harder to label (exposes ambiguity)

Closer to application

Multiple tasks

Labels



Unnaturally balanced (or nearly)

Easier to label (hides ambiguity)

Can be expanded

Single task



WHAT
DO YOU
THINK?

What should come first? Labels or Data



Labels



Data





sli.do

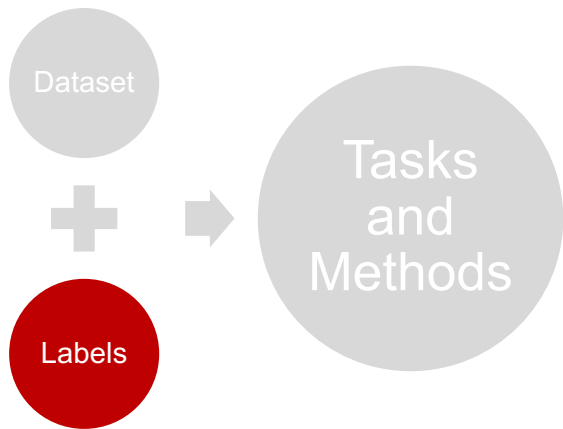
Joining as a participant?

#3639 120





Video source and data
collection approach heavily
influences video
understanding tasks



Part II: Labelling a Dataset

What type of labels can we provide?

- Temporal labels – Strong vs. Weak labels
- Semantic labels – Open-vocab. vs Closed-vocabulary
- Ranking labels – video-to-video comparisons
- Pixel-level labels – segmentation labels

What type of labels can we provide?

- Temporal labels – Strong vs. Weak labels
- Semantic labels – Open-vocab. vs Closed-vocabulary
- Ranking labels – video-to-video comparisons
- Pixel-level labels – segmentation labels

Action Recognition Challenge



Given a trimmed action segment:

$(t_{\text{start}}, t_{\text{stop}})$

classify the action within.

$\hat{y}_{\text{verb}} = \text{open}$

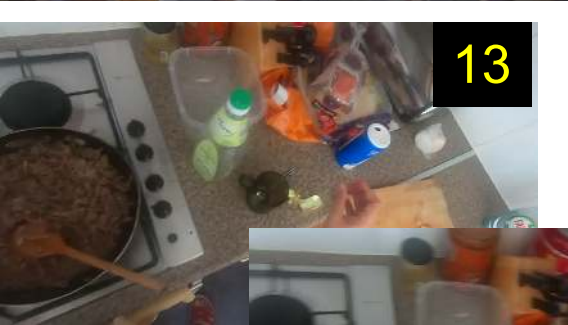
$\hat{y}_{\text{noun}} = \text{oven}$

$\hat{y}_{\text{action}} = (\text{open}, \text{oven})$



Do
It
Yourself





13

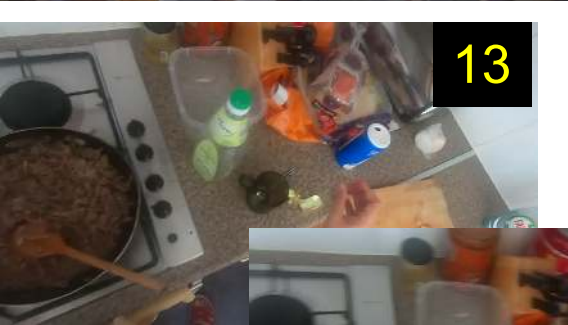
56

73

93

131

177



13

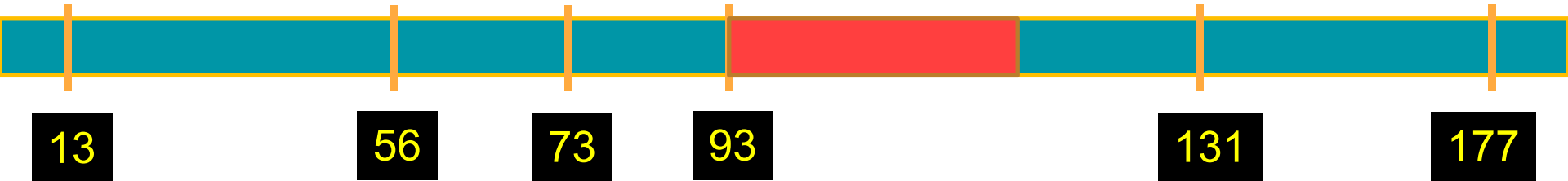
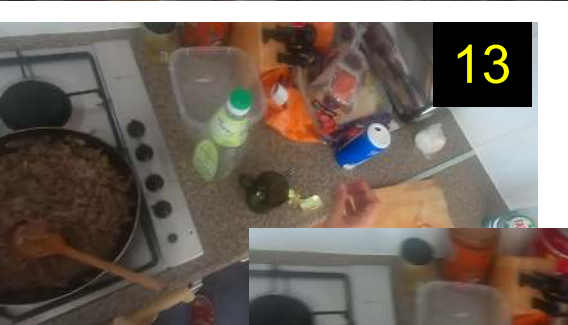
56

73

93

131

177



13

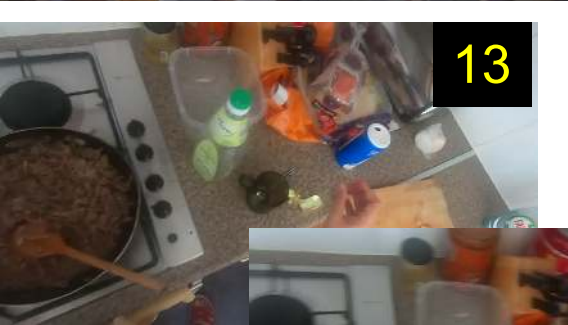
56

73

93

131

177



13

56

73

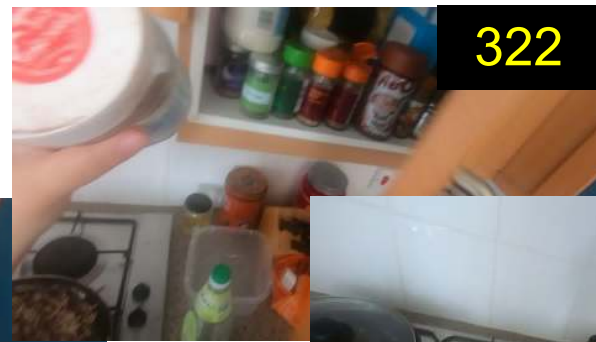
93

131

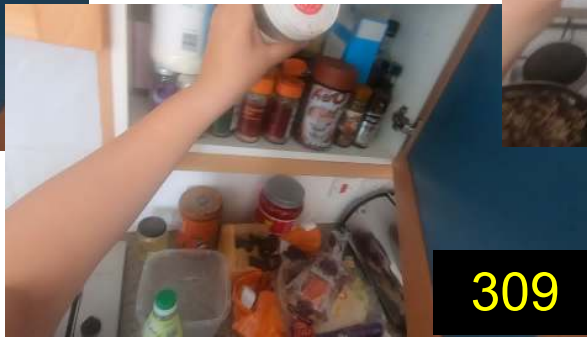
177



280



322



309



348



280

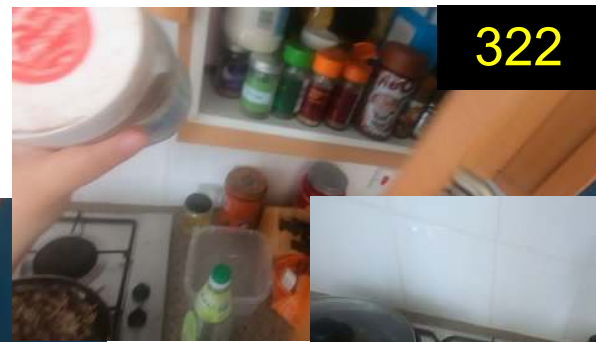
309

322

348



280



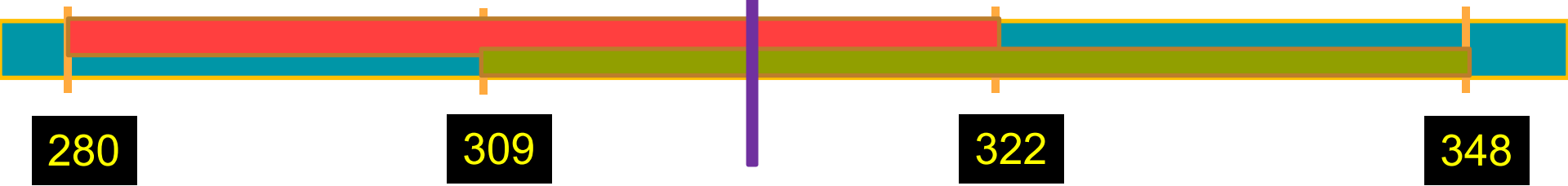
322



309



348



280

309

322

348

Inconsistencies of temporal bounds across datasets for the same action

BEOID: take cup





GTEA Gaze+

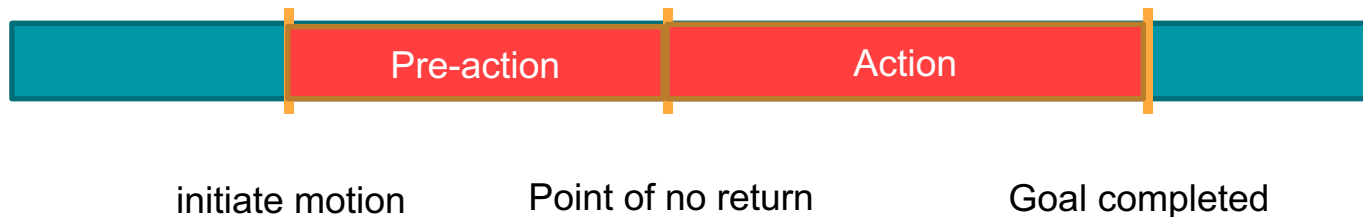
ground truth



predicted class: take knife



- [A] There are two stages of an action, separated by three boundary points
 - Pre-action stage:
 - Action stage:



[A] P. M. Gollwitzer (1990). Action phases and mind-sets. Handbook of motivation and cognition.

Cut pepper (GTEA Gaze+)

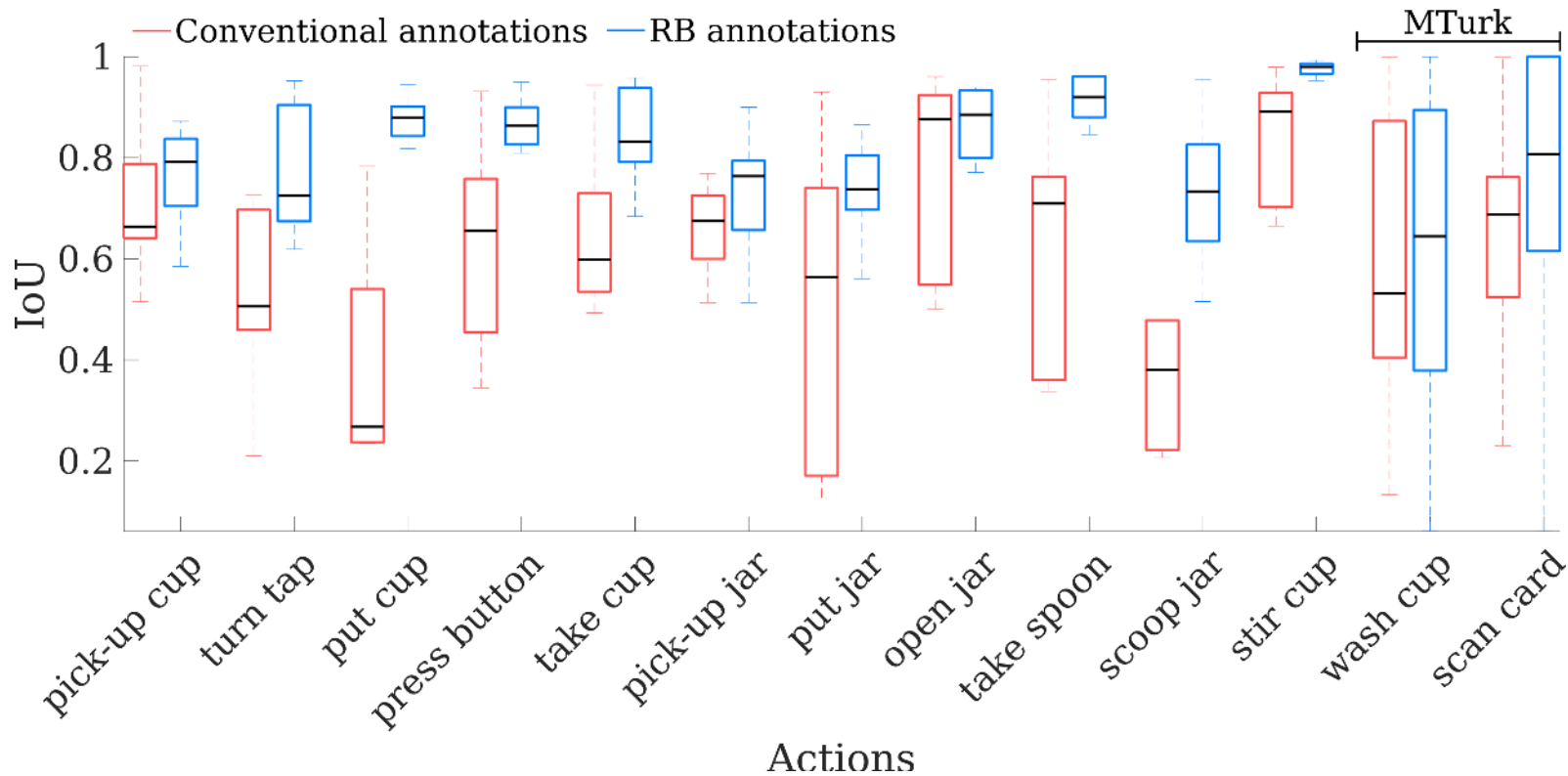


Rubicon Boundaries

Now we show some object interactions segmented by multiple annotators using conventional labeling, along with the same actions labeled by different annotators following the Rubicon Boundaries (ref. Figure 3).

The Rubicon Boundaries

with: Davide Moltisanti



The power of temporal labels

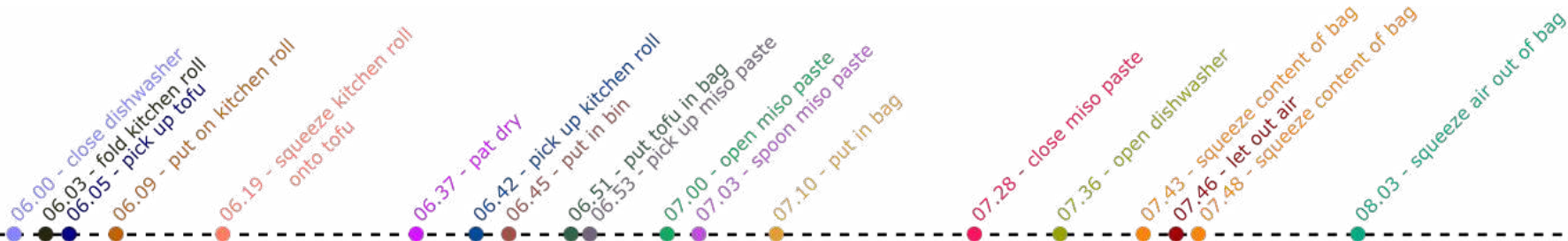


Other approaches to temporal boundaries

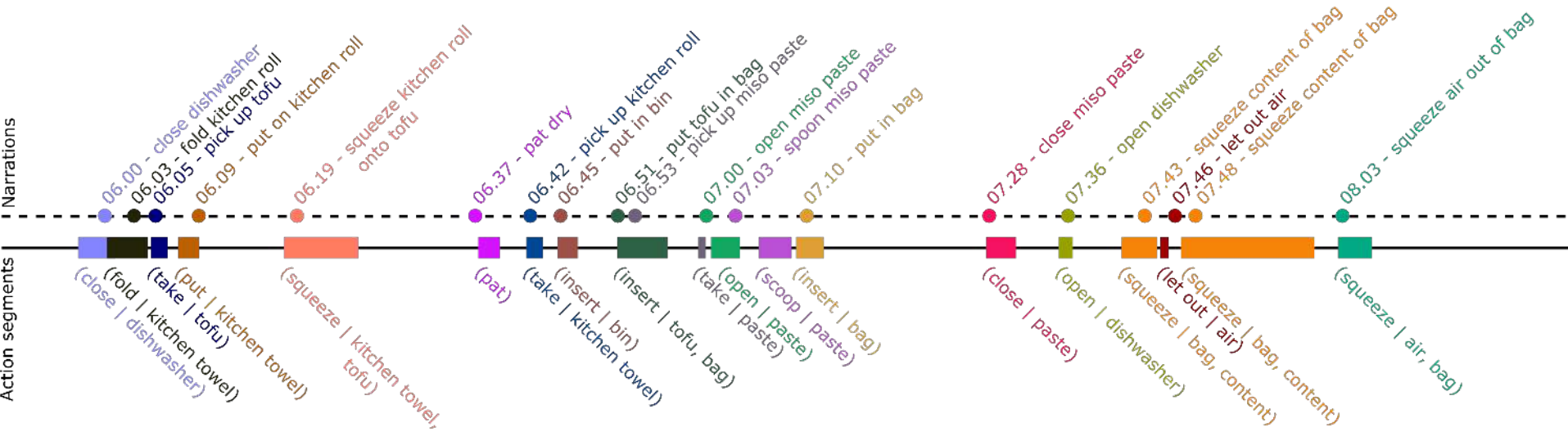
- Start-and-end times
 - Inconsistent
 - Consistent
- Fixed segment lengths
 - Kinetics Dataset -- 10 seconds videos
 - Moments in Time Dataset – 3 seconds videos
- No temporal annotations
 - Charades Dataset – Video-Level supervision (3-4 actions per video)
- Single-timestamp supervision

Scaling and Rescaling Egocentric Vision

Narrations



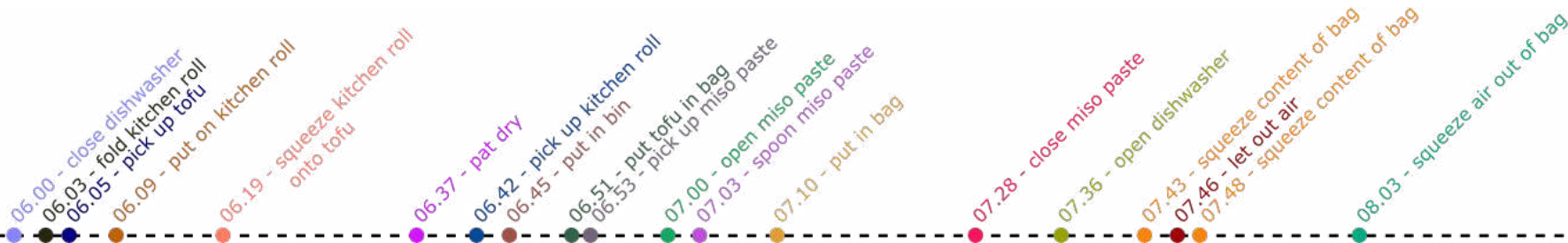
Scaling and Rescaling Egocentric Vision



Learning from a Single Timestamp

with: Davide Moltisanti
Sanja Fidler

Narrations



pick up
cup

turn
tap

rinse
cup

turn
tap

put
cup

press
button

take
cup

put
cup

pick-up
jar

put
jar

take
spoon

open
jar

scoop
spoon

pour
spoon

stir
spoon



video frames

Learning from a Single Timestamp

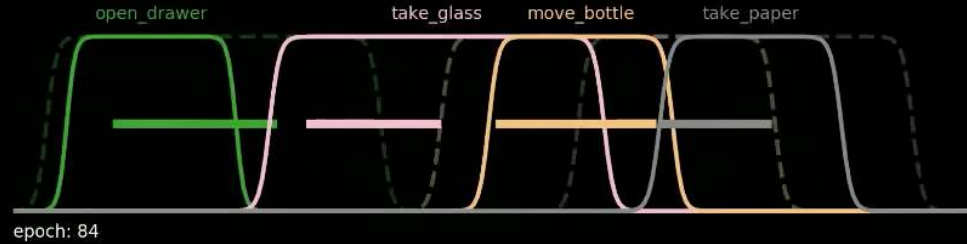
with: Davide Moltisanti
Sanja Fidler



Learning from a Single Timestamp

with: Davide Moltisanti
Sanja Fidler

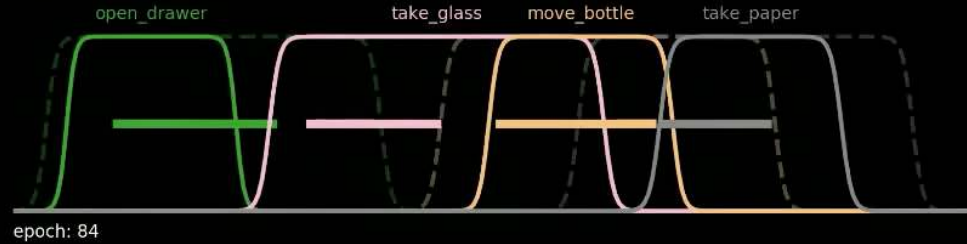
i) EPIC Kitchens (success)



Learning from a Single Timestamp

with: Davide Moltisanti
Sanja Fidler

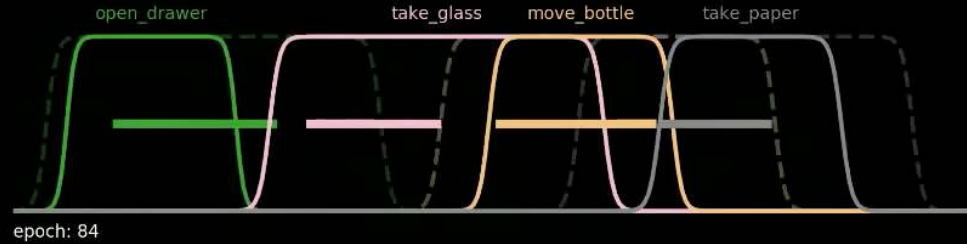
i) EPIC Kitchens (success)



Learning from a Single Timestamp

with: Davide Moltisanti
Sanja Fidler

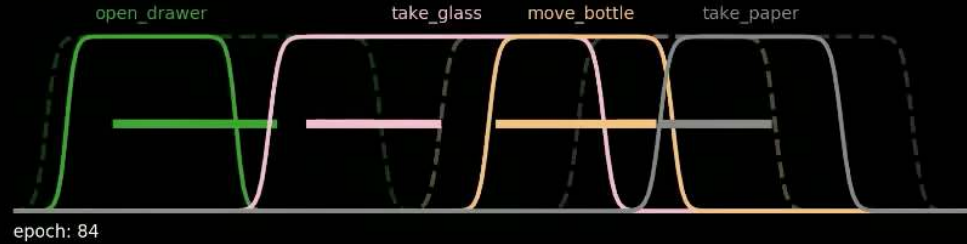
i) EPIC Kitchens (success)



Learning from a Single Timestamp

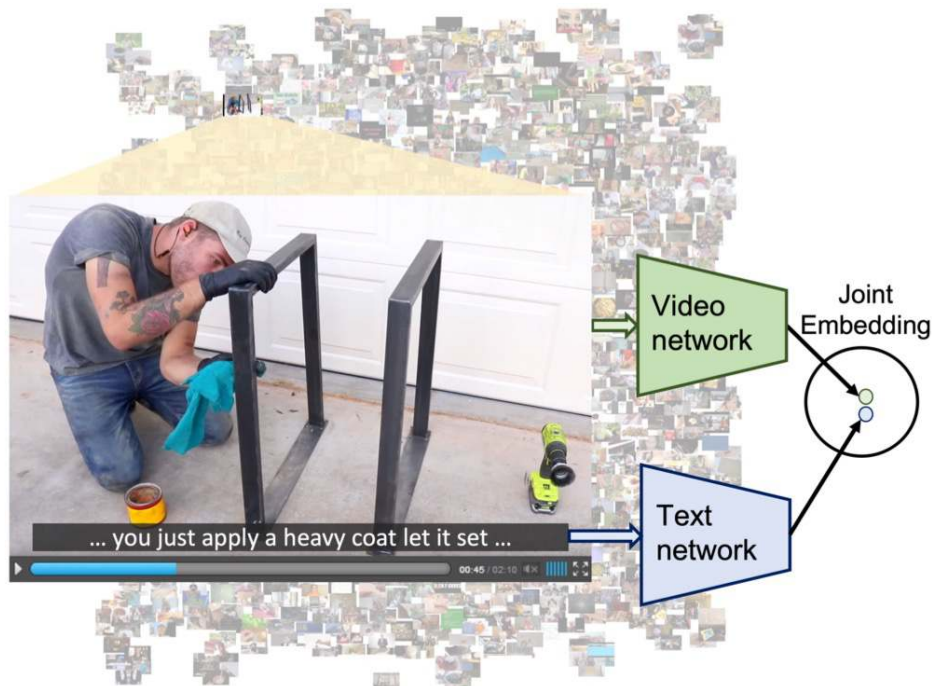
with: Davide Moltisanti
Sanja Fidler

i) EPIC Kitchens (success)



Learning from a Narration Timestamps

- Miech et al (2019). HowTo100M: Learning a Text-Video Embedding by Watching Hundred Million Narrated Video Clips



#C C picks trowel



#C C opens the decoration balls box



#C C walks out



#C C picks water bottle



#C C rolls the rope



#C C places the pen on the canvas paper



#C C folds pizza box



#C C wipes his hands with his trousers



#C C moves the soil.



#C C throws the coconut.



#C C tightens the knob of a lawn mower



#C C chops the cucumber



#C C fixes the pipe in the cable pass



#C C hold the piece of cloth



#C C moves the shovel.



#C C spreads the fabric



#C C climbs the ladder



#C C throws a ball



#O person E uses phone



#O A man X looks at the ceiling



#O man x talks to c



#X person o shows the tin to the child



#O Lady Y moves meat in a bowl



#O person A drops the chaff in the dust bin

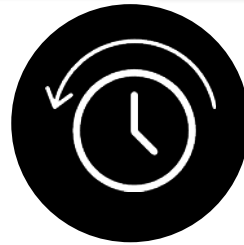




Temporal labels vary across
datasets... and should be
consistent

Reversing Time

with: Will Price



W Price, D Damen (2019). Retro-Actions: Learning 'Close' by Time-Reversing 'Open' Videos. ICCV MDALC Workshop

Reversing Time



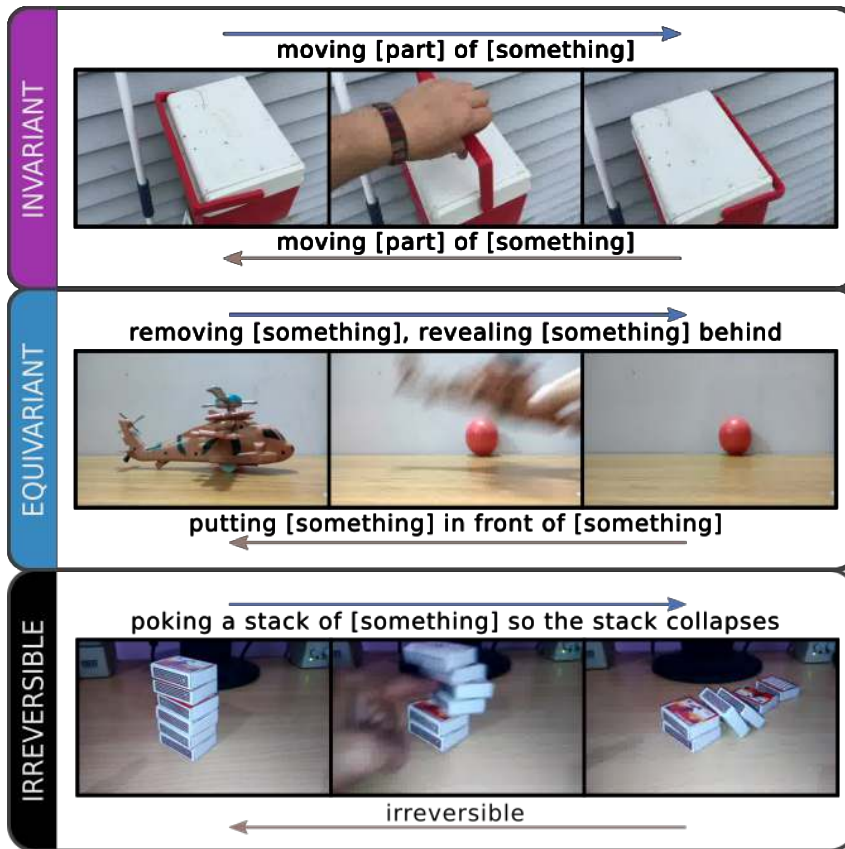
This CVPR2014 paper is the Open Access version, provided by the Computer Vision Foundation.
The authoritative version of this paper is available in IEEE Xplore.

Seeing the Arrow of Time

Lyndsey C. Pickup¹ Zheng Pan² Donglai Wei³ YiChang Shih³ Changshui Zhang²
Andrew Zisserman¹ Bernhard Schölkopf⁴ William T. Freeman³

Reversing Time

with: Will Price



Reversing Time

with: Will Price



Open

Close



Pull →

Push ←



Camera ←

Camera →



Moving _ and _ so they pass each other



Trying to bend _ unbendable



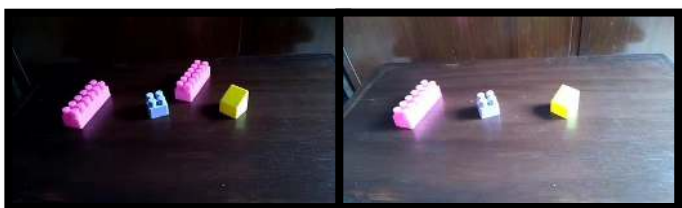
Removing _
revealing _
behind

Putting _ in
front of _



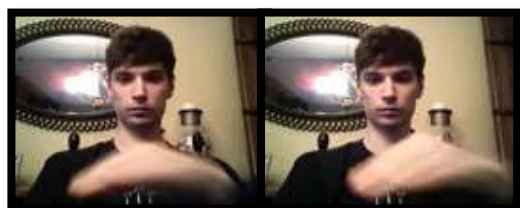
Swipe

Swipe



Take one of many
similar things on the
table

Put something similar
to other things already
on the table



Roll hand
forward

Roll hand
backward
Dina Daren
July 26, 2024



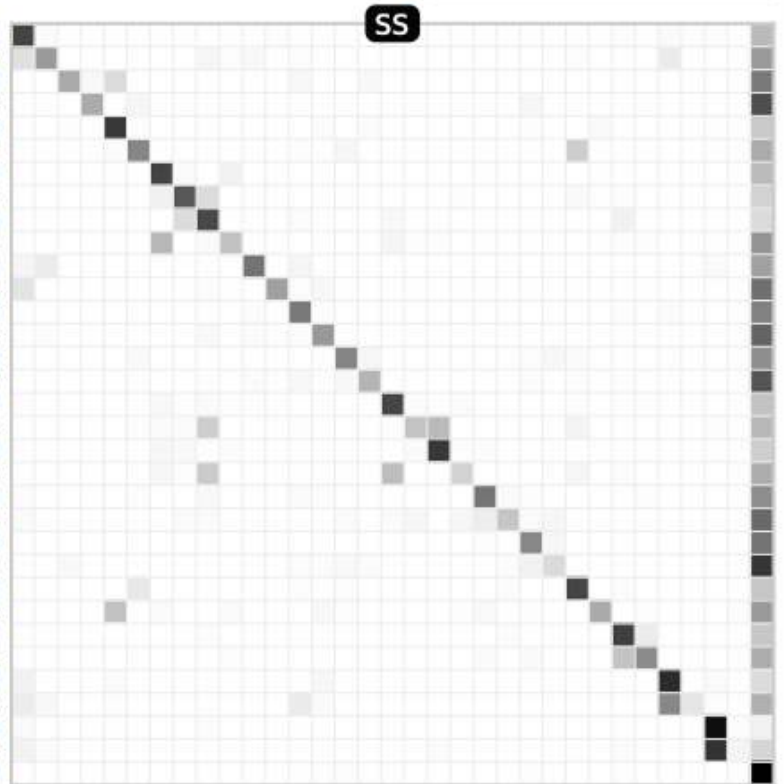
Can we train
Learning 'Close' by time-reversing 'Open'

Reversing Time

with: Will Price

TR

- Approaching something with your camera ●
 - Moving away from something with your camera ●
 - Burying something in something ●
 - Digging something out of something ●
 - Covering something with something ●
 - Uncovering something ●
 - Moving something and something closer to each other ●
 - Moving something and something away from each other ●
 - Moving something away from something ●
 - Moving something closer to something ●
 - Moving something away from the camera ●
 - Moving something towards the camera ●
 - Moving something up ●
 - Moving something down ●
 - Opening something ●
 - Closing something ●
 - Pushing something from left to right ●
 - Pulling something from right to left ●
 - Pushing something from right to left ●
 - Pulling something from left to right ●
 - Putting something behind something ●
 - Pulling something from behind of something ●
 - Putting something into something ●
 - Pulling something out of something ●
 - Removing something, revealing something behind ●
 - Putting something in front of something ●
 - Taking one of many similar things on the table ●
 - Putting something similar to other things that are already on the table ●
 - Turning the camera downwards while filming something ●
 - Turning the camera upwards while filming something ●
 - Turning the camera left while filming something ●
 - Turning the camera right while filming something ●
 - Other ●
- Many-shot
● Zero-shot



Now, results from a model supervised by
time-reversal example synthesis



The Arrow of Time is Critical



Multi-modal learning...

with: Vangelis Kazakos
Arsha Nagrani.
Andrew Zisserman

Jaesung Huh
Jacob Chalk

- The magic of audio-visual understanding...
- Object-Object interactions



Multi-modal learning...

with: Vangelis Kazakos
Arsha Nagrani.
Andrew Zisserman

Jaesung Huh
Jacob Chalk

- The magic of audio-visual understanding...
- Object-Object interactions
- Material sounds



Multi-modal learning...

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

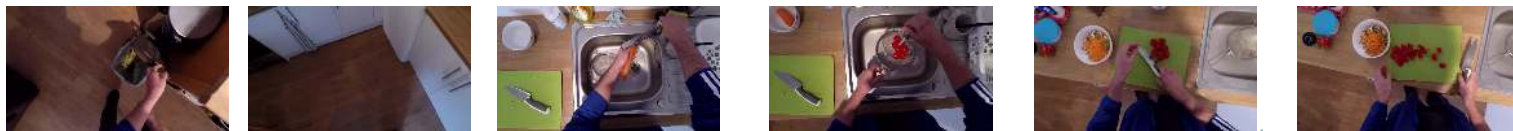
- The magic of audio-visual understanding...
- Object-Object interactions
- Material sounds
- Sound-emitting objects



Motivation

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

Video



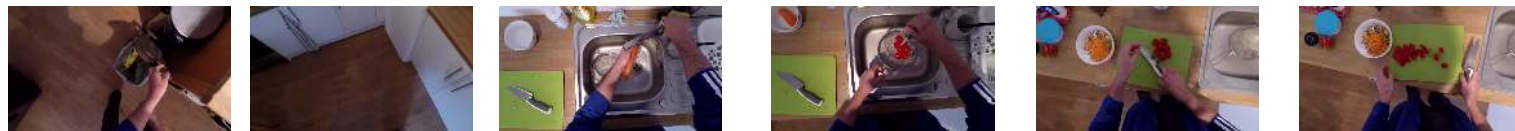
Audio



Motivation

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

Video



Close bin Close bag

Wash carrot

Wash tomato

Take knife

Cut tomato

Audio



Motivation

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

Video



Close bin Close bag

Wash

wash tomato

Take knife

Cut tomato

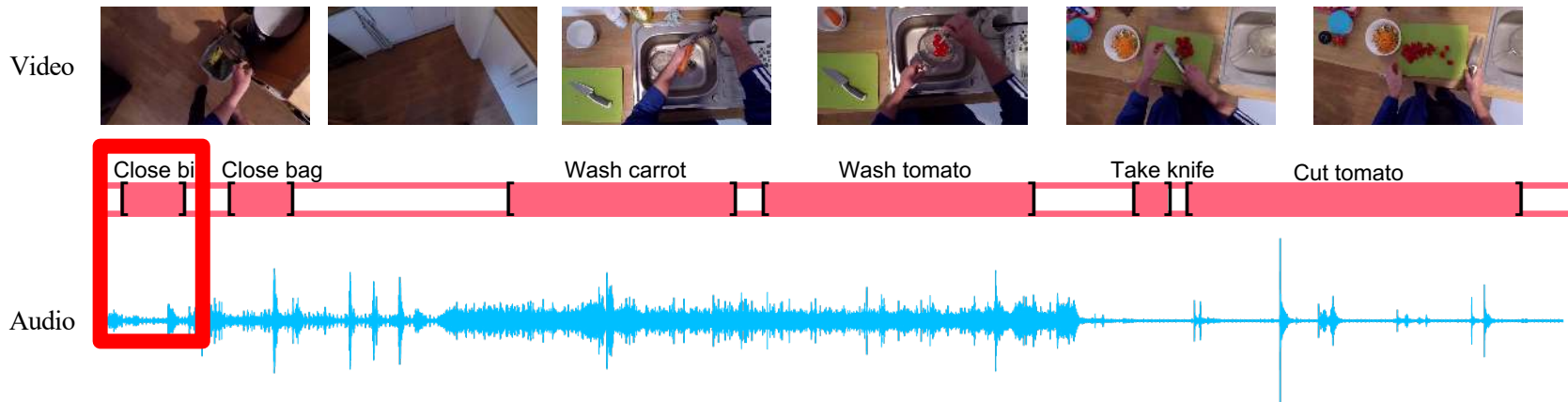
Incorrect assumption

Audio



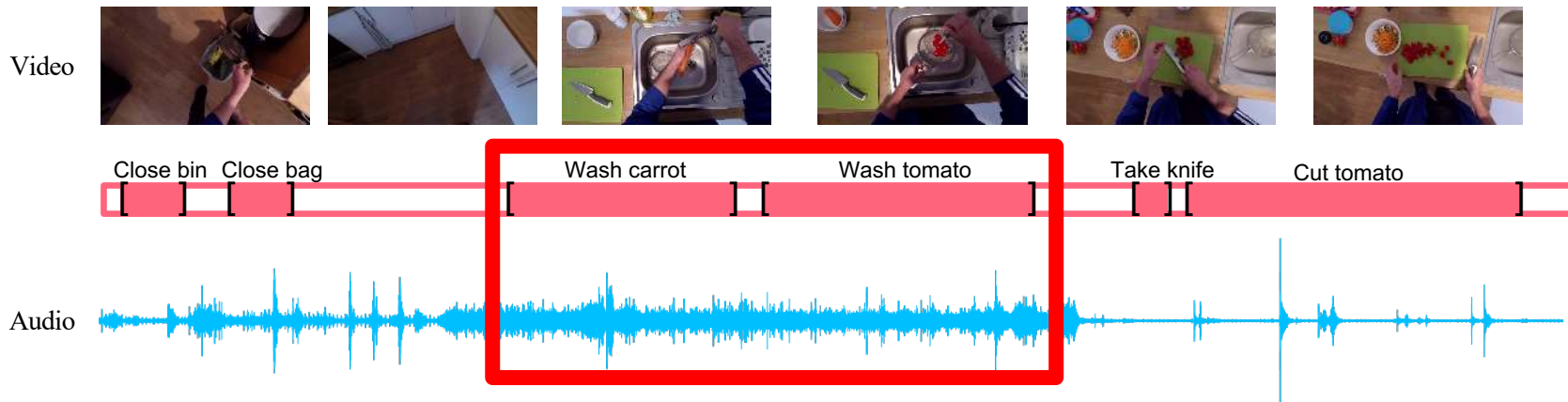
Motivation

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



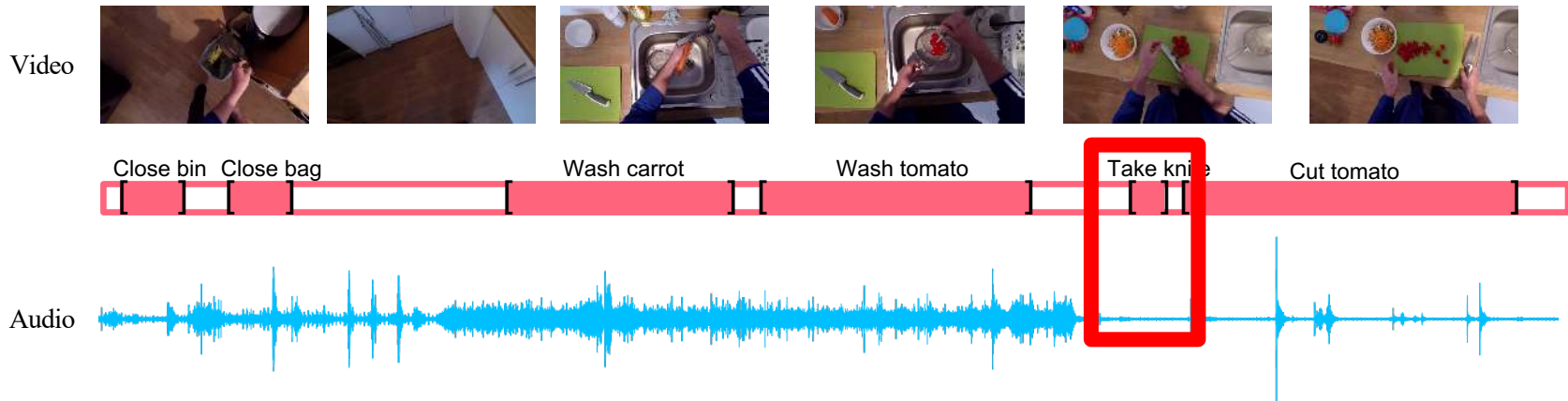
Motivation

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



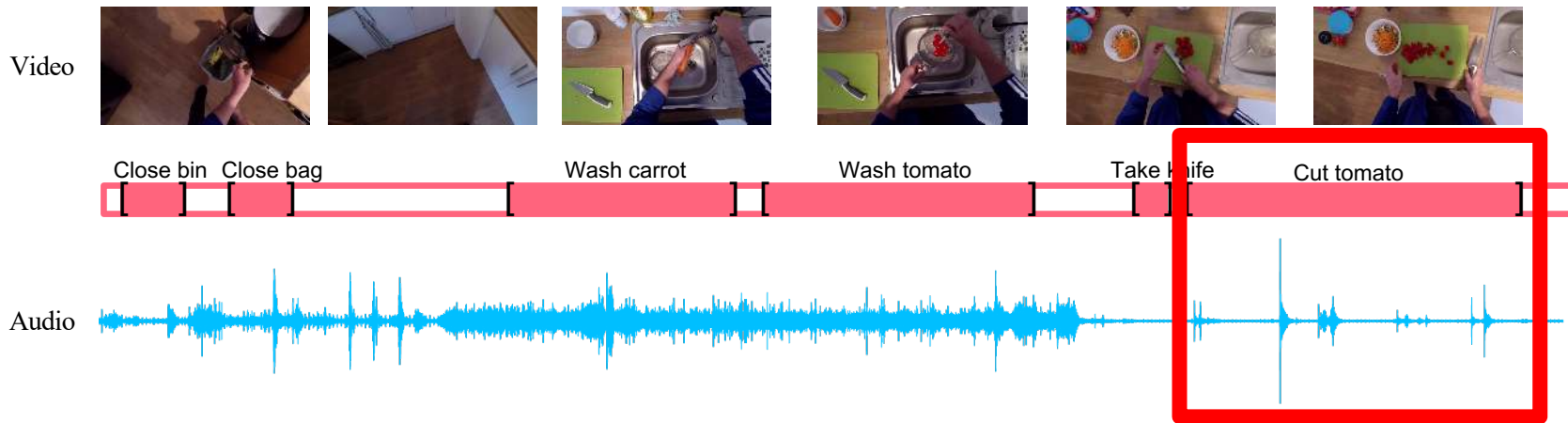
Motivation

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



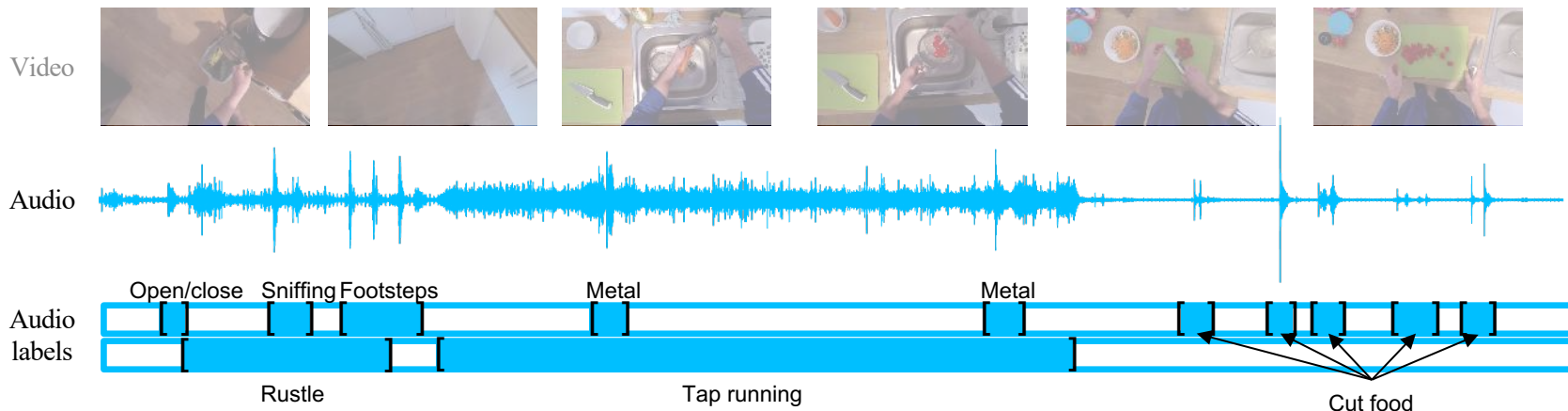
Motivation

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



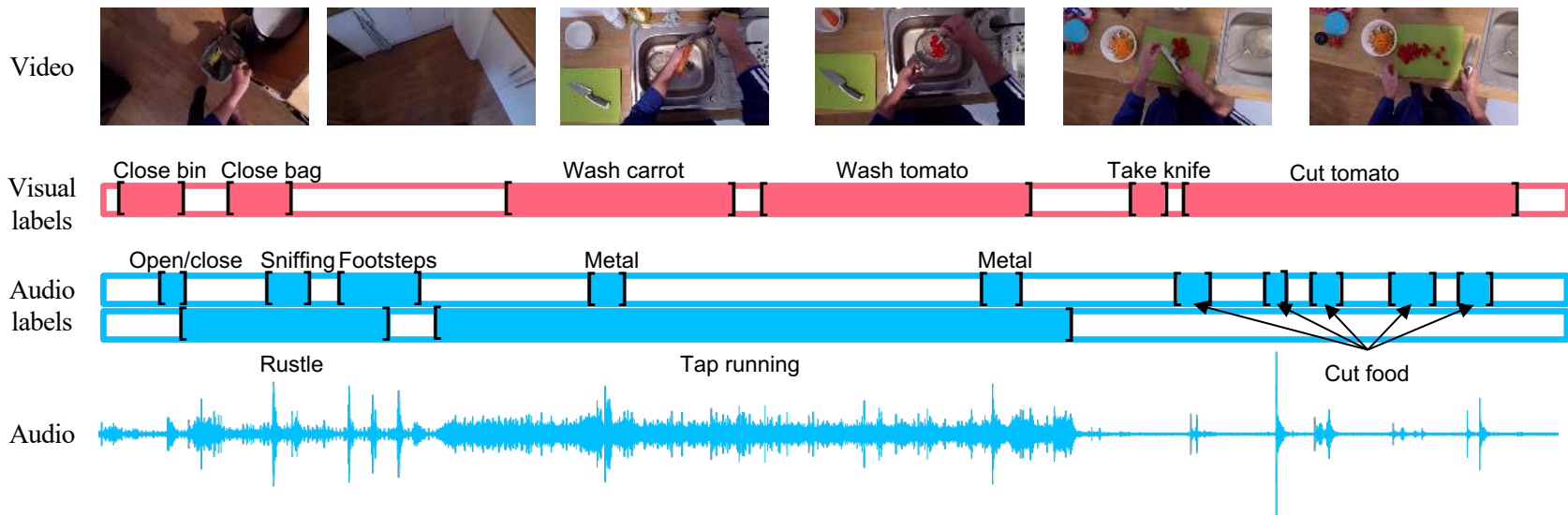
Motivation

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



Motivation

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

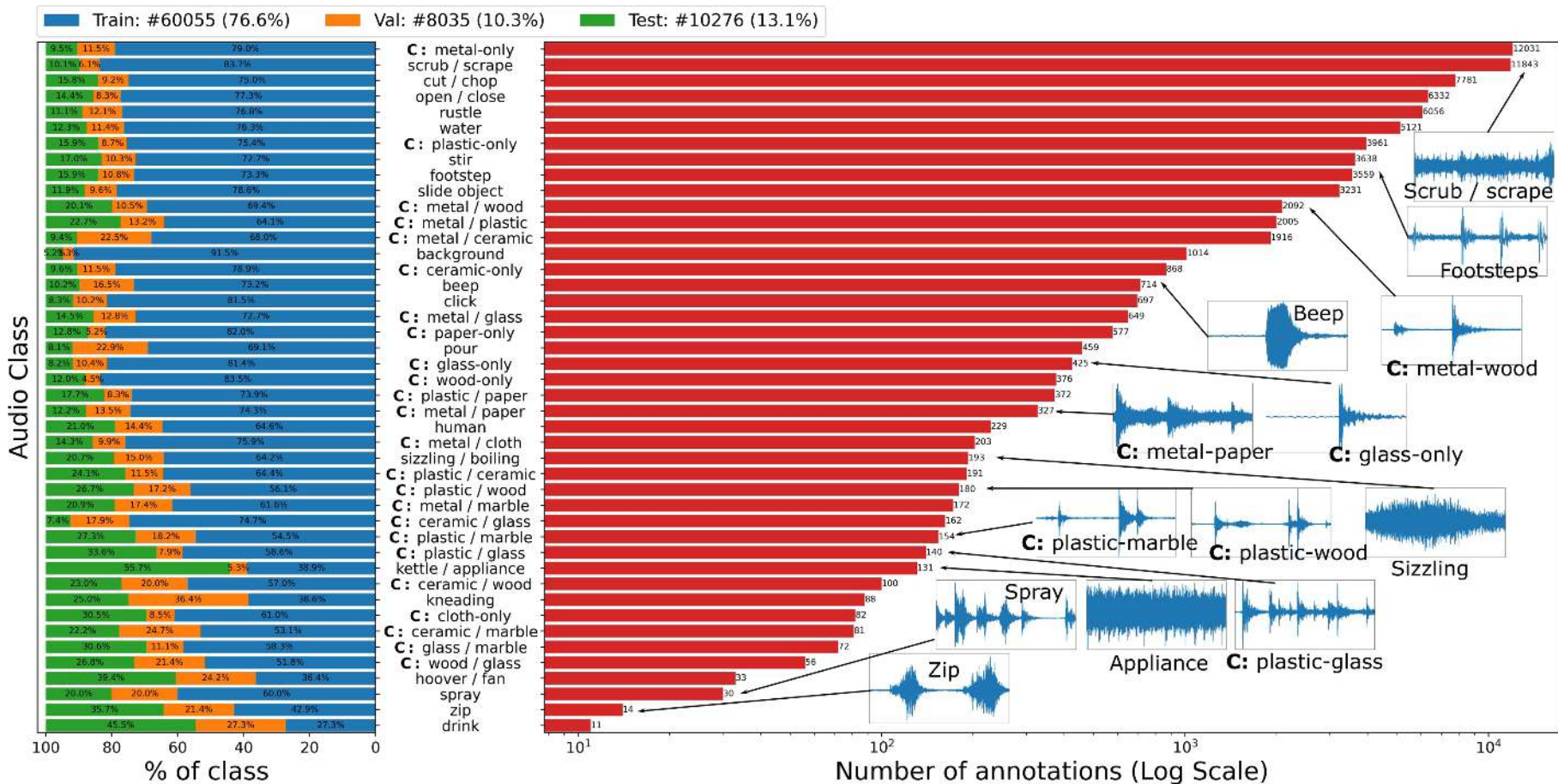


spray



EPIC-SOUNDS

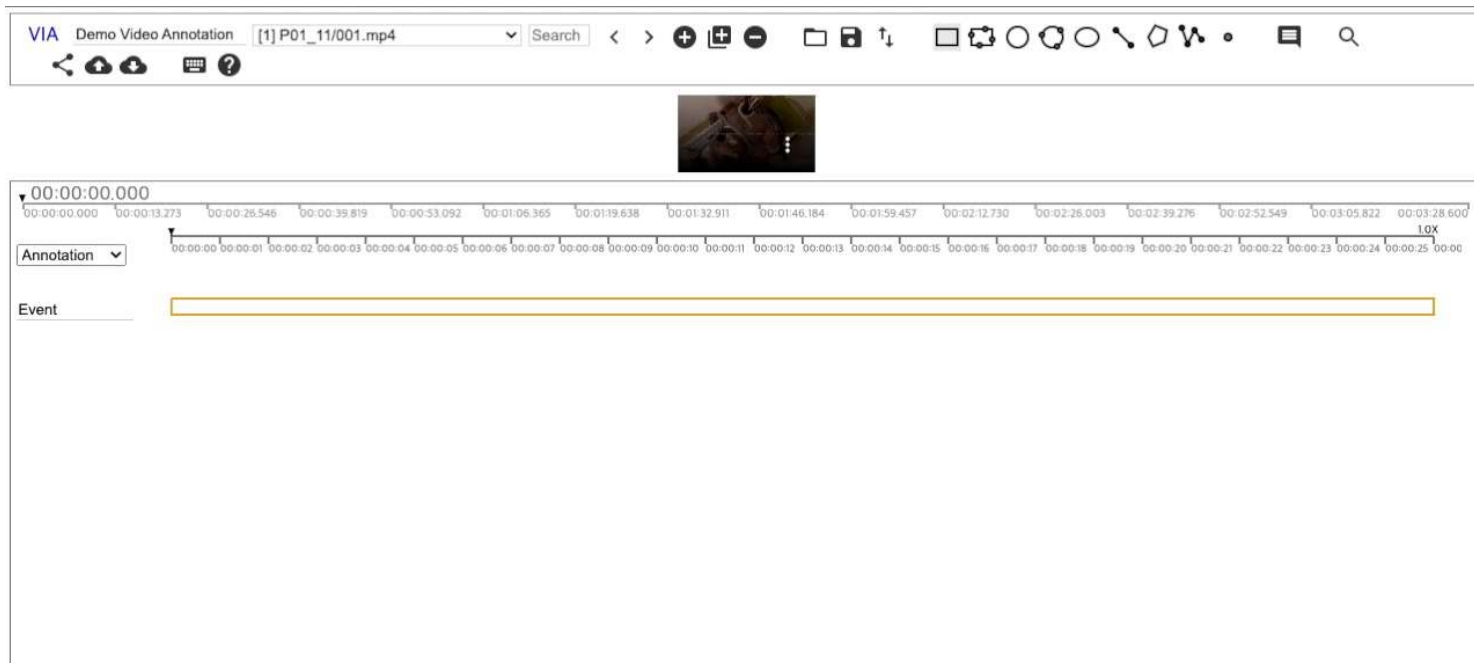
with: Jaesung Huh* & Jacob Chalk*
Vangelis Kazakos Andrew Zisserman



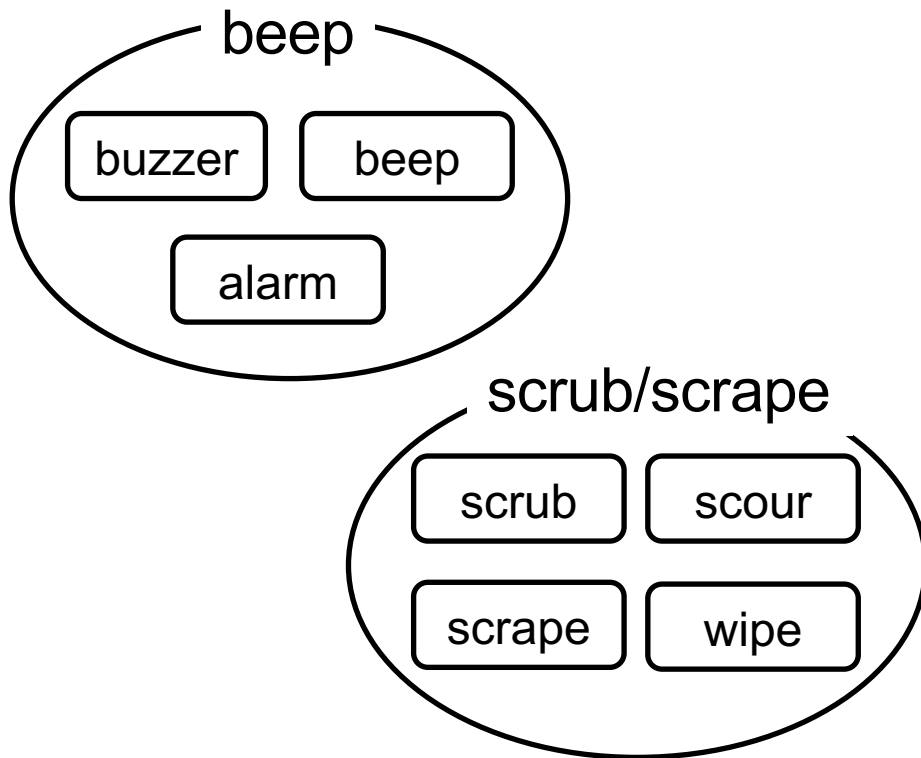
Annotations Pipeline

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 Image annotator tool



- From free-form descriptions to categories



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



Ceramic



Cloth



Metal



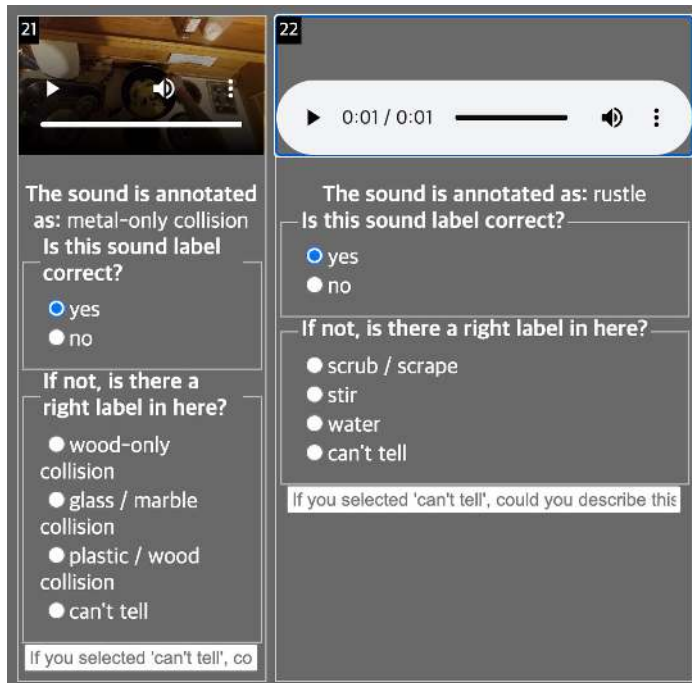
Plastic



Glass

- Manual check on validation / test set

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





Temporal labels are
modality-specific!

What type of labels can we provide?

- Temporal labels – Strong vs. Weak labels
- Semantic labels – Open-vocab. vs Closed-vocabulary
- Ranking labels – video-to-video comparisons
- Pixel-level labels – segmentation labels





Verb?

Noun?



sli.do

Joining as a participant?

#3639 120





Verbs:
add
pour
sprinkle
salt
season



Nouns:
salt
sea salt
seasoning
salt granules



sprinkle salt
season meat



Think of an example of an
opening action

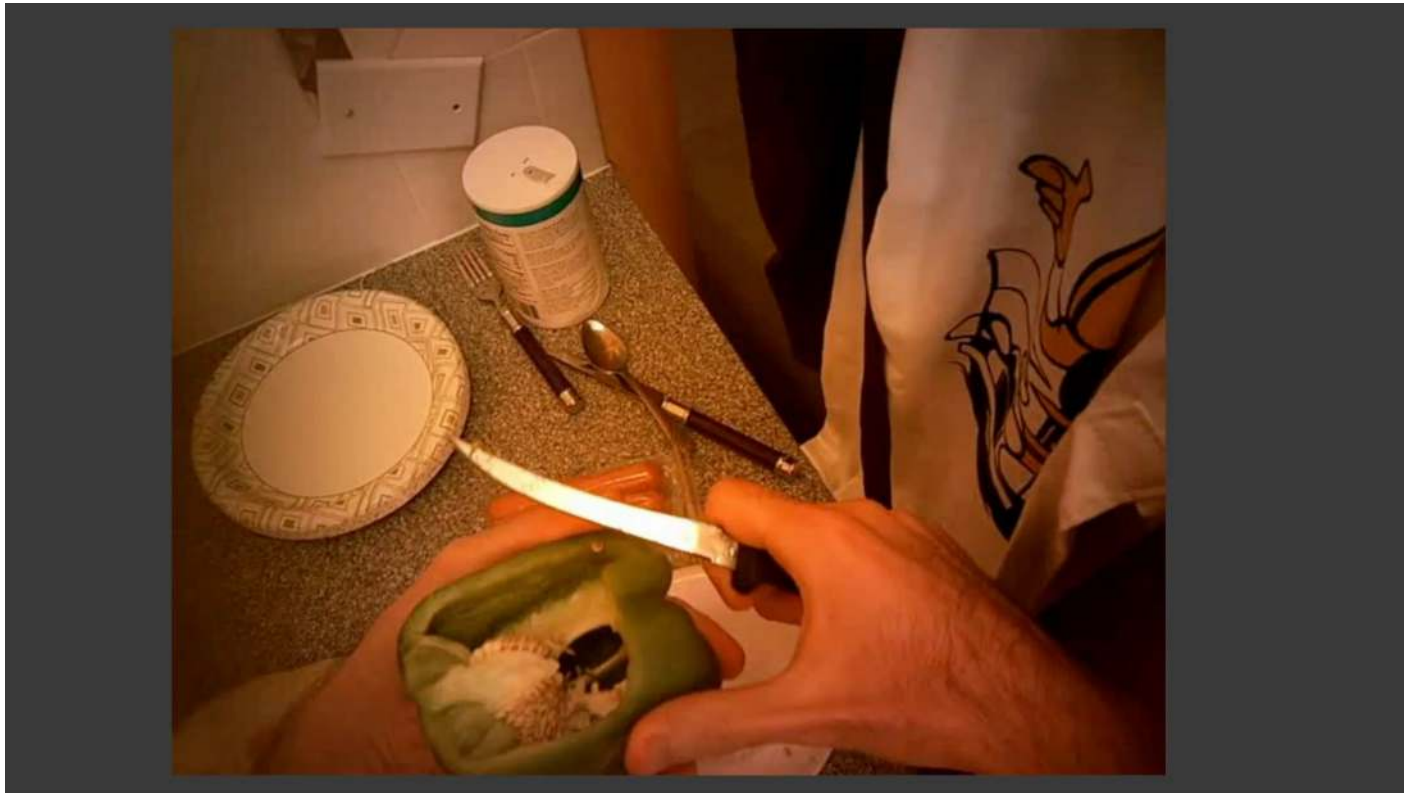


with: Michael Wray



Open





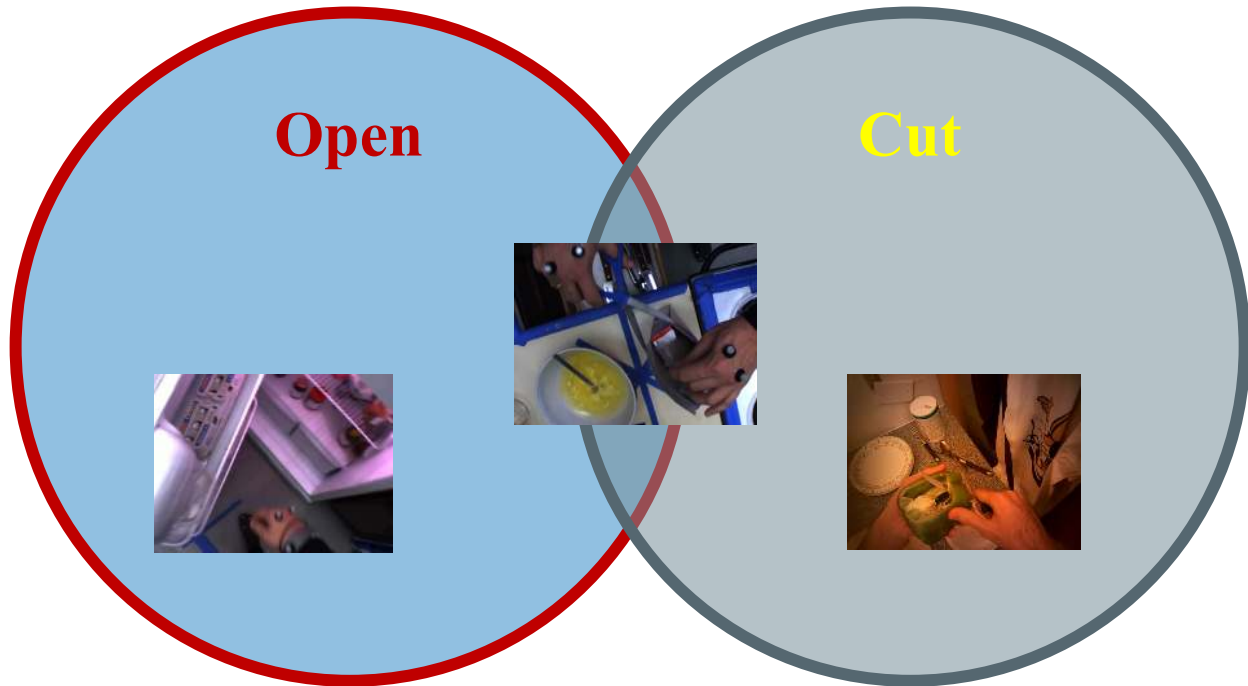
Open



Cut







Open

Cut



Towards an unequivocal rep. of actions

with: Michael Wray

- Action representations using a single verb is highly-ambiguous
 - Solution1: pre-selected non-overlapping verbs (SL)
 - run, walk, open, close
 - Solution2: Using nouns to disambiguate actions (V-N)
 - open-drawer, open-bottle, open-fridge
 - actions constrained to known nouns
 - Solution3: Multi-verb labels (ML, SAML)
 - open, hold, pull

Towards an unequivocal rep. of actions

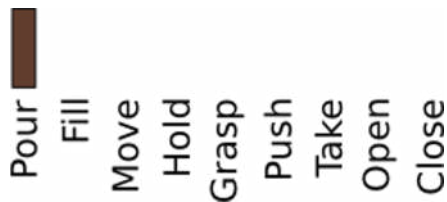
with: Michael Wray

- Collected from AMT



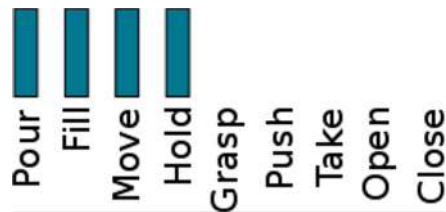
SL

- Majority Vote.
- One-hot vector.



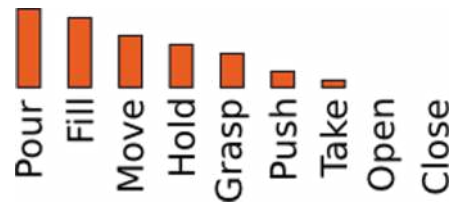
ML

- Threshold of 0.5.
- Binary Vector



SAML

- Full Annotation.
- Continuous Vector.



Top 3 retrieved classes across all datasets.

Turn On/Off
Press
Rotate



Turn On/Off
Press
Rotate



Labelling Method can differentiate turn On/Off tap by pressing and by rotating.



Semantics are harder than
you think...
There are significant
ambiguities

What type of labels can we provide?

- Temporal labels – Strong vs. Weak labels
- Semantic labels – Open-vocab. vs Closed-vocabulary
- Ranking labels – video-to-video comparisons
- Pixel-level labels – segmentation labels

Quality of Actions...



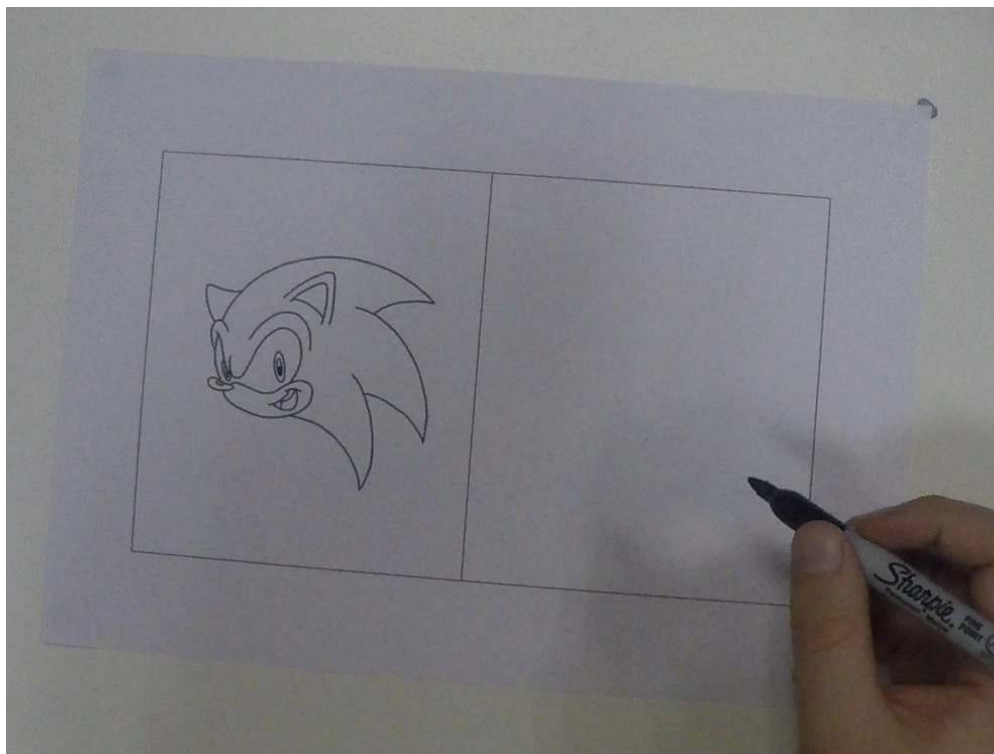
Pirsiavash et al, ECCV 2014



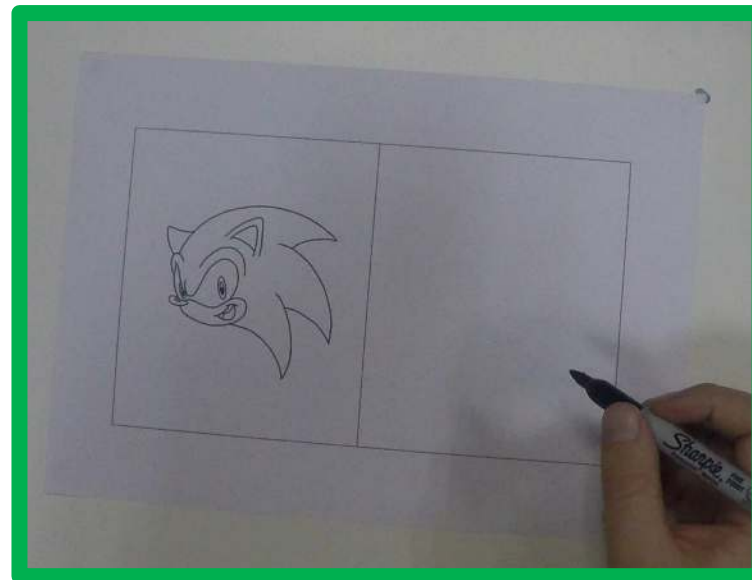
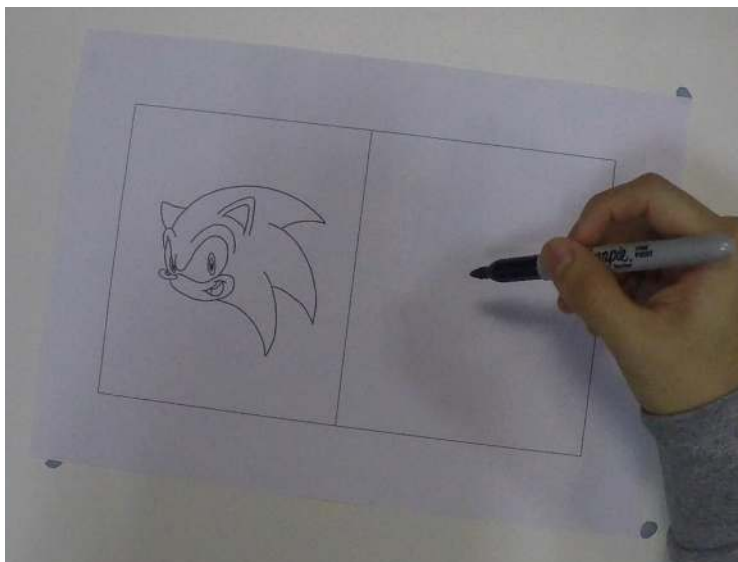
Shao et al, CVPR 2020

Quality of Actions...

with: Hazel Doughty
Walterio Mayol-Cuevas



Pairwise annotations of videos, indicating higher skill or no skill preference



What type of labels can we provide?

- Temporal labels – Strong vs. Weak labels
- Semantic labels – Open-vocab. vs Closed-vocabulary
- Ranking labels – video-to-video comparisons
- Pixel-level labels – segmentation labels

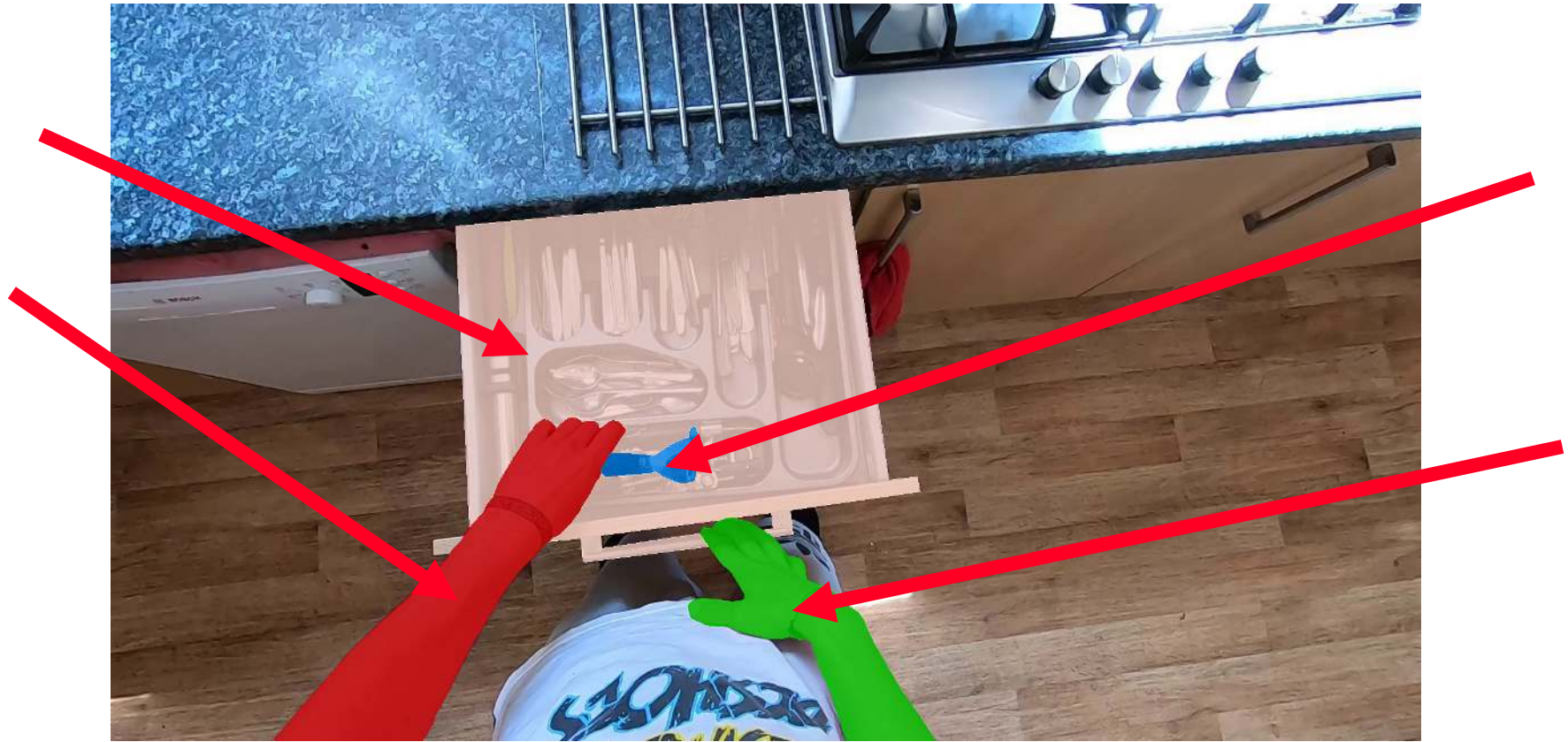
EPIC-KITCHENS VISOR

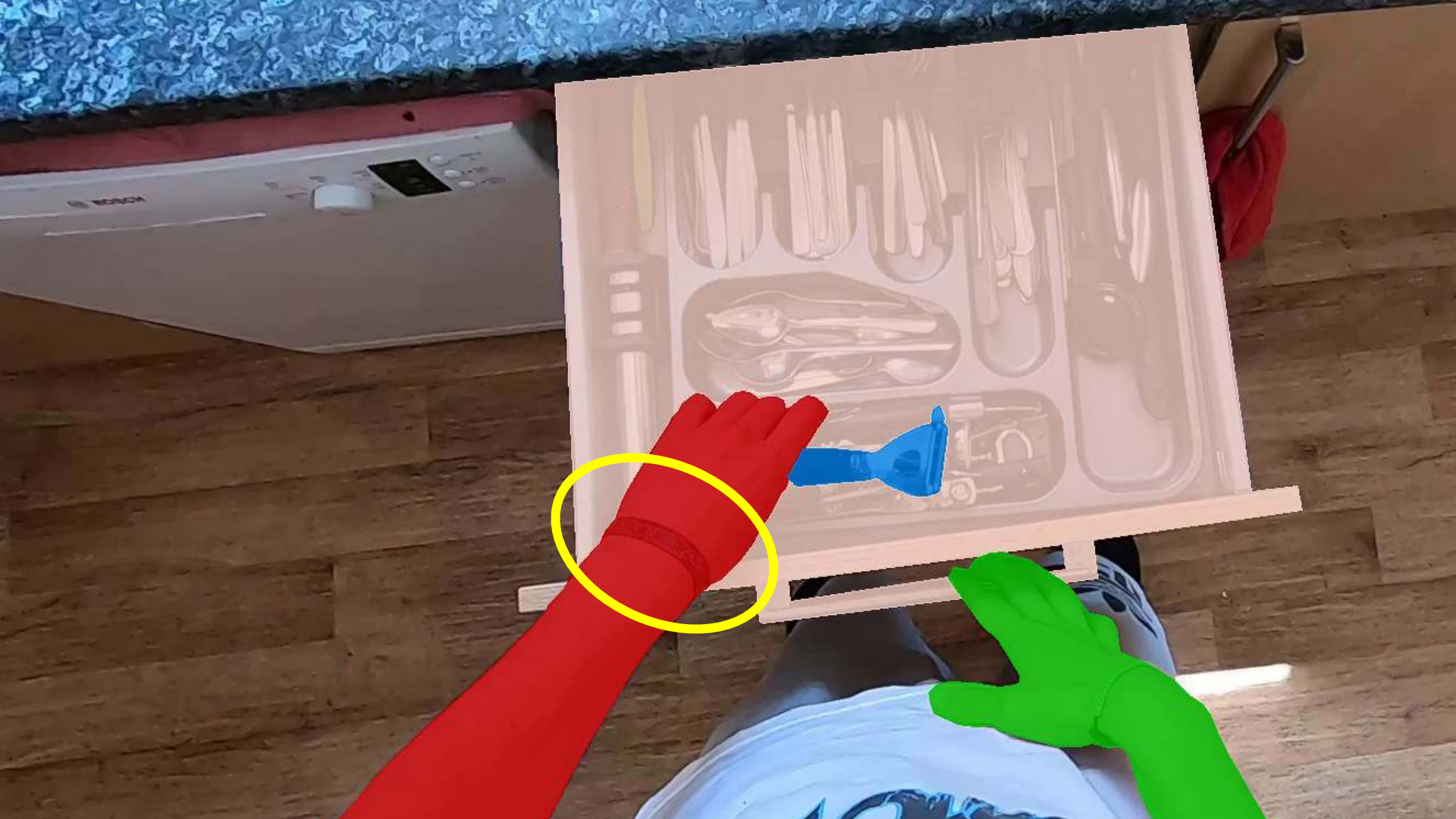
with: Ahmad Darkhalil, Dandan Shan, Bin Zhu, Jian Ma,
Amlan Kar, Richard Higgins, David Fouhey, Sanja Fidler




EPIC-KITCHENS VISOR

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







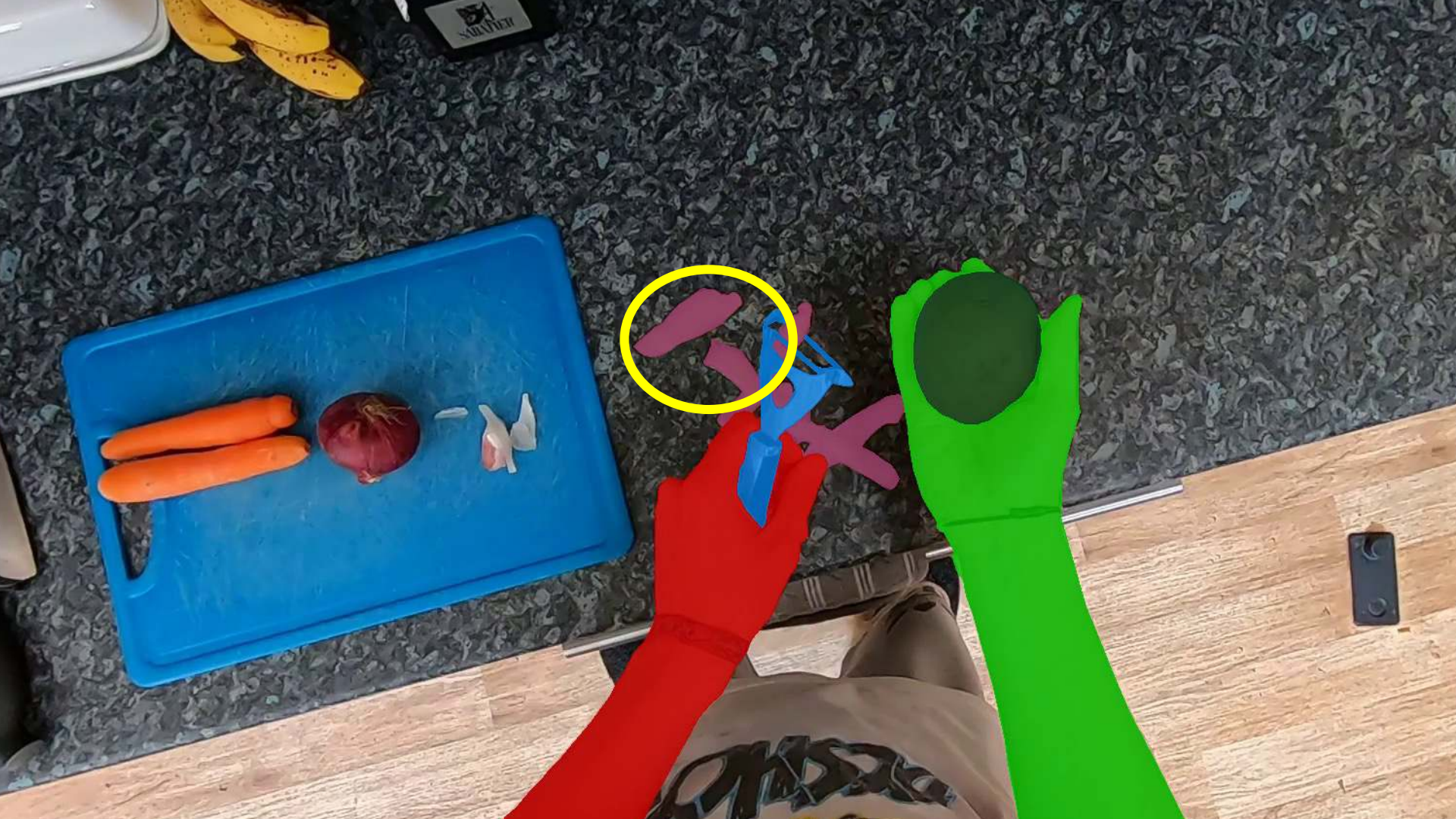


Depth 0 : peeler
Depth -1 : drawer

EPIC-KITCHENS VISOR

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





EPIC-KITCHENS VISOR

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



EPIC-KITCHENS VISOR

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



EPIC-KITCHENS VISOR

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



EPIC-KITCHENS VISOR

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



EPIC-KITCHENS VISOR

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



pour spice



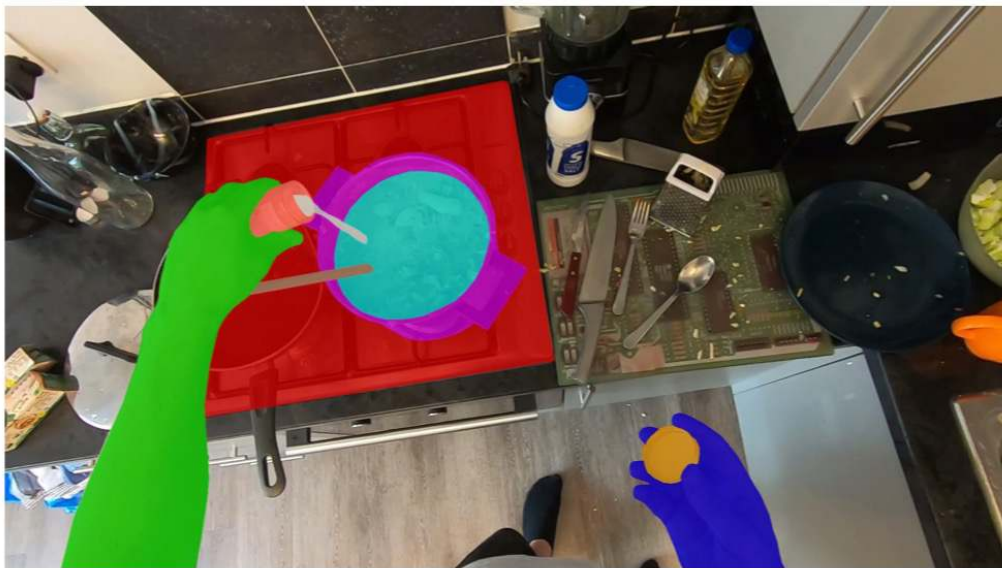
- left hand
- right hand
- hob
- saucepan
- spice
- spice container
- spoon
- soup
- pepper container lid

pour spice



- left hand
- right hand
- hob
- saucepan
- spice
- spice container
- spoon
- soup
- pepper container lid

pour spice



- left hand
- right hand
- hob
- saucepan
- spice
- spice container
- spoon
- soup
- pepper container lid

saucepan → pan → cookware
spoon → spoon → cutlery

pour spice ← action



- left hand
- right hand
- hob
- saucepan
- spice
- spice container
- spoon
- soup
- pepper container lid

in-contact (spice container) in-contact (container lid)

pour spice



- left hand
- right hand
- hob
- saucepan
- spice
- spice container
- spoon
- soup
- pepper container lid

spoon (non-exhaustive)

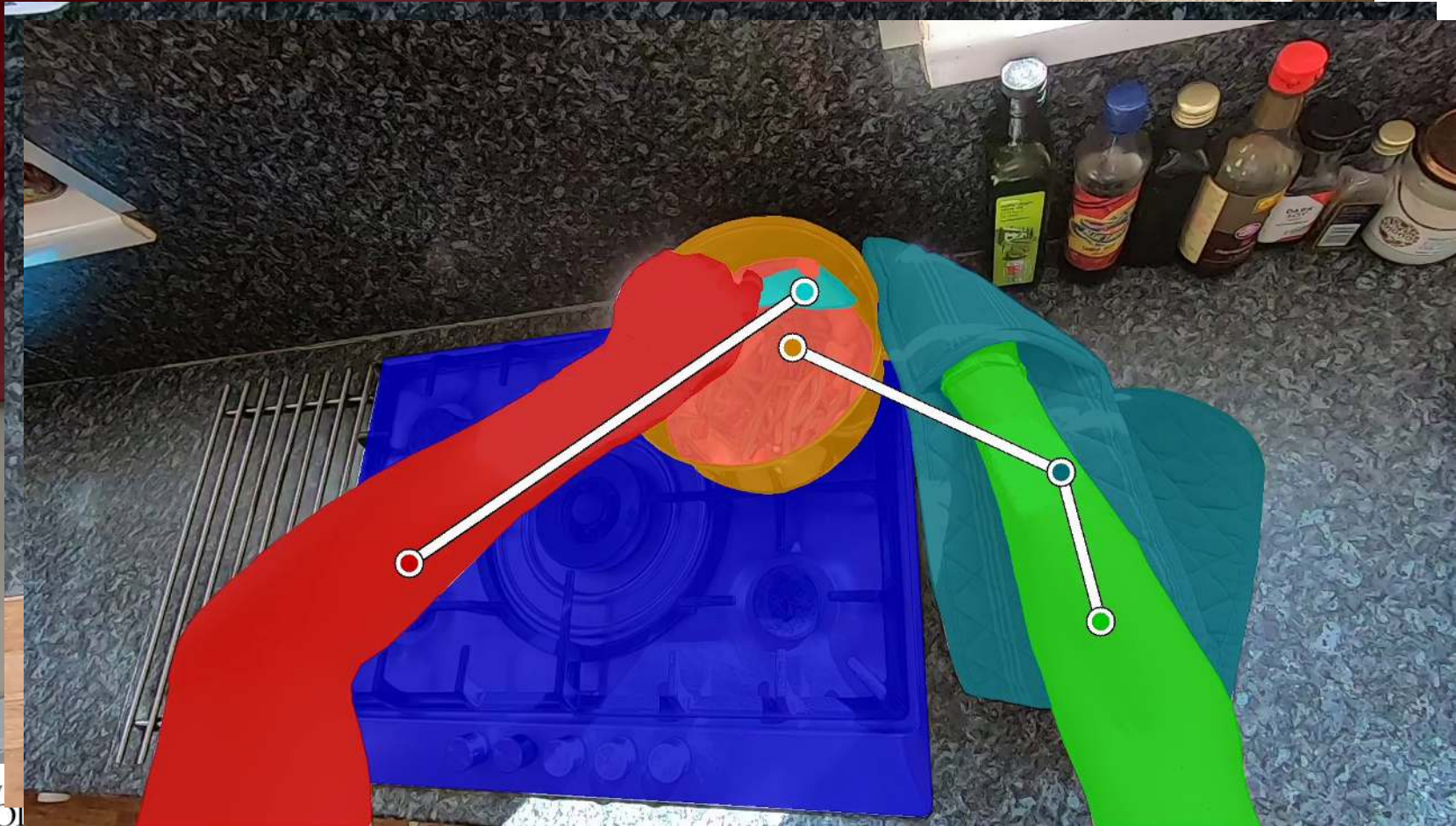
EPIC-KITCHENS VISOR

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



VISOR Relations

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



Object relation stats

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

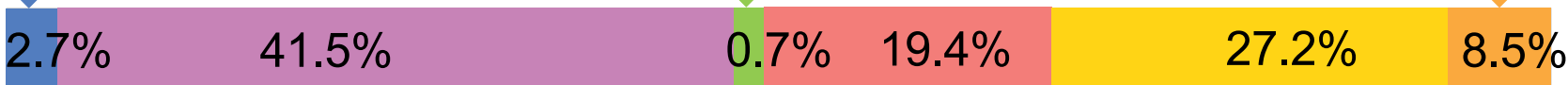
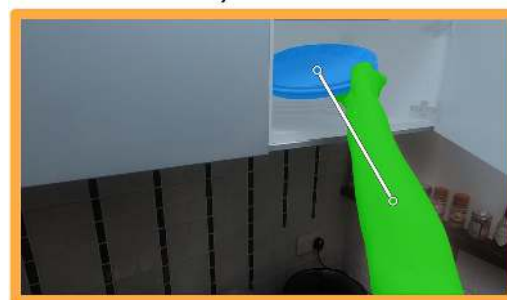
1 Hand, No Contact



2 Hands, No Contact



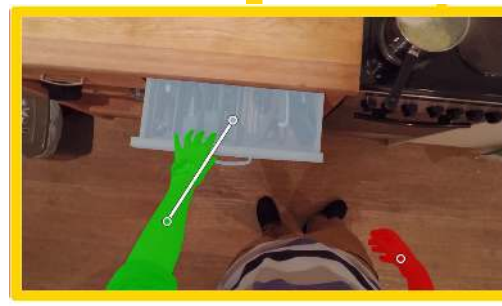
1 Hand, In Contact



2 Hands, 2 Obj Contacts



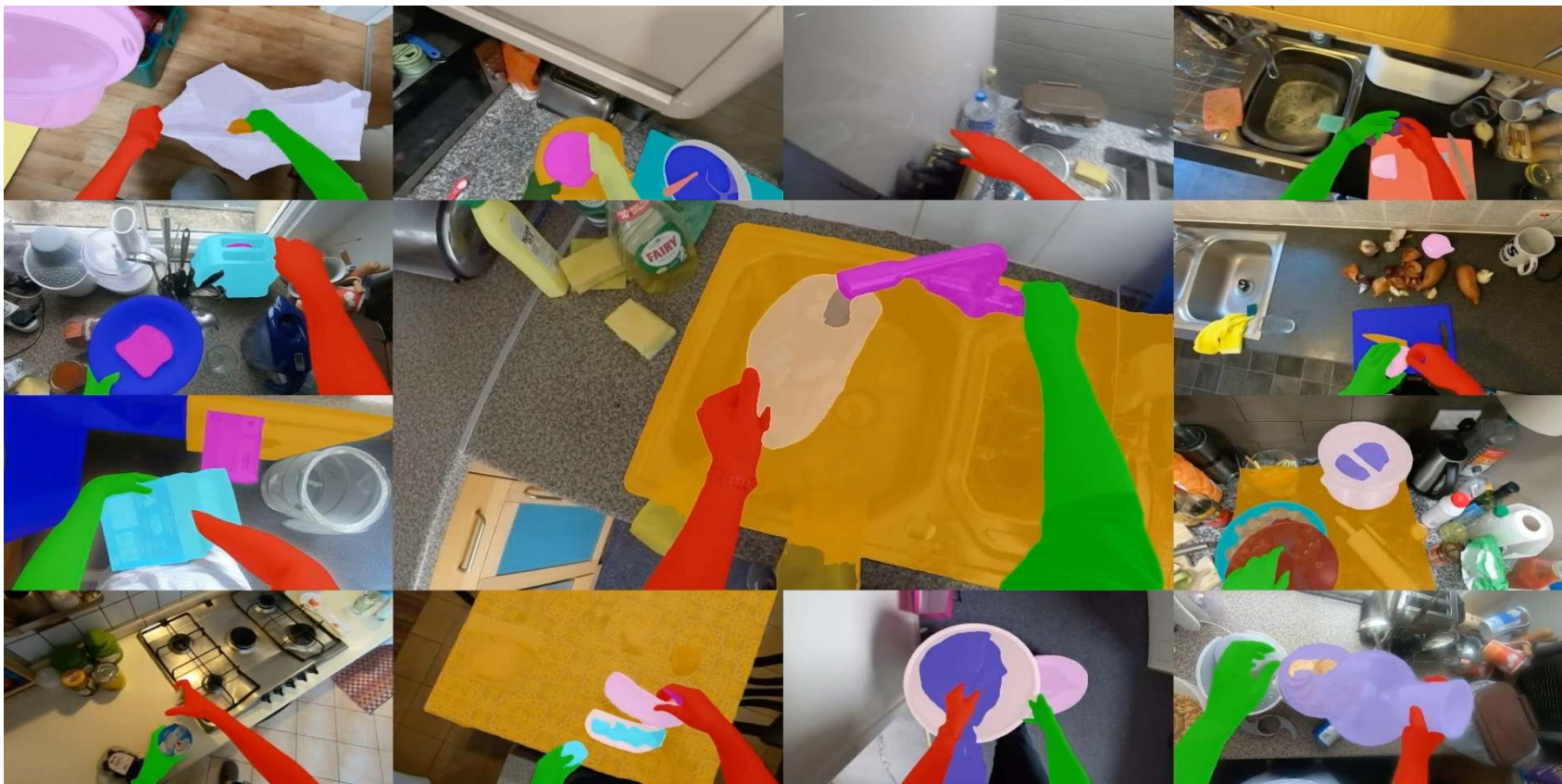
2 Hands, Same Contact



2 Hands, 1 In Contact

EPIC-KITCHENS VISOR

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





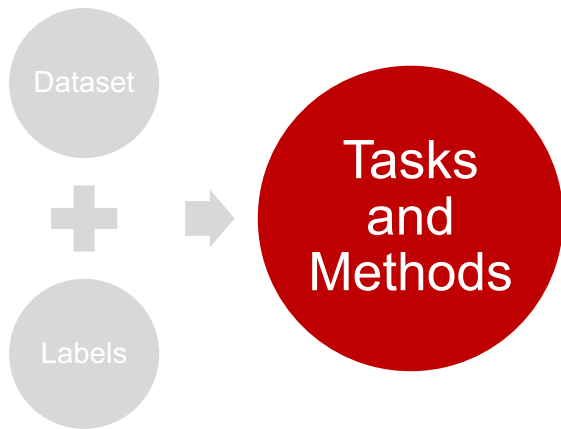
Semantic segmentations
during transformations is...
very challenging



There is no Ground Truth – only biased labels

Data
bi

are always
researchers'



Part III: Tasks and Methods

Sampling frames...



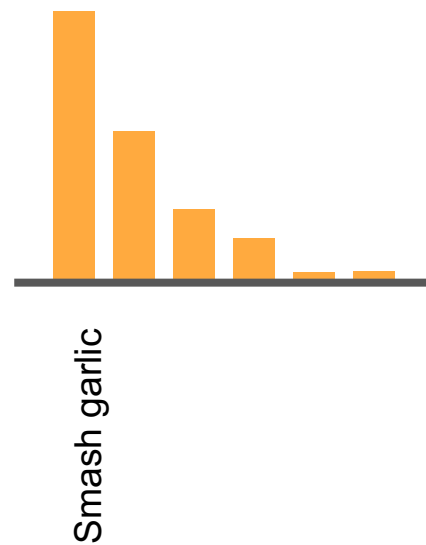
Sampling frames...



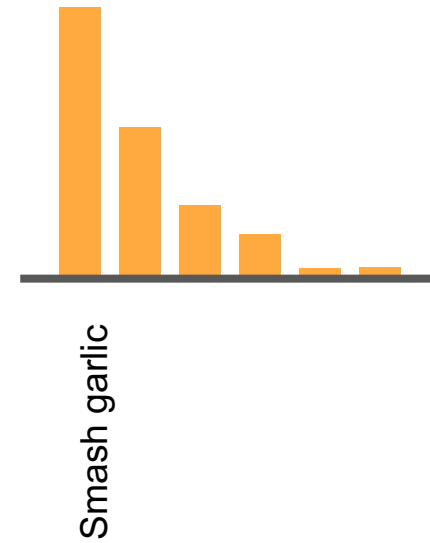
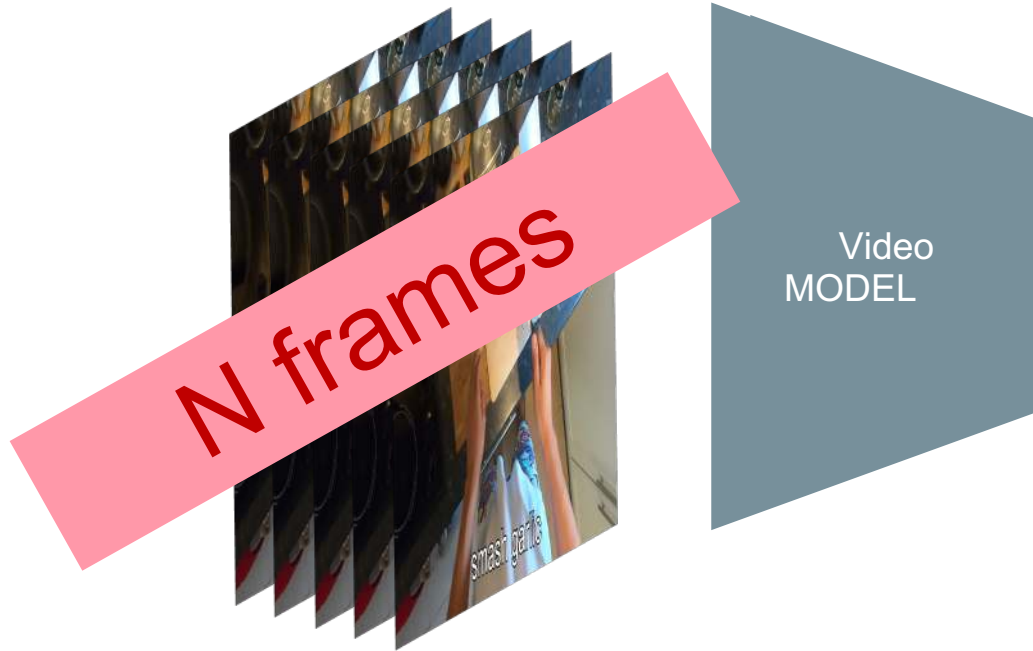
Sampling frames...



MODEL



Sampling frames...



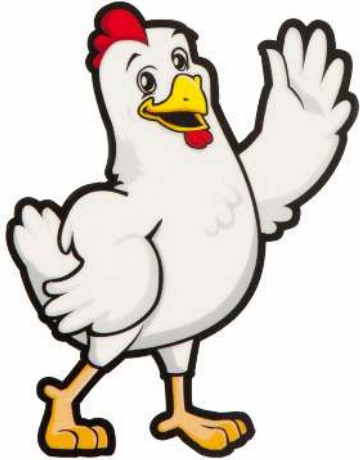
Sampling frames...



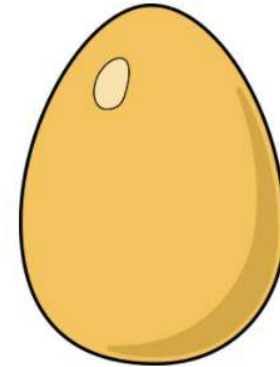
Sampling frames...



Frames to select



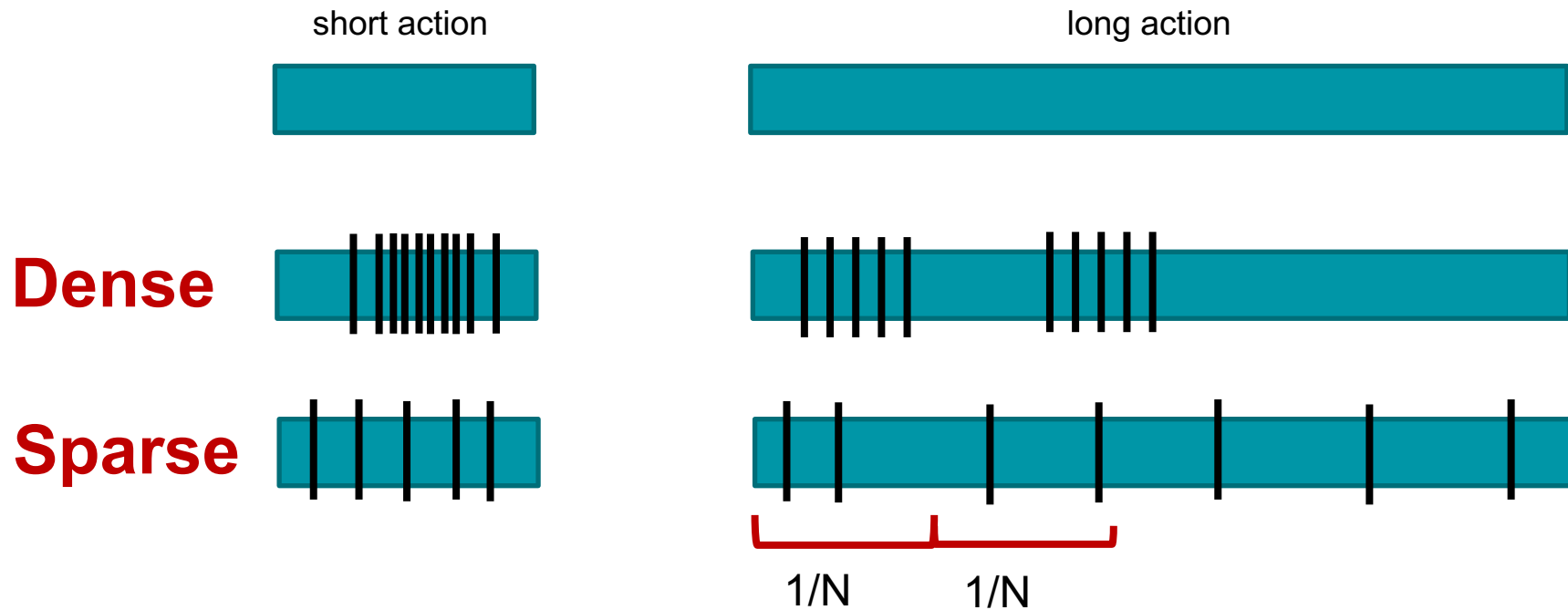
Action to recognise





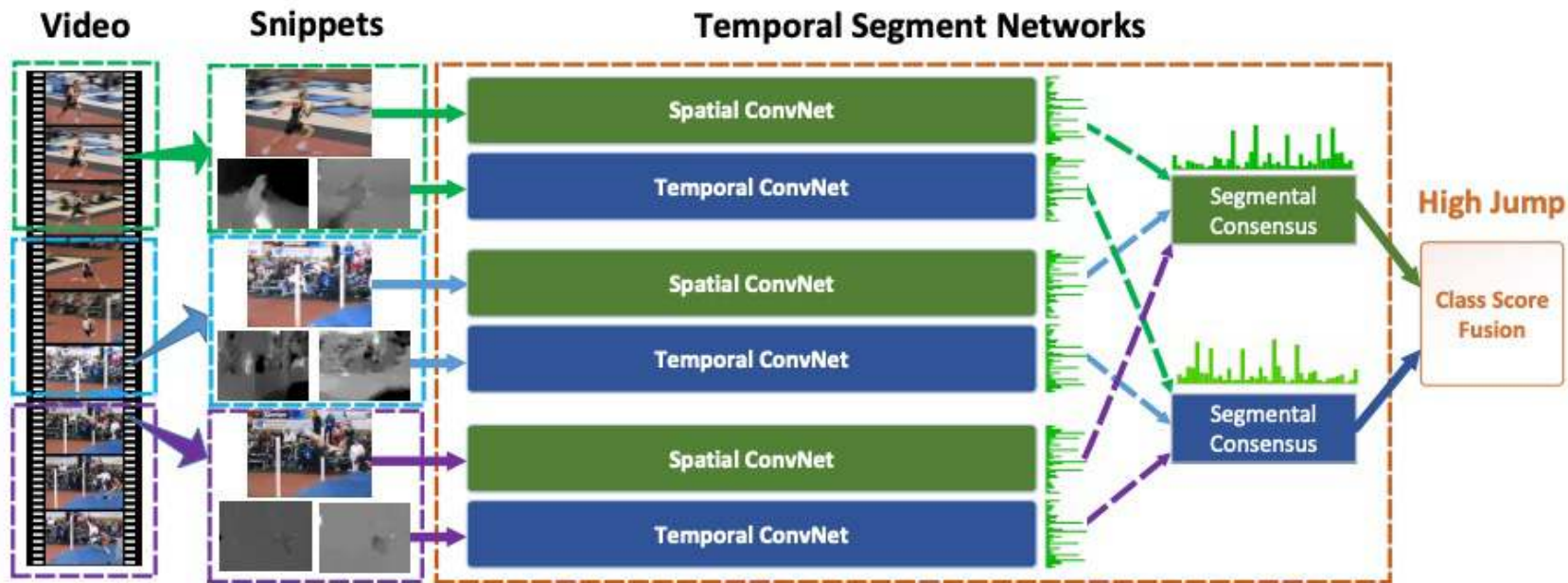
All models and methods
sample frames...
Sampling is often hidden in
implementation details...
It is **critical** ...

Two sampling approaches



Sparse sampling

Temporal Segment Networks (TSN) – Wang et al, ECCV 2016



Two sampling approaches

Dense

Better motion features

Short motion signature

Easier to implement

Better for cross-dataset generalisation

Sparse

More complex features

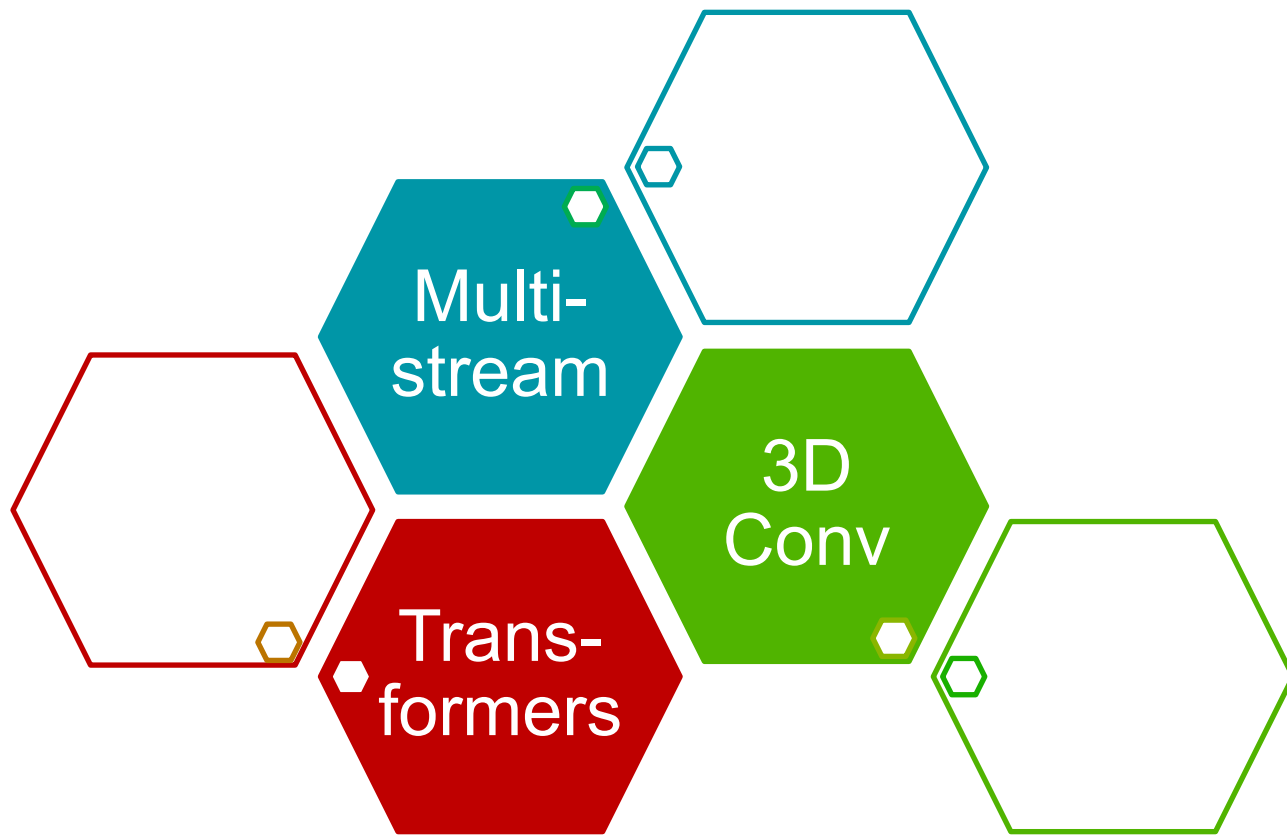
Complete action representation

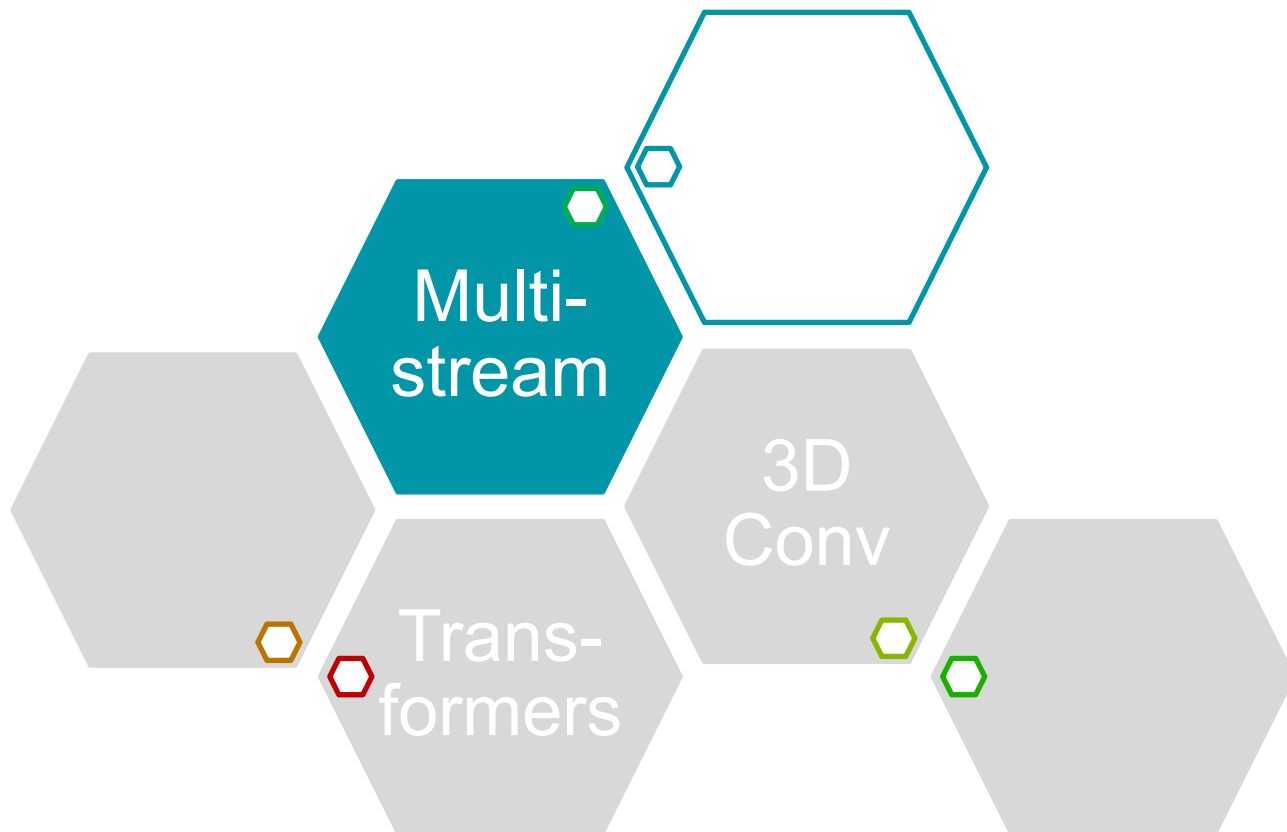
More augmentations

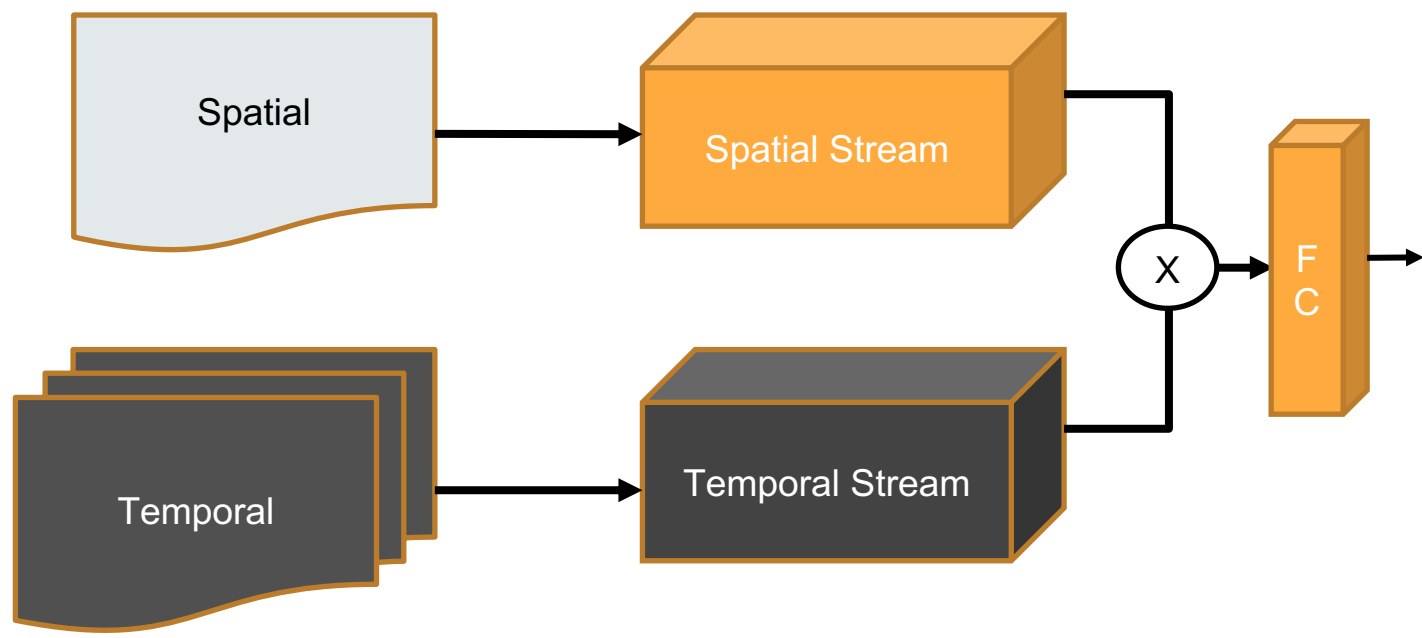
Better for *temporal* datasets

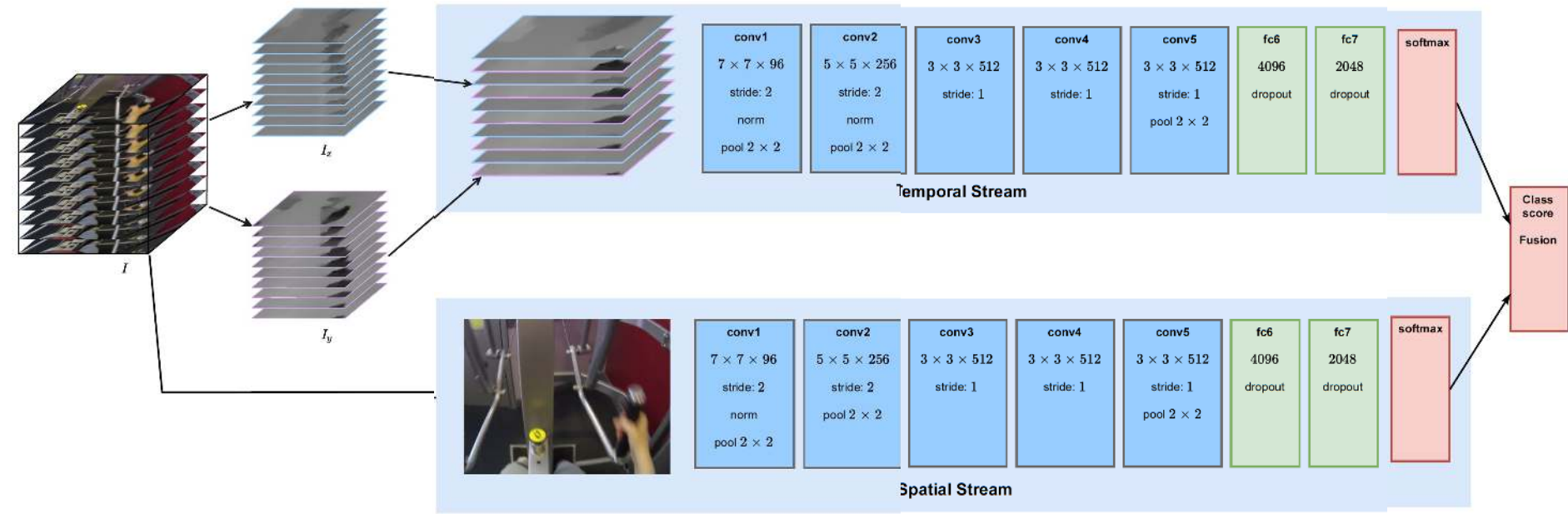


Models









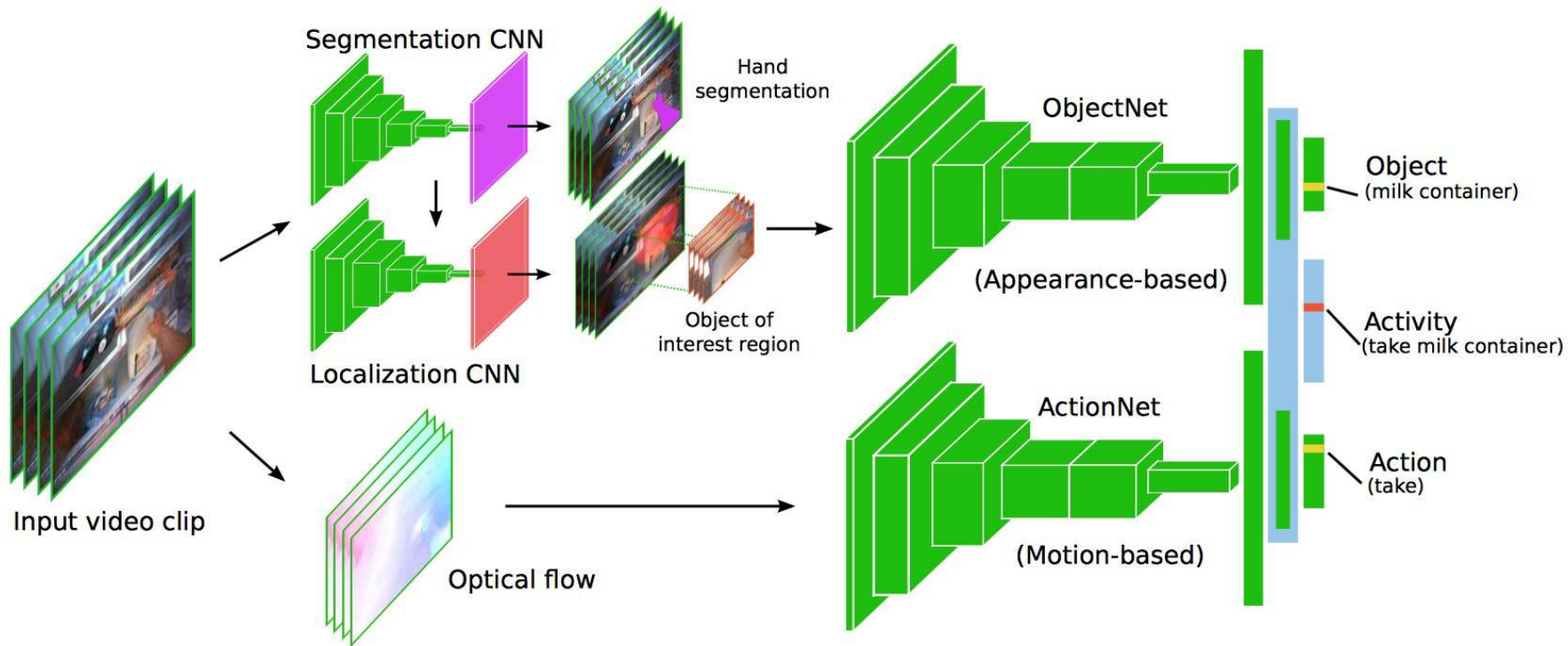
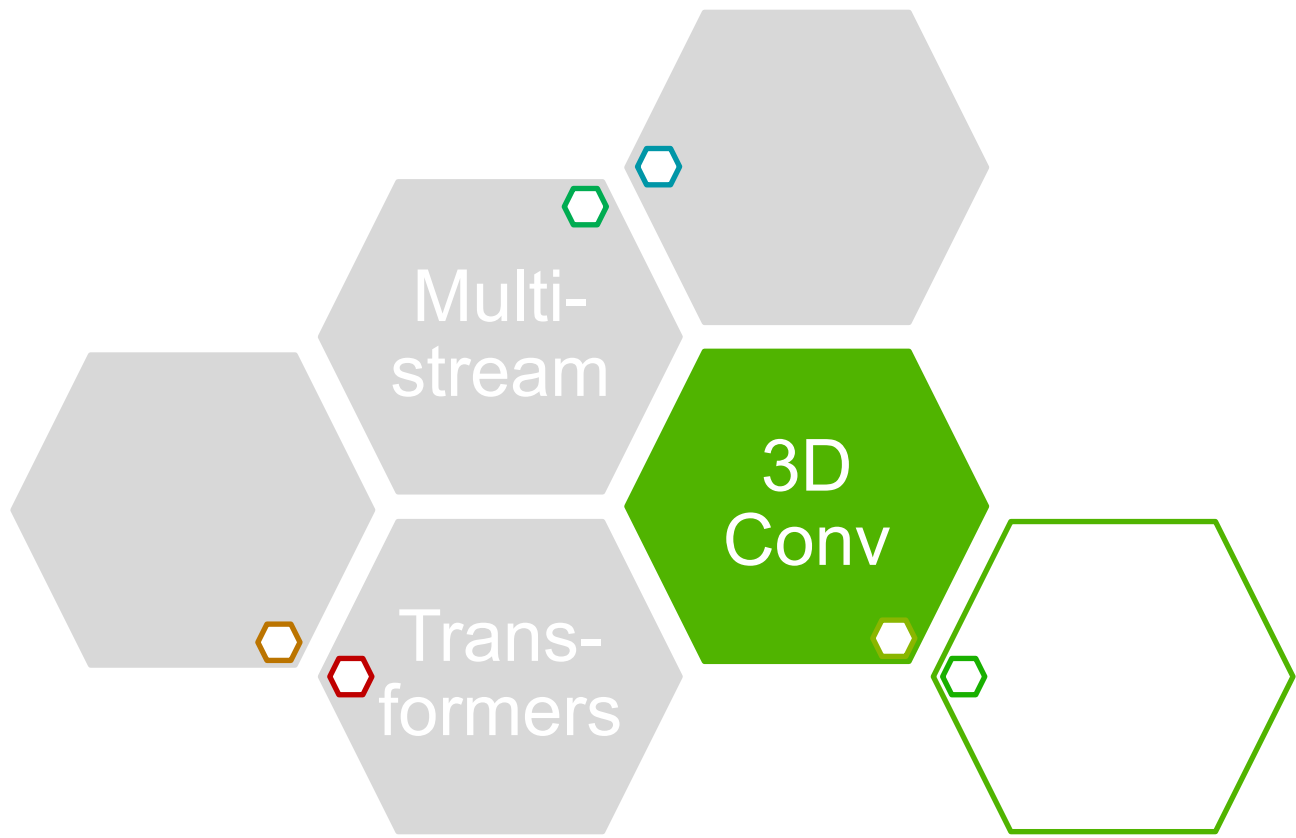
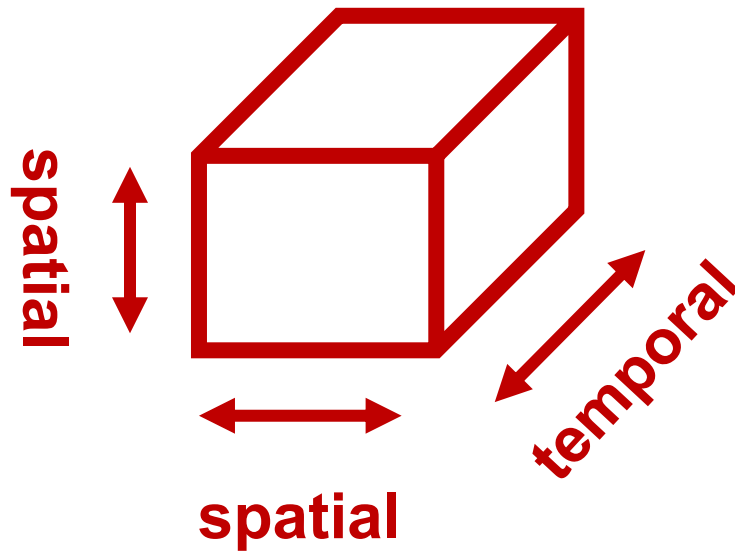


Figure from: Ma et al. Going Deeper into First-Person Activity Recognition. CVPR 2016



3D Convolutional



3D Convolutional

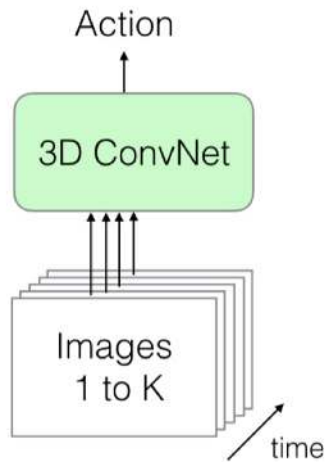


- Initial attempts required initialisation from 2D Networks (i.e. ImageNet)
 - No presence of large scale video dataset
 - Inflated networks (I3D) Carriera et al

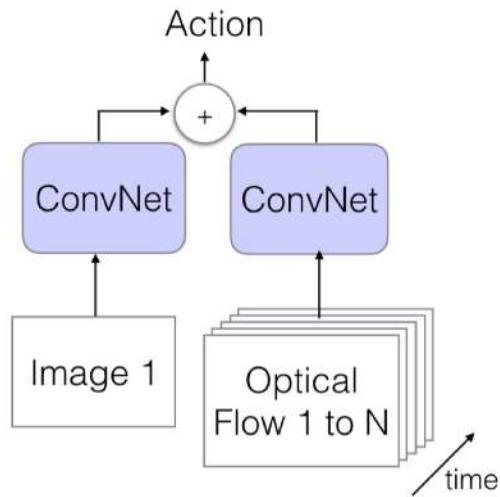
- The objective was to remove the need for optical flow...

3D Convolutional

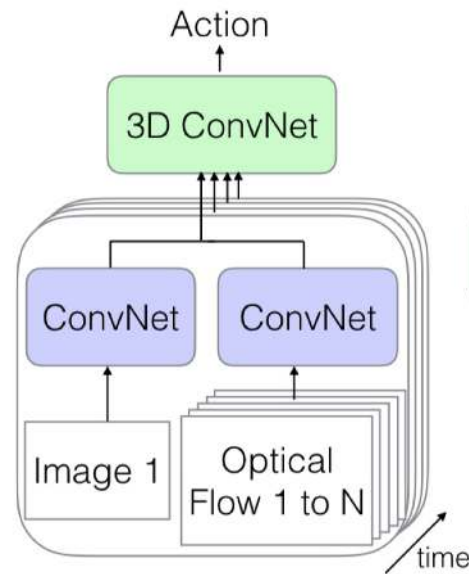
b) 3D-ConvNet



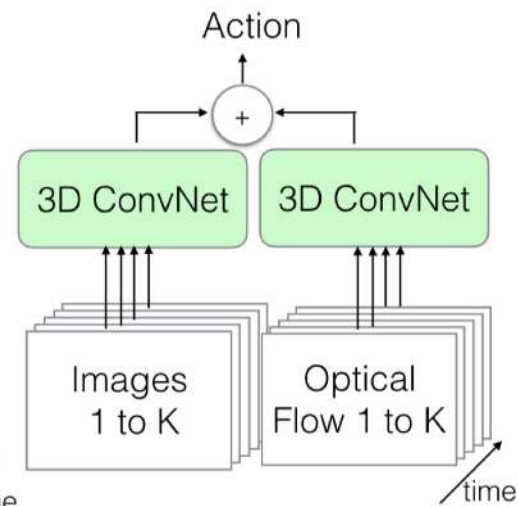
c) Two-Stream



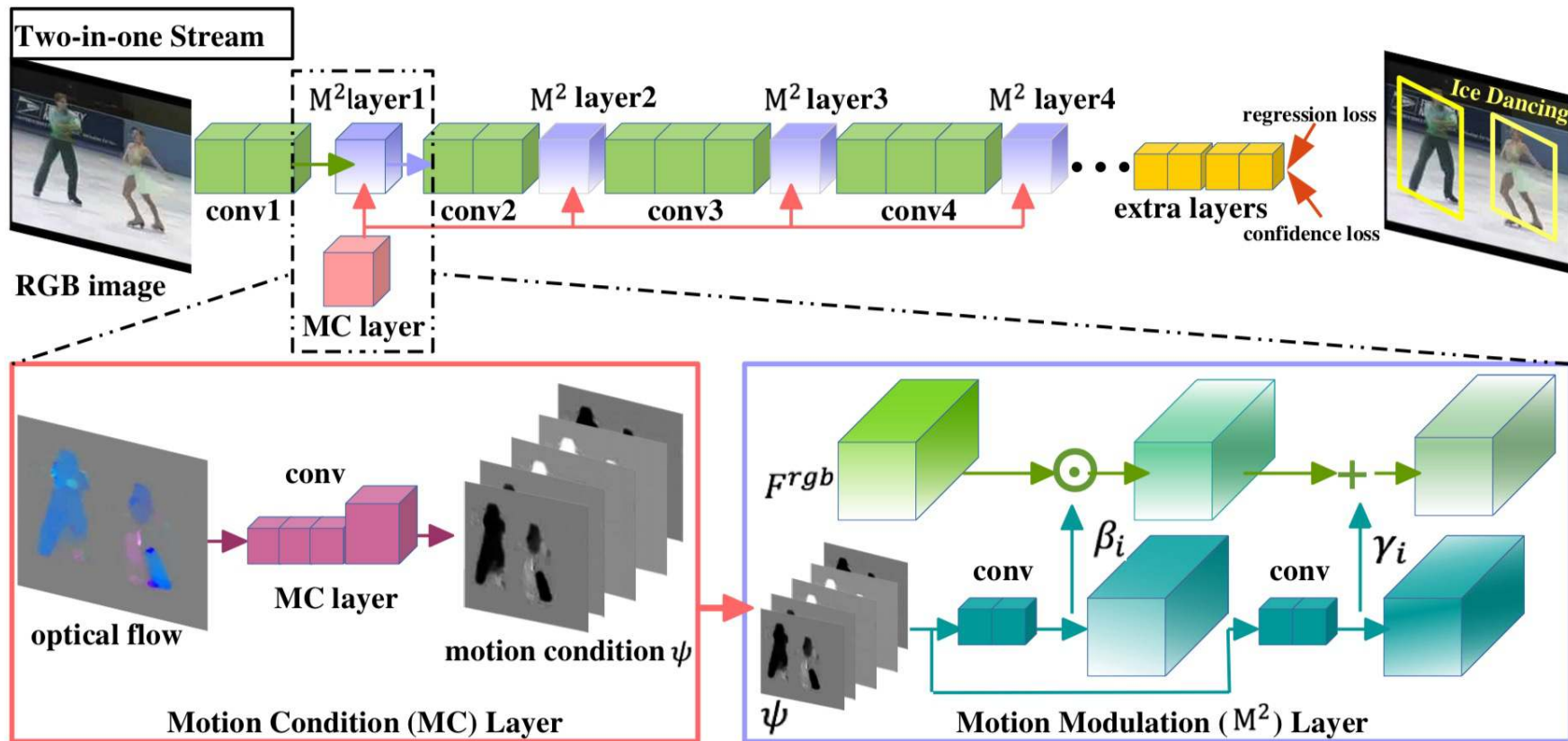
d) 3D-Fused Two-Stream

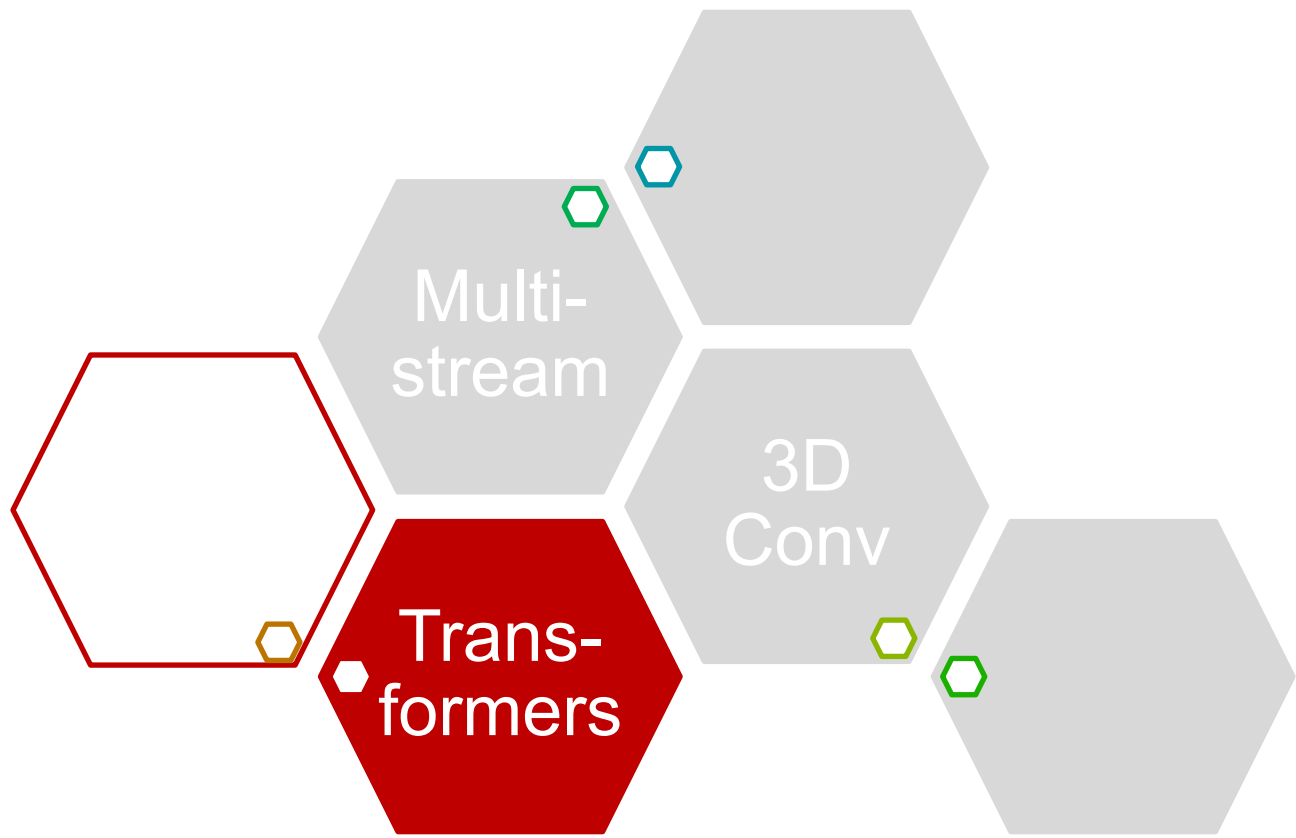


e) Two-Stream 3D-ConvNet



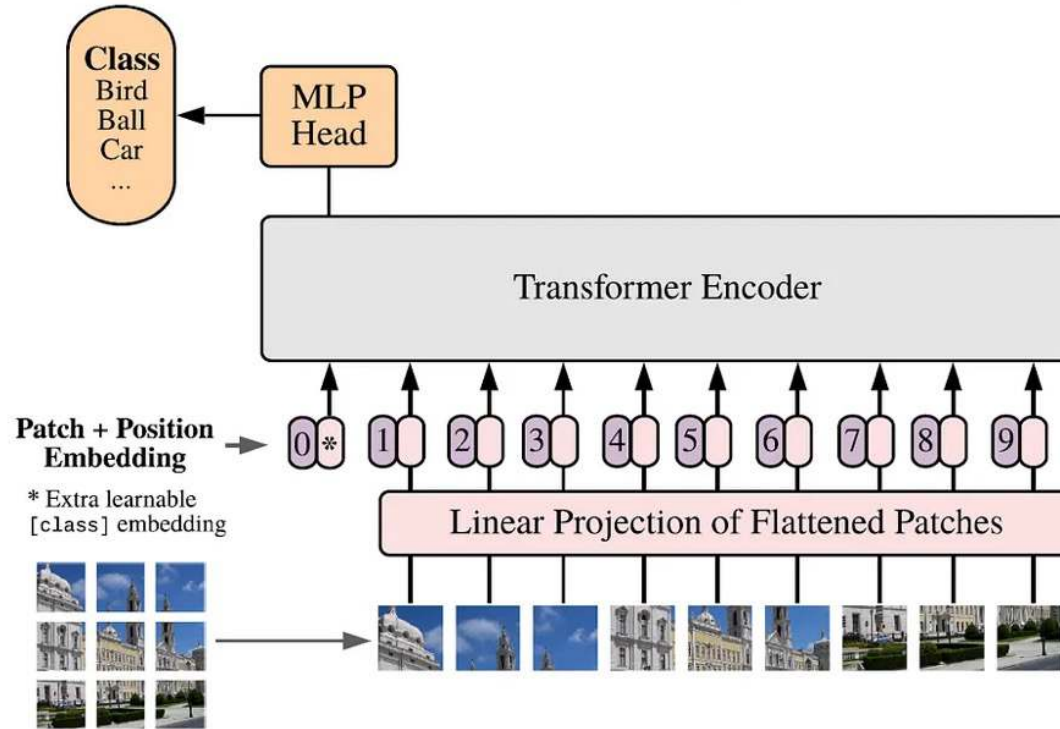
Via knowledge distillation



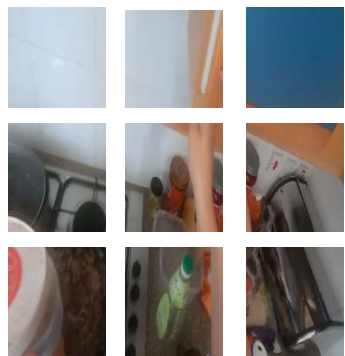




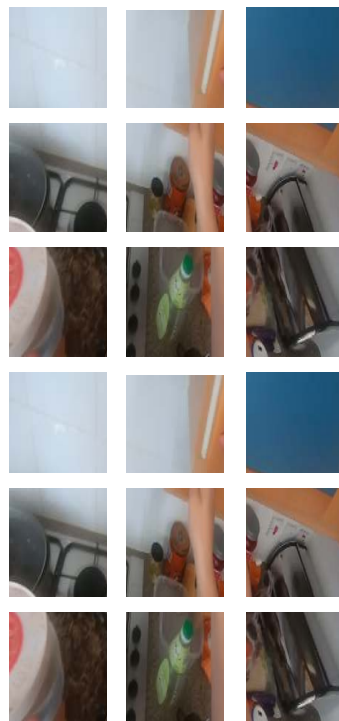
Vision Transformer (ViT)



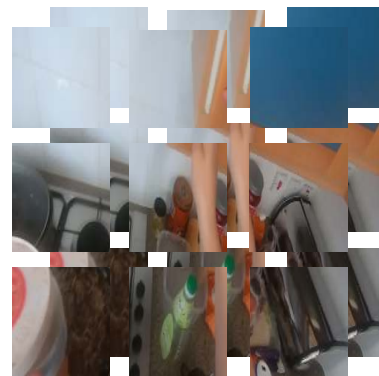
How to Patch-ify a Video?



$H \times W \times C$

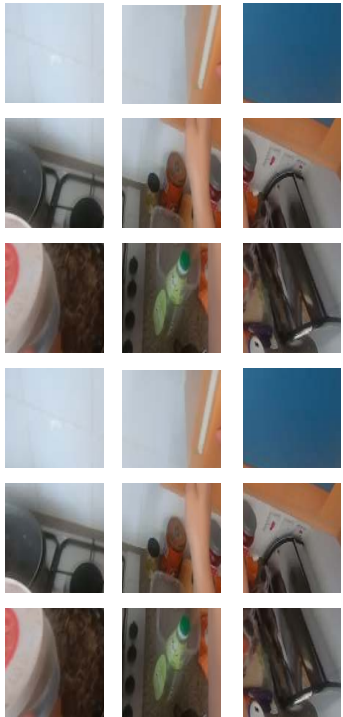


$[TH] \times W \times C$

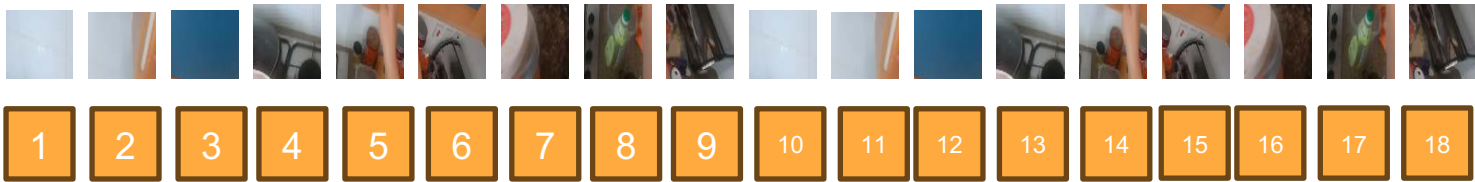


$H \times W \times [CT]$

How to Patch-ify a Video?

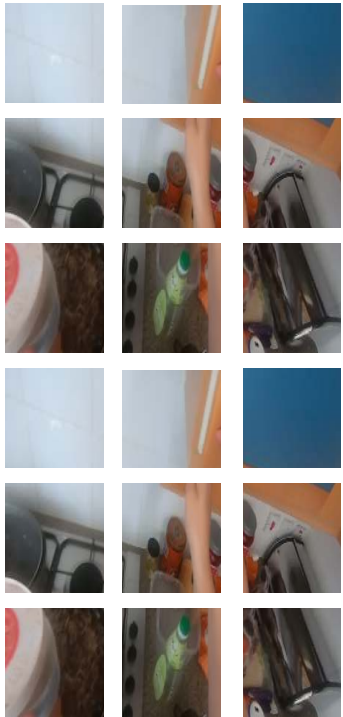


Flatten

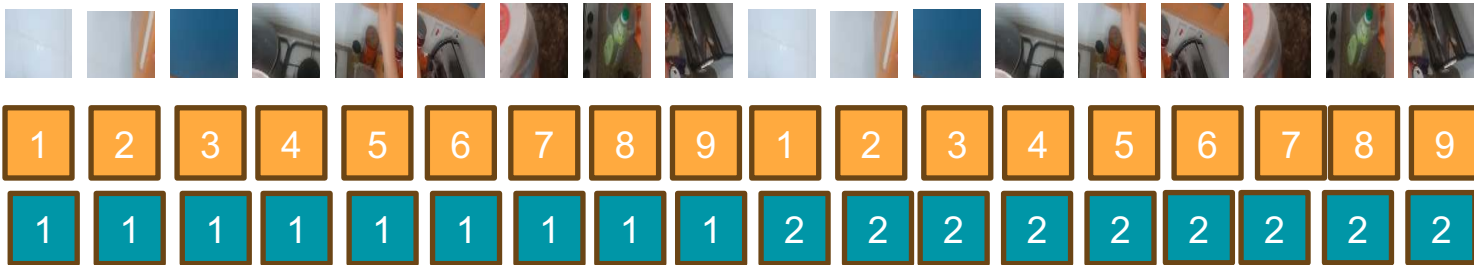


[TH] x W x C

How to Patch-ify a Video?

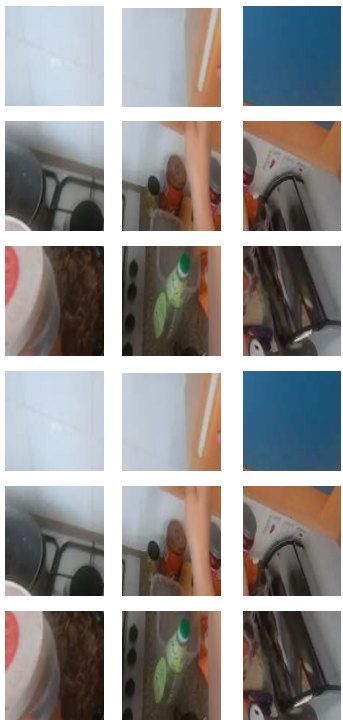


Flatten



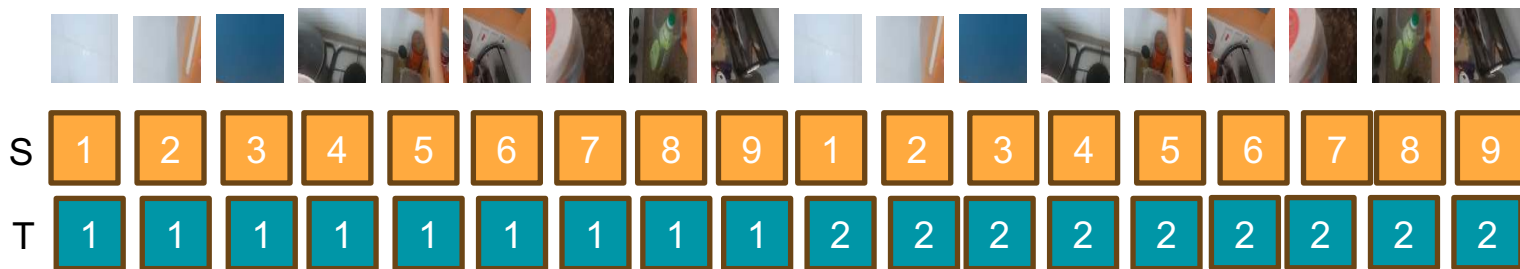
[TH] x W x C

How to Patch-ify a Video?



[TH] x W x C

fully-connected layer



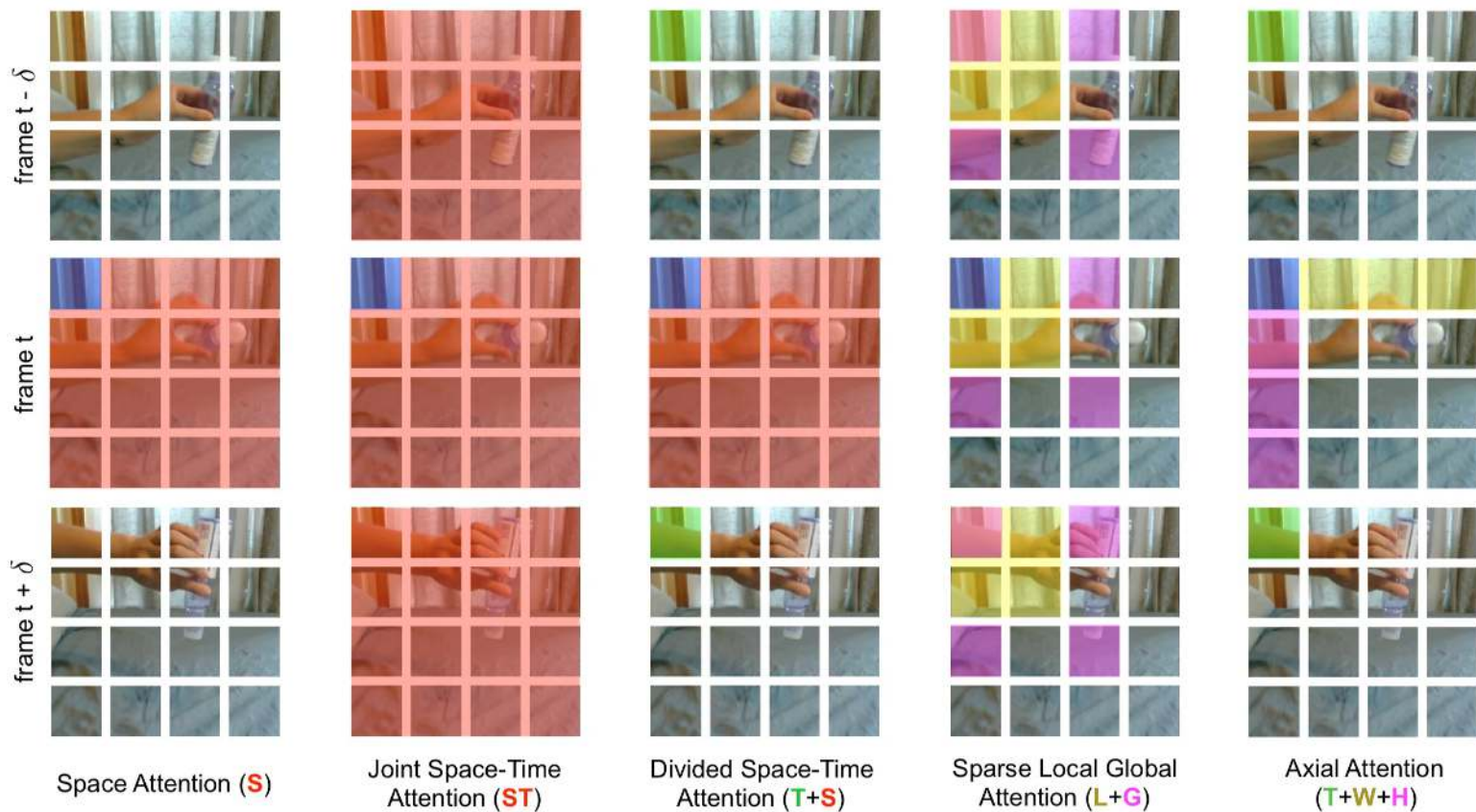
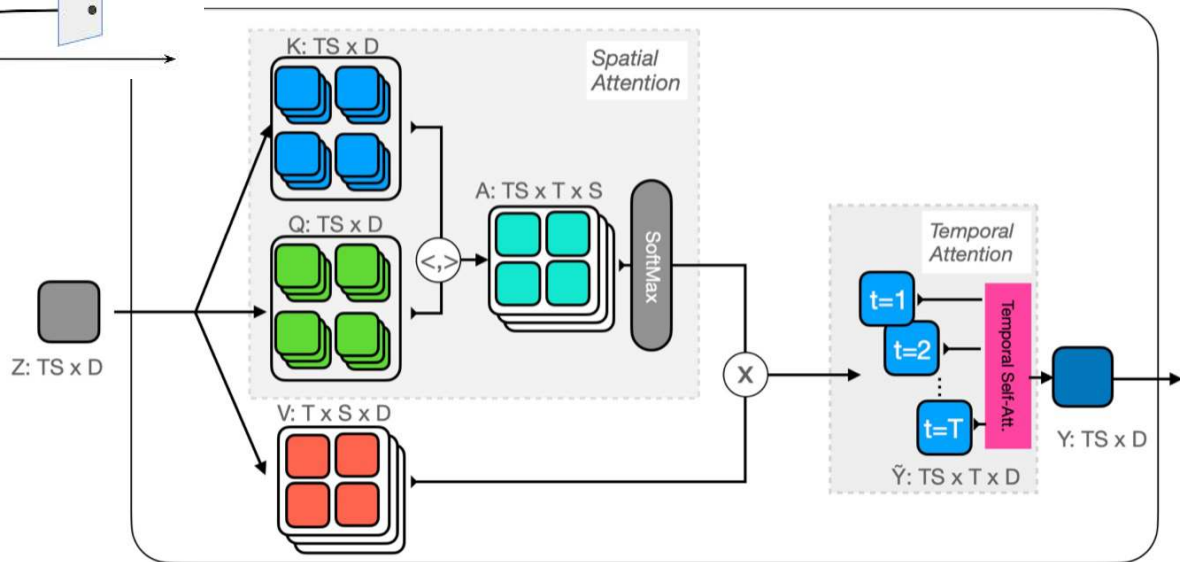
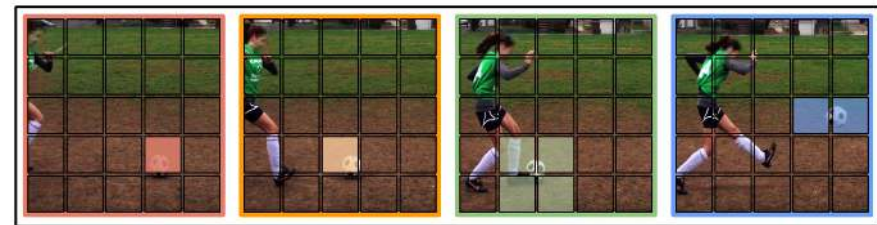


Figure from: Bertasius et al. Is Space-Time Attention All You Need for Video Understanding? ArXiv 2021

MotionFormer



VideoMAE

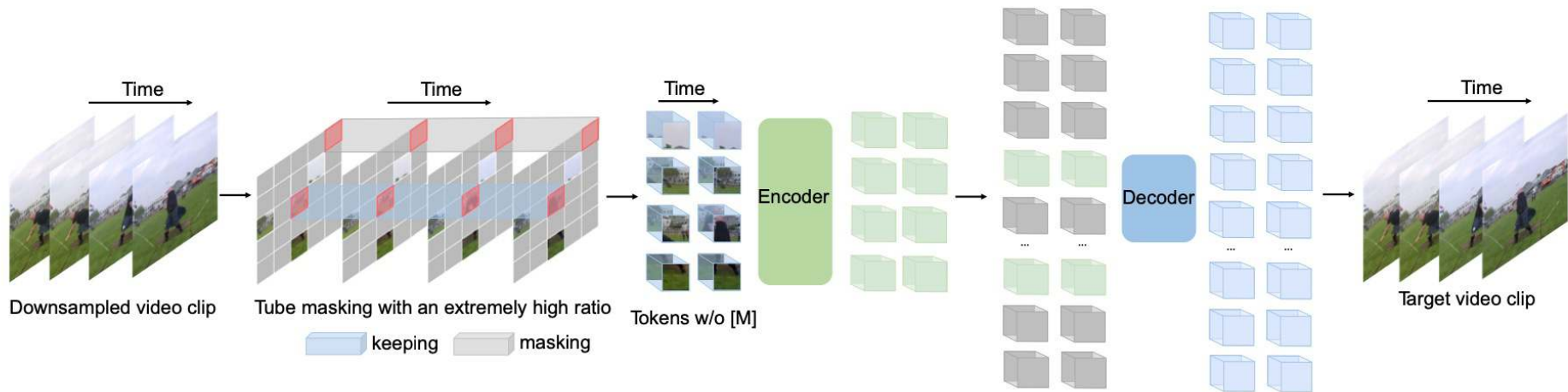


Figure from: Tong et al (2022). VideoMAE: Masked Autoencoders are Data-Efficient Learners for Self-Supervised Video Pre-Training. NeurIPS 2022

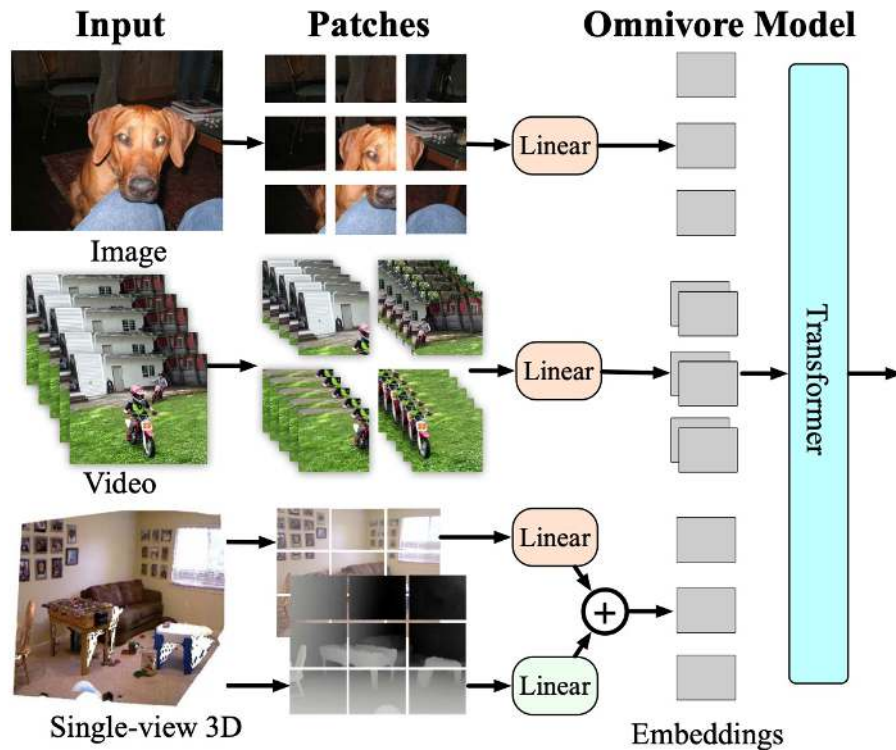


Figure 2. Multiple visual modalities in the OMNIVORE model.

Figure from: Gridhar et al (2022). OMNIVORE: A Single Model for Many Visual Modalities. CVPR

ImageBind



Web Image-Text



Sheep basking in the sun

Depth Sensor Data



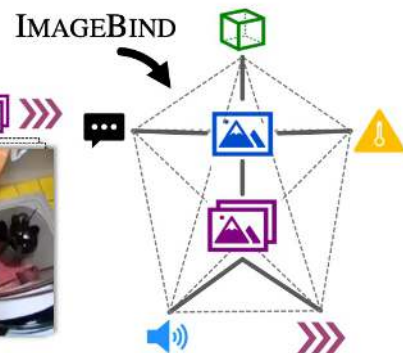
Web Videos



Thermal Data

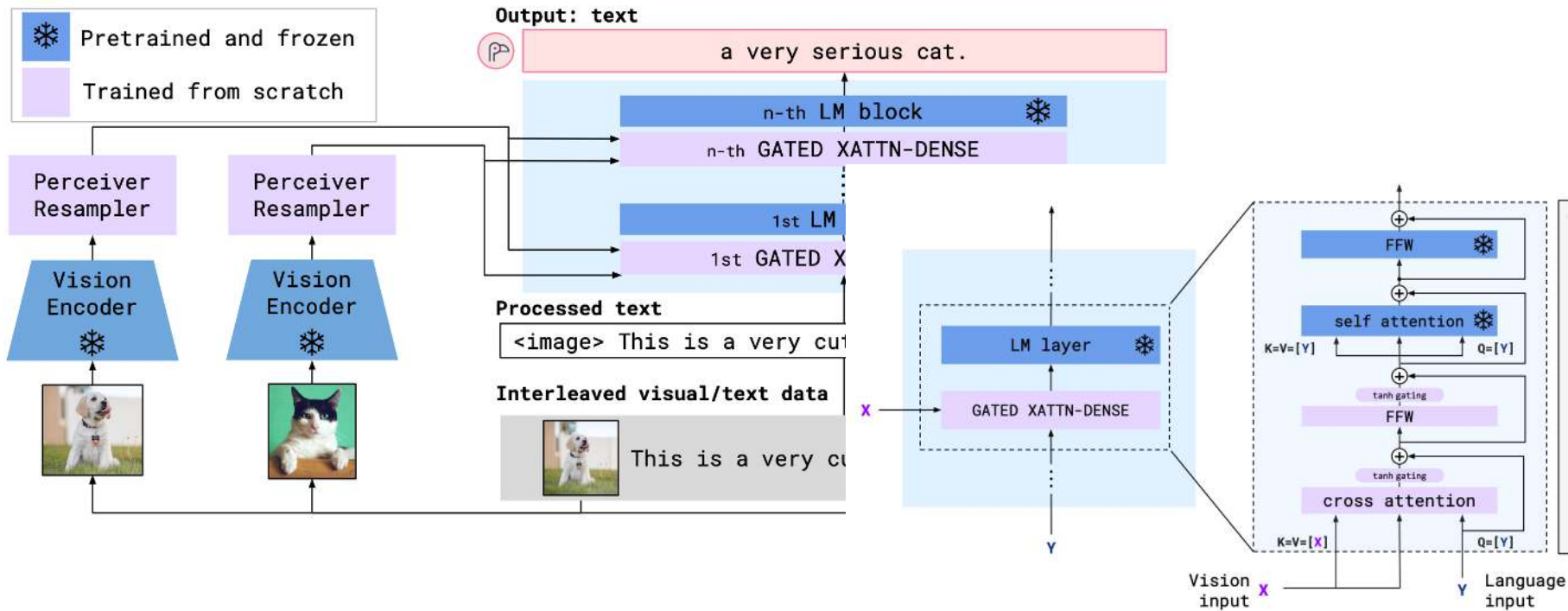


Egocentric Videos



$$L_{\mathcal{I}, \mathcal{M}} = -\log \frac{\exp(\mathbf{q}_i^\top \mathbf{k}_i / \tau)}{\exp(\mathbf{q}_i^\top \mathbf{k}_i / \tau) + \sum_{j \neq i} \exp(\mathbf{q}_i^\top \mathbf{k}_j / \tau)}$$

Flamingo



InternVideo

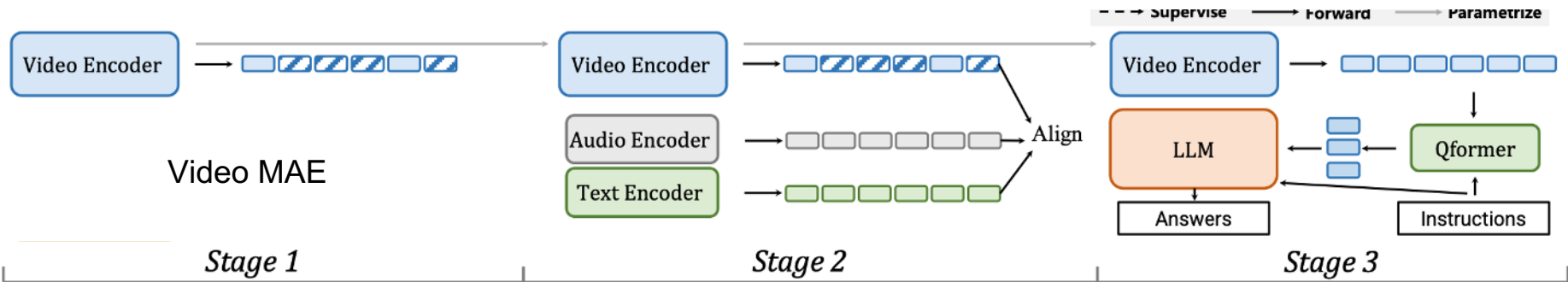


Figure 2: Framework of **InternVideo2**. It consists of three consecutive training phases: unmasked video token reconstruction, multimodal contrastive learning, and next token prediction. In stage 1, the video encoder is trained from scratch, while in stages 2 and 3, it is initialized from the version used in the previous stage.

InternVideo

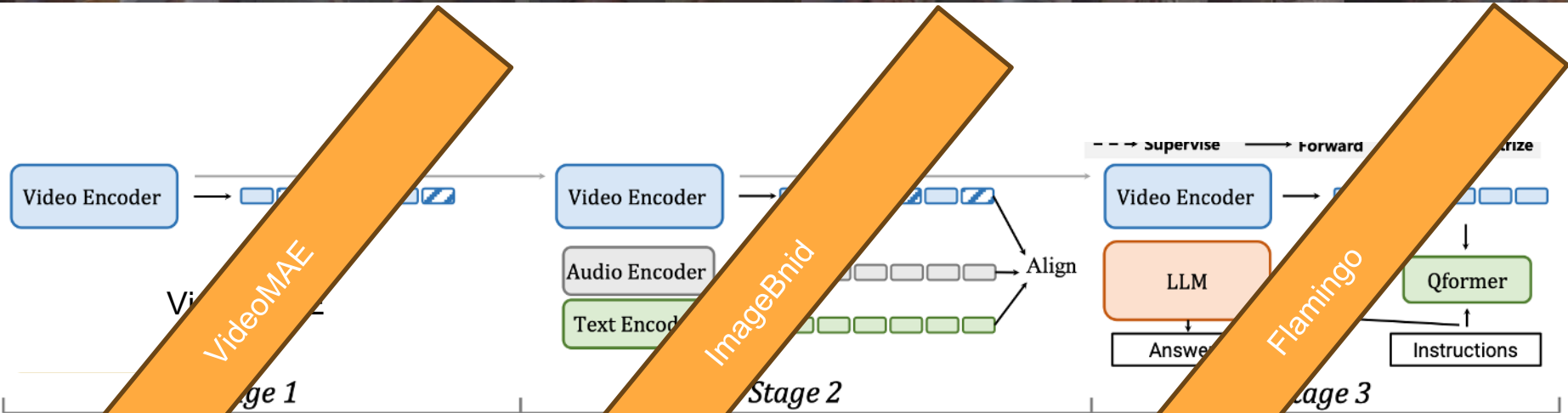


Figure 1: Framework of InternVideo2. It consists of three consecutive training stages. In stage 1, unmasked video token reconstruction, multimodal contrastive learning, and next token prediction. In stage 2, the video encoder is trained from scratch while in stages 2 and 3, it is initialized from the version used in the previous stage.



We are still lacking the right
models for video
understanding

From Clip to Video

Video



Wash carrot

From Clip to Video

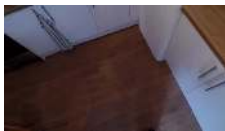


Wash carrot



From Clip to Video

Video



Visual labels





Most models work only within
the clip... [ignoring the
context]

Long-Term Feature Bank

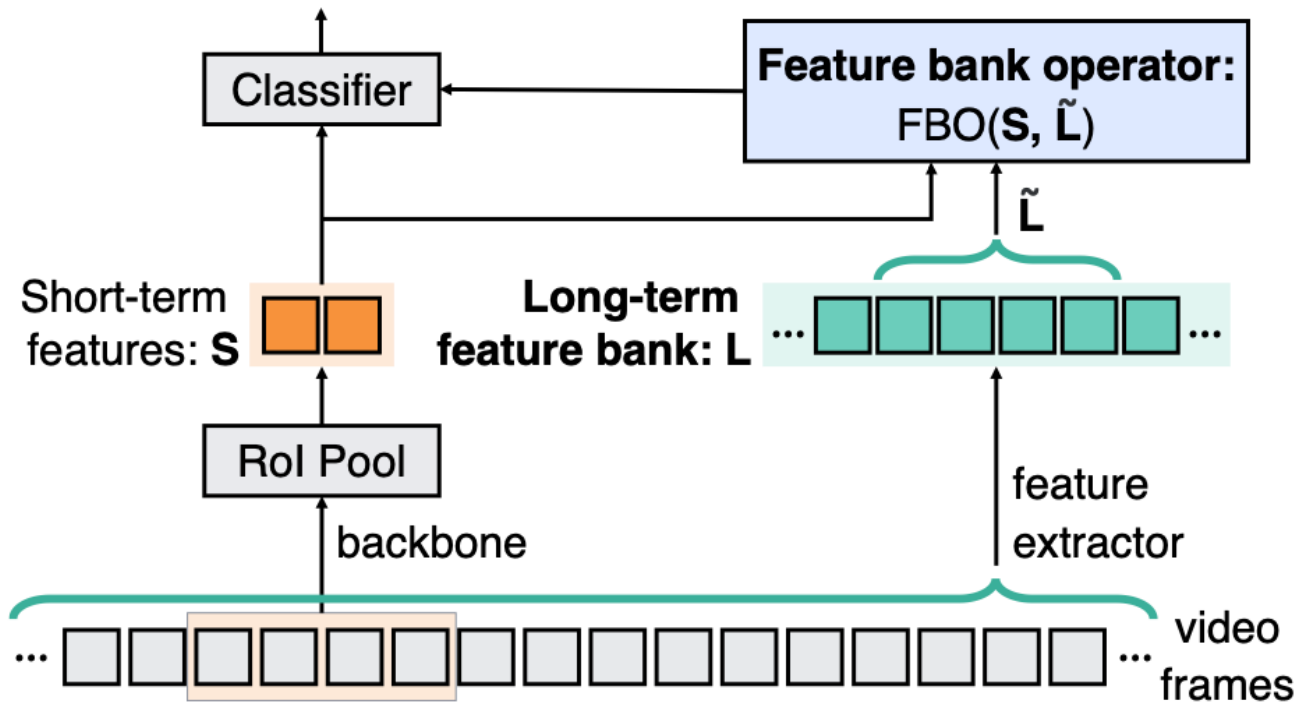
Target frame



←————— Input clip (4 seconds) —————→

Figure 1. What are these people doing? Current 3D CNN video models operate on short clips spanning only ~ 4 seconds. Without observing longer-term context, recognition is difficult. (Video from the AVA dataset [14]; see next page for the answer.)

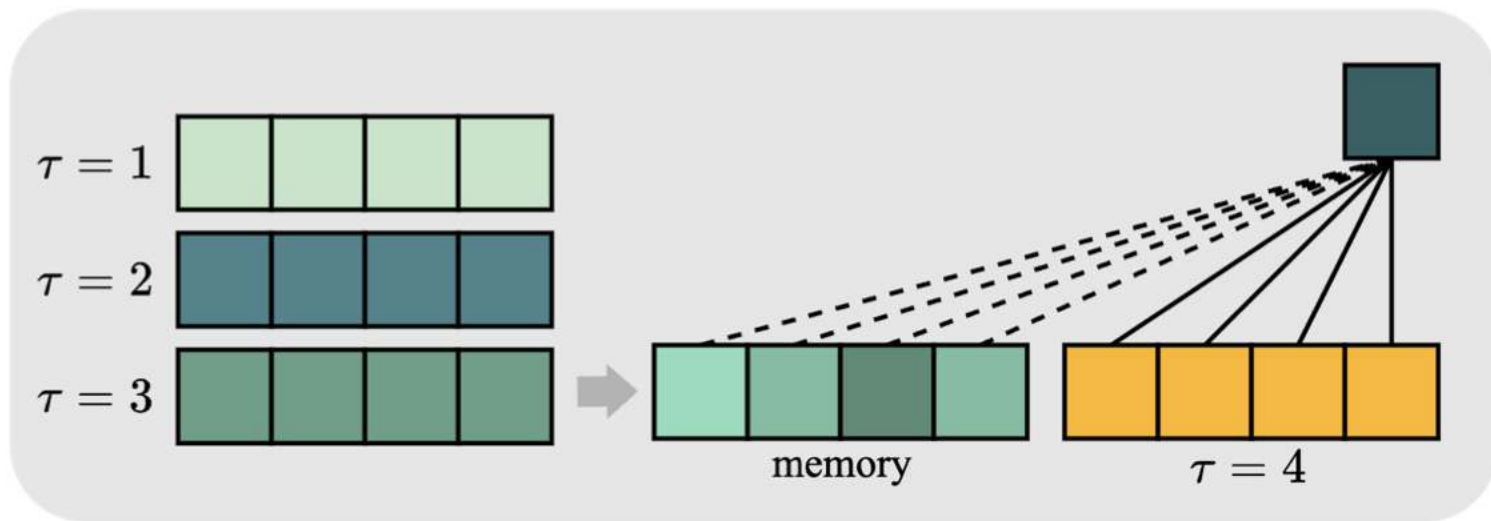
Long-Term Feature Bank



(b) **3D CNN with a Long-Term Feature Bank (Ours)**

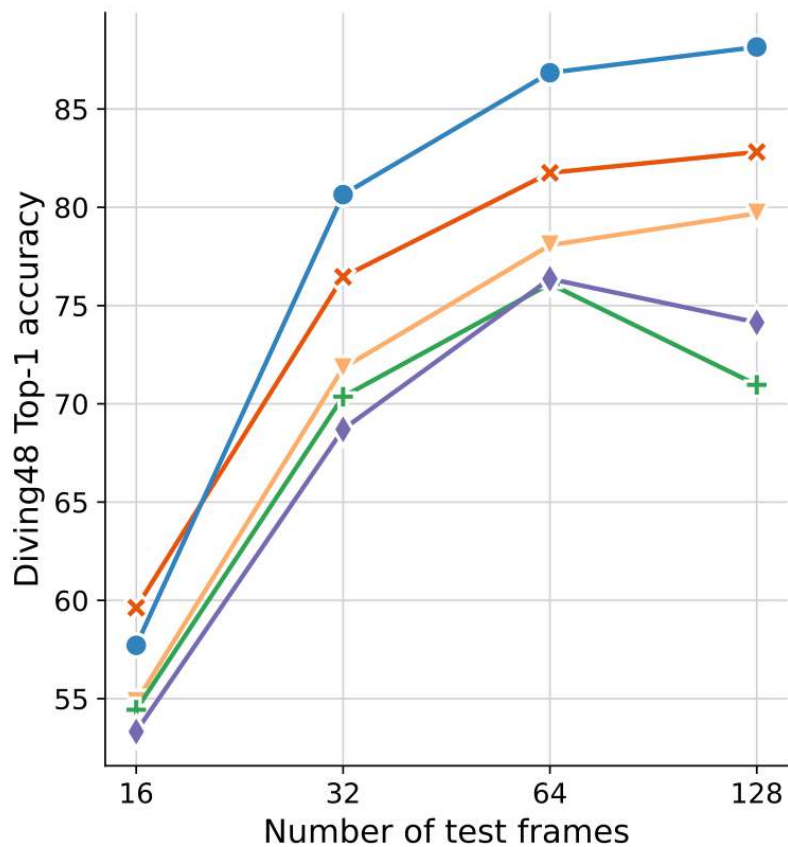
Memory Consolidation

— self-attention
- - - cross-attention



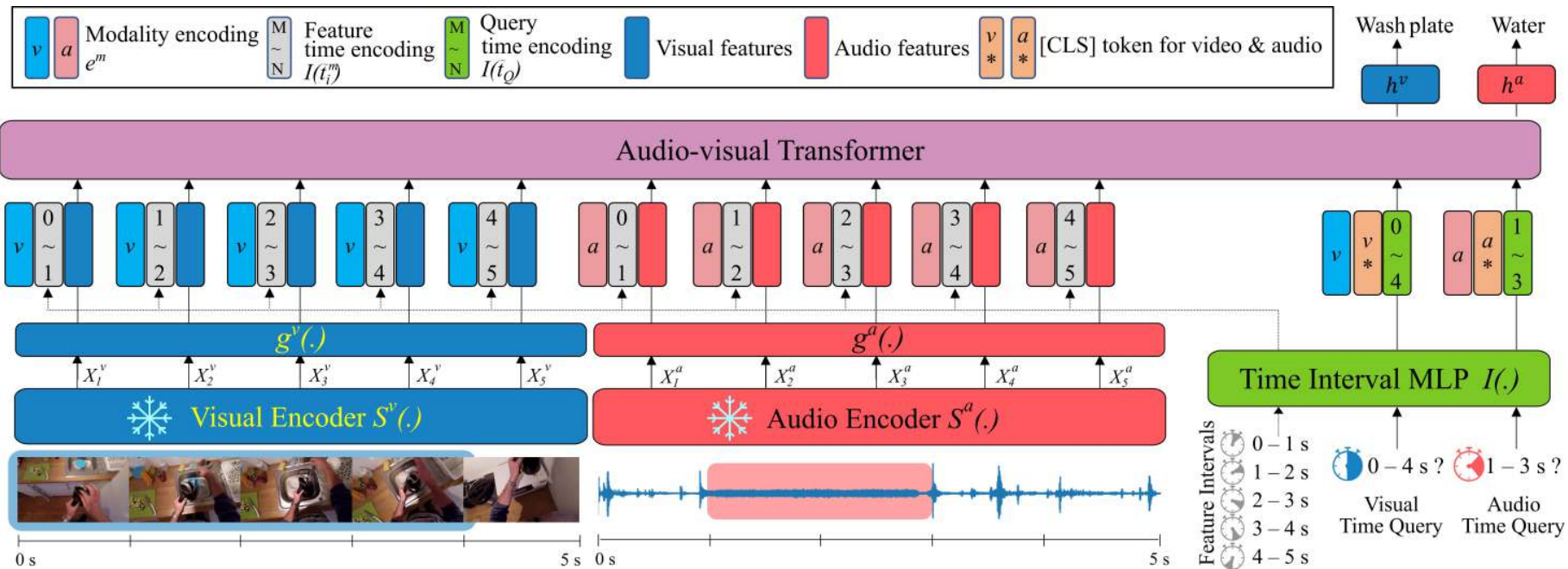
Memory-Consolidated ViT

Memory Consolidation



TIM: A Time-Interval Audio-Visual Machine

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





*Do you think video is a
unique modality?!!*



On Tasks...

Two types of video understanding tasks

Video Understanding Tasks

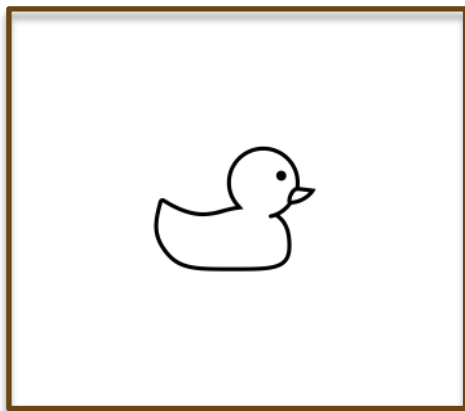
Analogous to Image-based Tasks

Temporal tasks (not relevant for images)

Analogous Tasks

Image

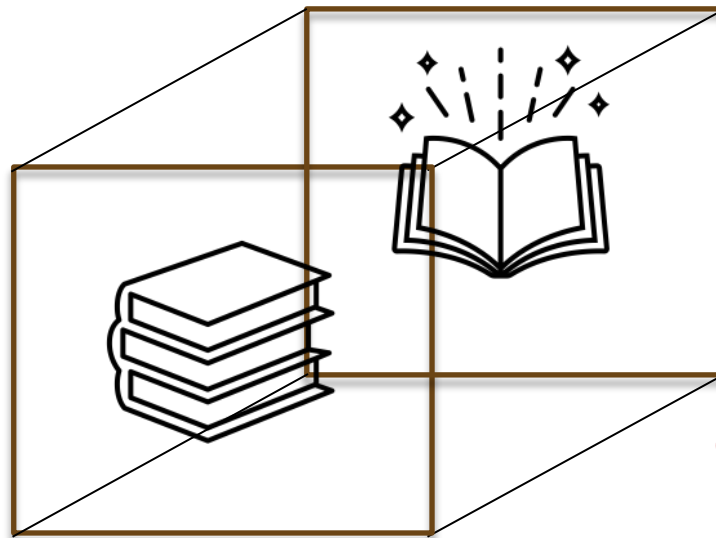
- Object Recognition



Duck

Video

- Action Recognition

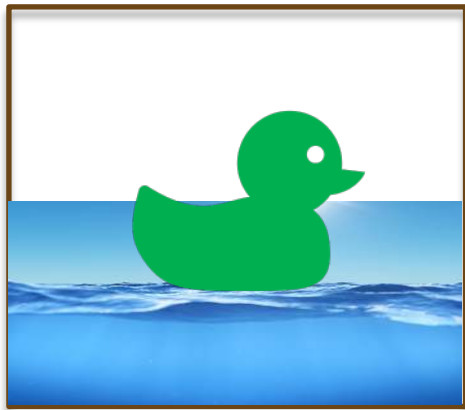


Open
Book

Analogous Tasks

Image

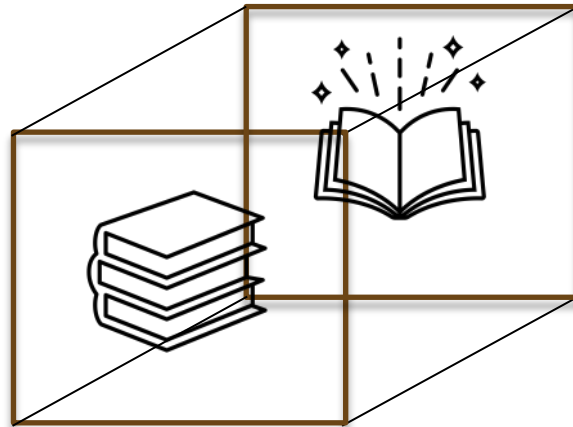
- Image Captioning



A green duck swimming
In clear water

Video

- Video Captioning

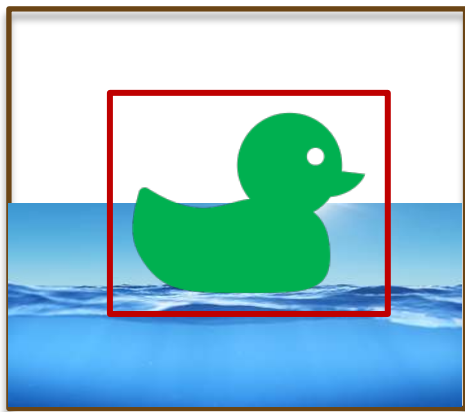


A book picked from top of the pile
and opened to a page in the middle

Analogous Tasks

Image

- Object Detection



Duck

Video

- Action Detection



Open Book

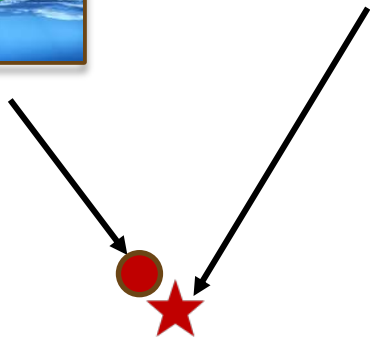
Analogous Tasks

Image

- Image Retrieval

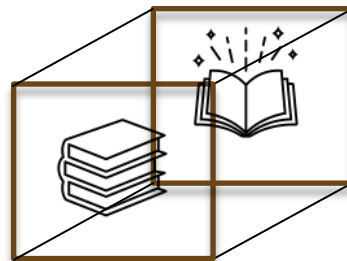


Duck

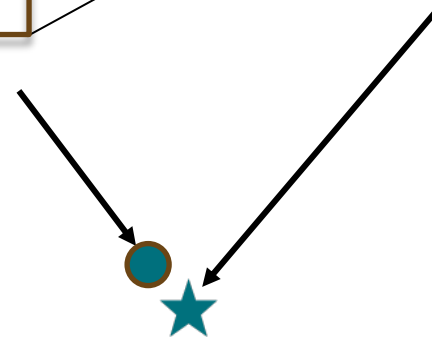


Video

- Video Retrieval



Open Book



What is a Cross-Modal Video Retrieval?

Video-to-Text Retrieval Task

Q



Ranked Text – Gallery (or Retrieval Set)



put garlic down

Text-to-Video Retrieval Task

Q put garlic down

Ranked Video – Gallery (or Retrieval Set)



In this work we focus on
Fine-Grained Action Retrieval

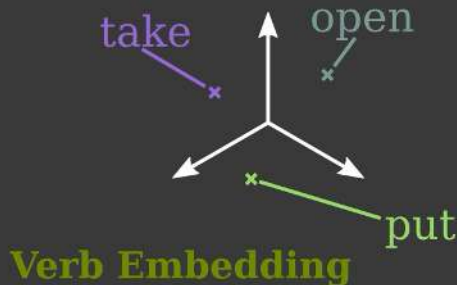
I put meat on a
ball of dough



Fine-Grained Action Retrieval

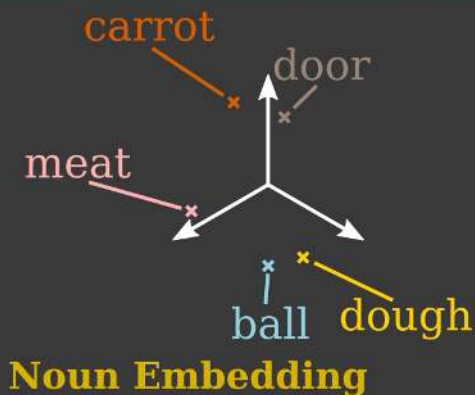
with: Michael Wray
Gabriela Csurka
Diane Larlus

We embed the video
and representations



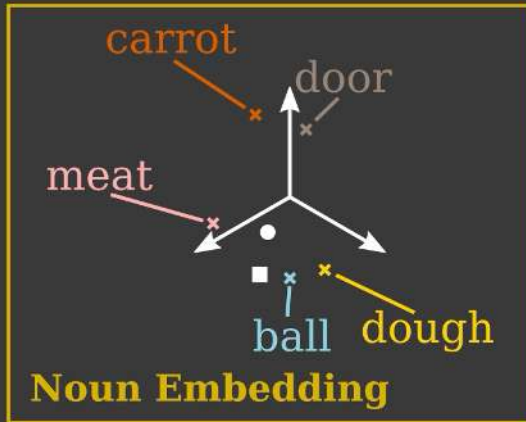
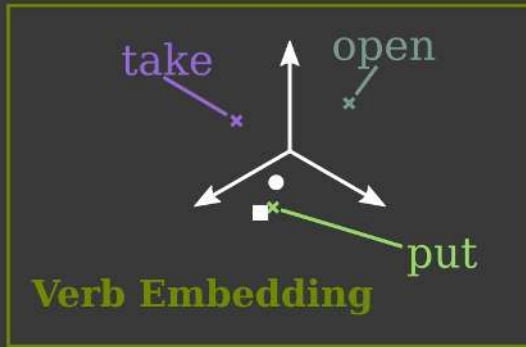
[put]

[meat, ball, dough]

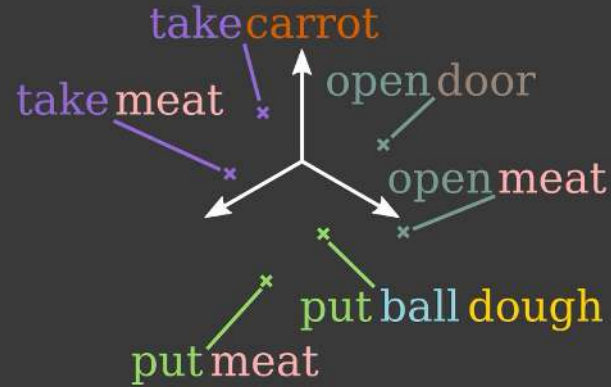


Fine-Grained Action Retrieval

with: Michael Wray
Gabriela Csurka
Diane Larlus

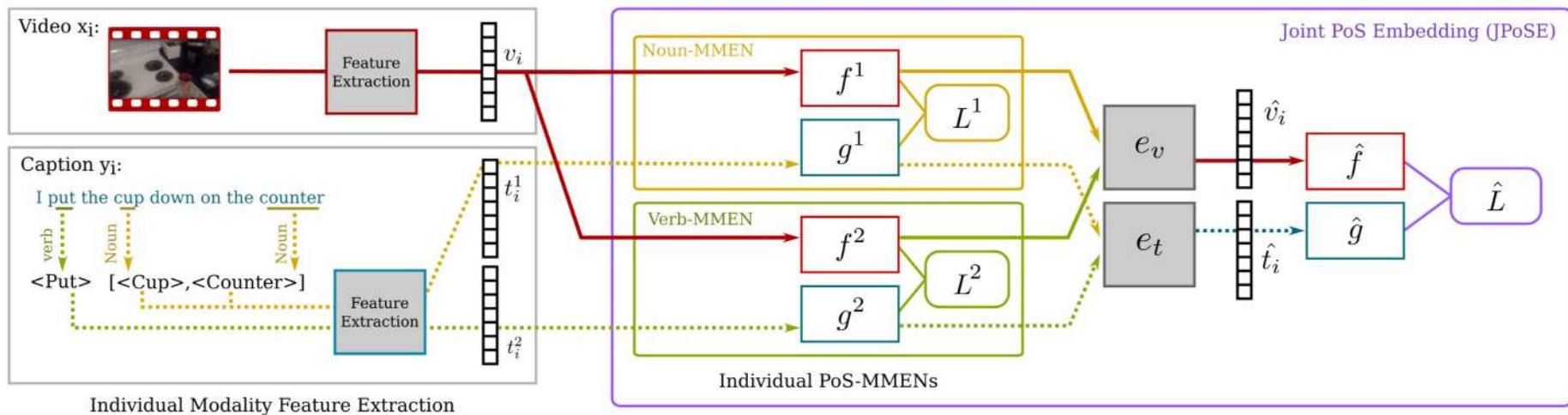


Finally, we combine the outputs and embed these into an action space



Fine-Grained Action Retrieval

with: Michael Wray
Gabriela Csurka
Diane Larlus



Individual Modality Feature Extraction

Individual PoS-MMENs

Joint PoS Embedding (JPoSE)

Maximum activation examples for a neuron in a noun PoS Embedding (Cutting Board) - Figure 4



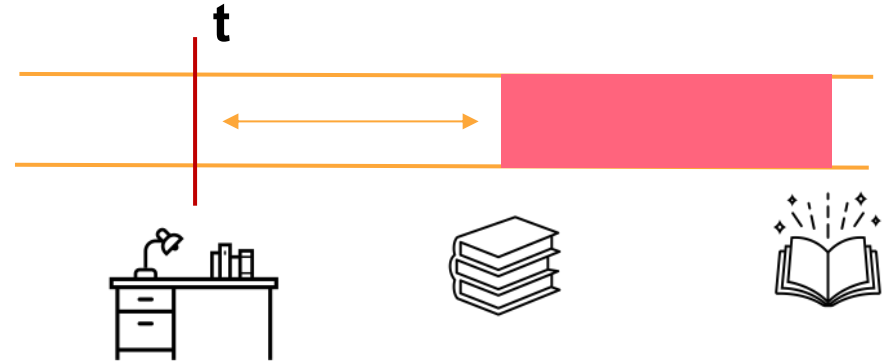
Non-Analogous Tasks

Image



Video

- Action Anticipation
What will happen after 1 second?



Open Book

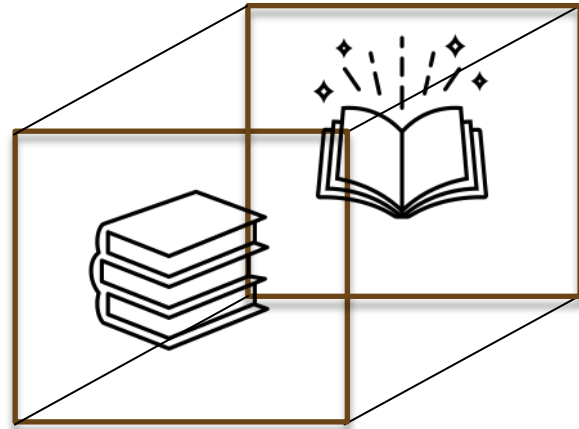
Non-Analogous Tasks

Image



Video

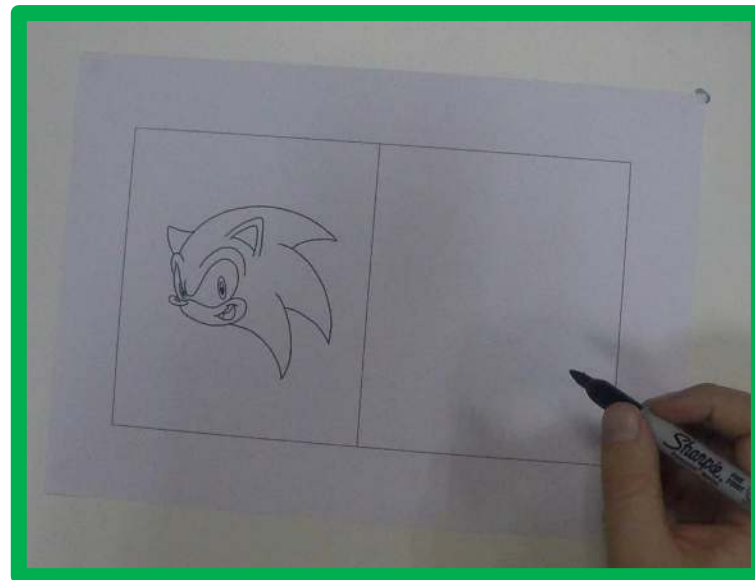
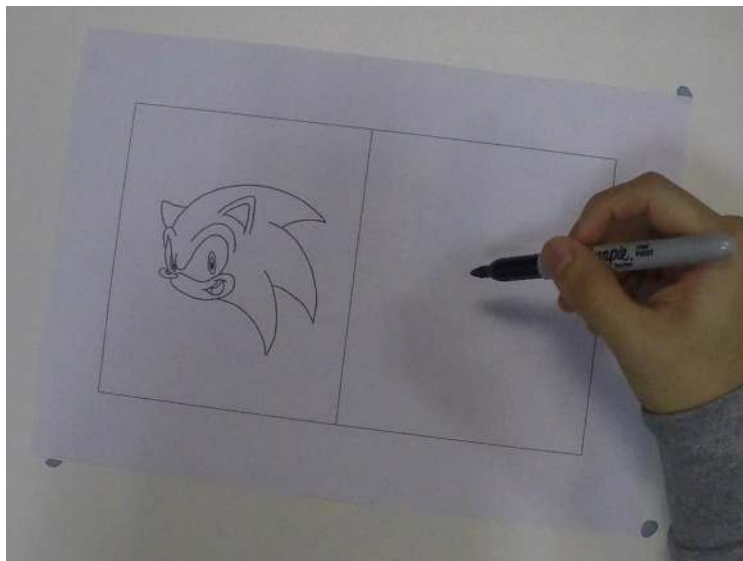
- Skill Understanding
How did you open the book?



Skill determination in video

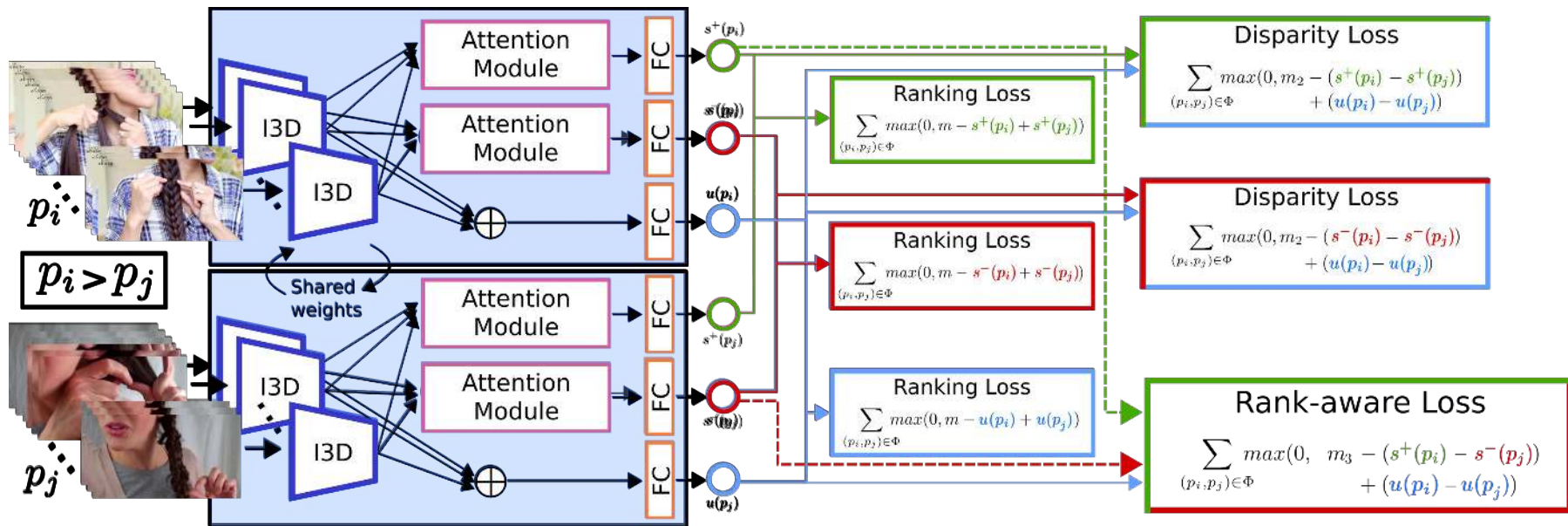
with: Hazel Doughty
Walterio Mayol-Cuevas

Pairwise annotations of videos, indicating higher skill or no skill preference



Skill determination in video

with: Hazel Doughty
Walterio Mayol-Cuevas



Low-skill Attention Module

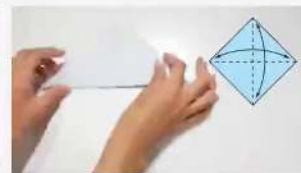
Surgery



Apply Eyeliner

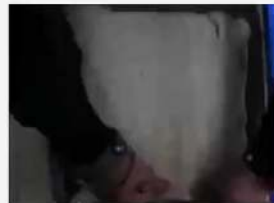


Origami

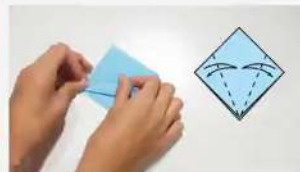
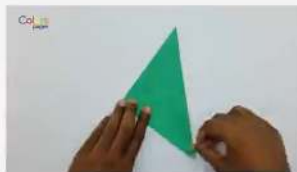


High-skill Attention Module

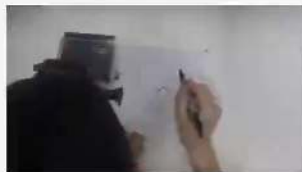
Dough
Rolling



Origami



Drawing



Analogous Tasks

Image

- Object Counting

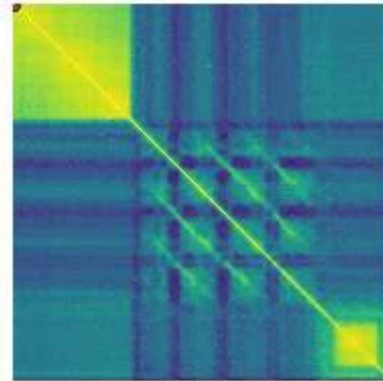
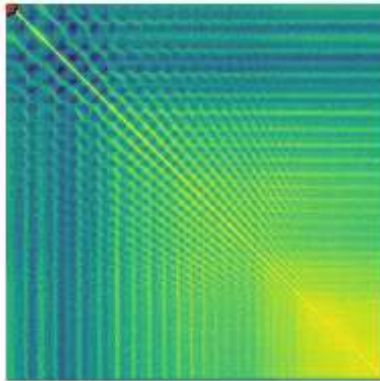
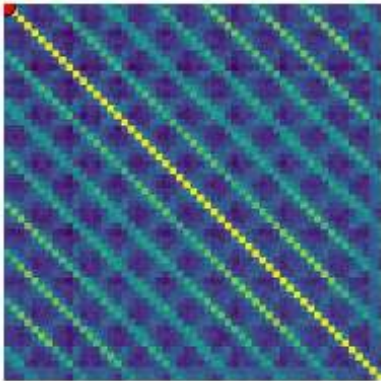
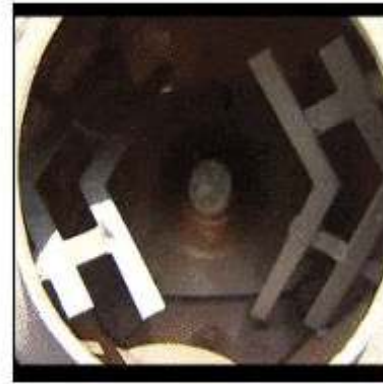


Video

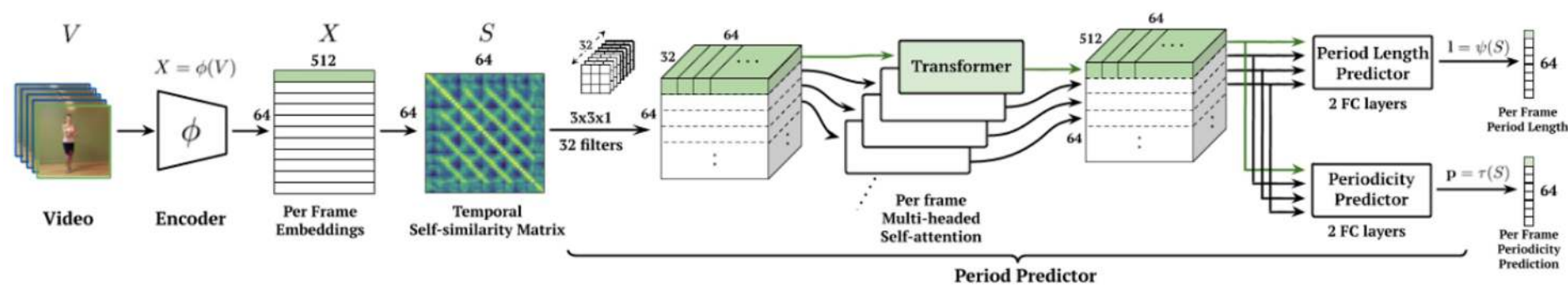
- Action Counting



Countix



Countix



Every Shot Counts

with: Saptarshi Sinha
Alexandros Stergiou

RepCount



Countix

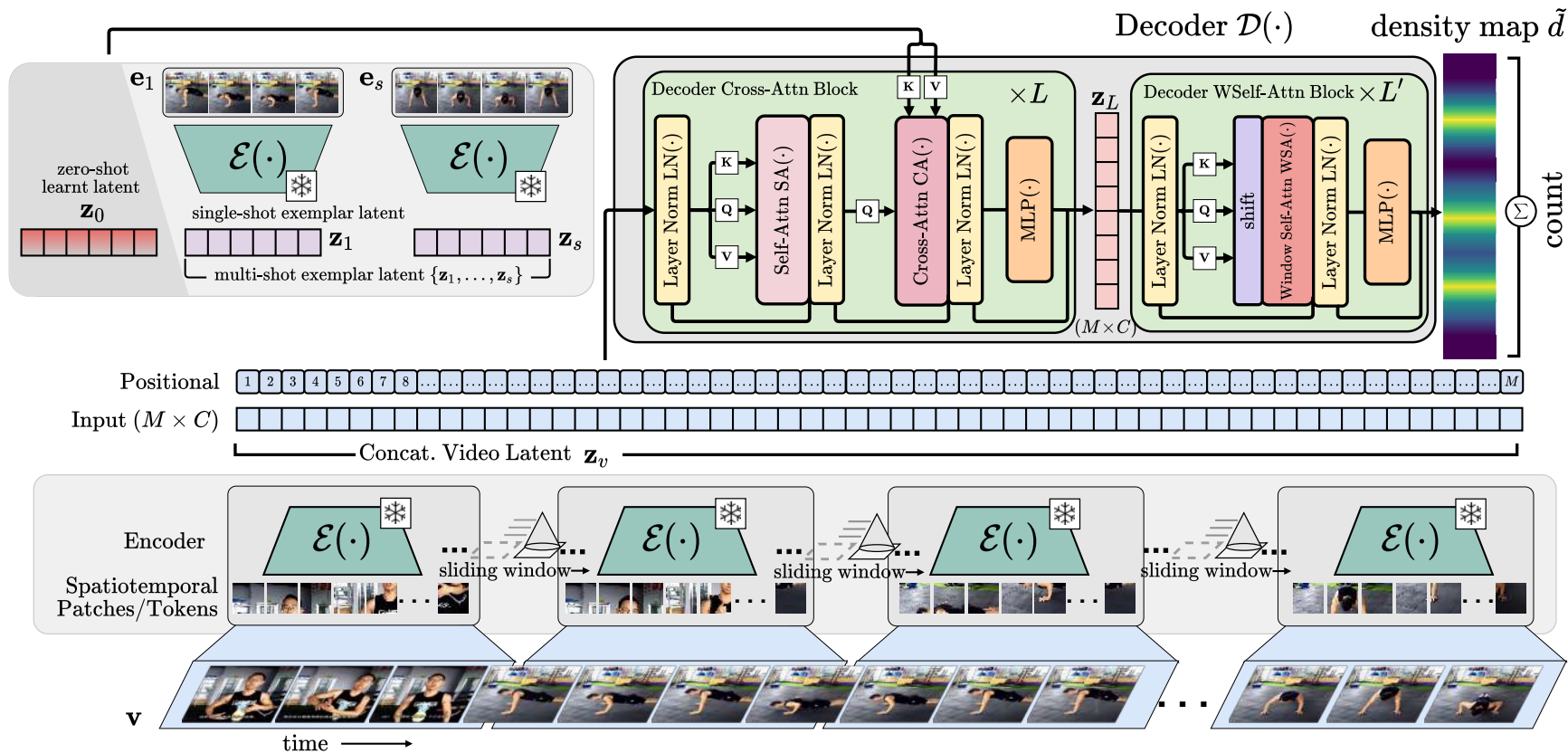


RepCount



Every Shot Counts

with: Saptarshi Sinha
Alexandros Stergiou



Every Shot Counts

with: Saptarshi Sinha
Alexandros Stergiou

(a) RepCount

Method	Encoder	RMSE↓	MAE↓	OBZ↑	OBO↑
RepNet [15]	R2D50	-	0.995	-	0.013
TransRAC [18]	VSwinT	9.130*	0.443	0.085*	0.291
MFL [27]†	VSwinT	-	0.384	-	0.386
ESCounts	VSwinT	6.905	0.298	0.183	0.403
ESCounts	VMAE	4.455	0.213	0.245	0.563

(c) UCFRep

Method	Encoder	RMSE↓	MAE↓	OBZ↑	OBO↑
Levy & Wolf [25]	RX3D101	-	0.286	-	0.680
RepNet [15]	R2D50	-	0.998	-	0.009
Context (F) [62]	RX3D101	5.761*	0.653*	0.143*	0.372*
TransRAC [18]	VSwinT	-	0.640	-	0.324
MFL [27]†	RX3D101	-	0.388	-	0.510
ESCounts	RX3D101	2.004	0.247	0.343	0.731
ESCounts	VMAE	1.972	0.216	0.381	0.704

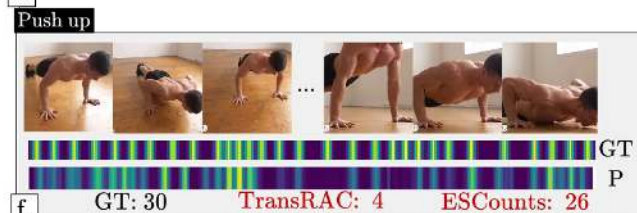
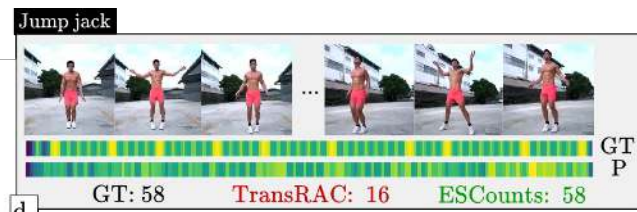
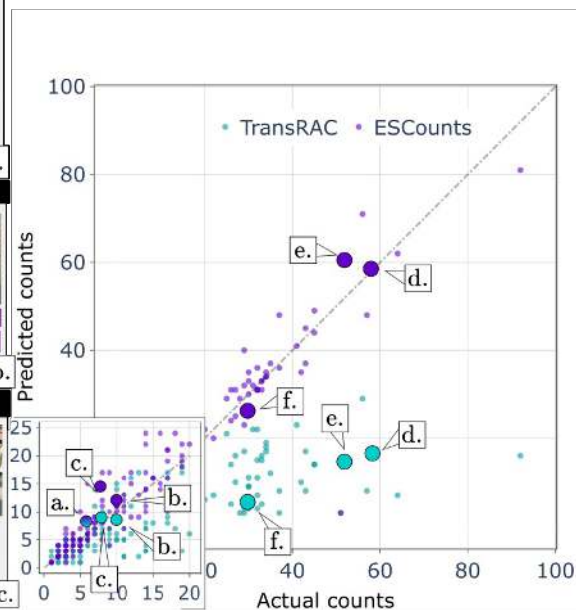
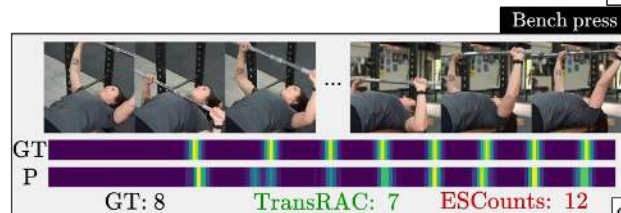
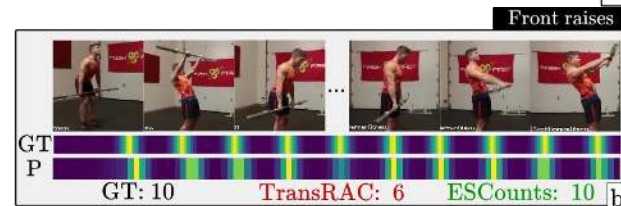
(b) Countix

Method	Encoder	RMSE↓	MAE↓	OBZ↑	OBO↑
RepNet [15]	R2D50	-	0.364	-	0.697
Sight & Sound [64]†	R(2+1)D18	-	0.307	-	0.511
ESCounts	R(2+1)D18	3.536	0.293	0.286	0.701
ESCounts	VMAE	3.029	0.276	0.319	0.673

Every Shot Counts

with: Saptarshi Sinha
Alexandros Stergiou

RepCount



Image

- Text-to-image Generation



Stable Diffusion

Video

- Text-to-Video Generation



SORA

Text-to-Video Generation



Text-to-Video Generation



Prompt: A grandmother with neatly combed grey hair stands behind a colorful birthday cake with numerous candles at a wood dining room table, expression is one of pure joy and happiness, with a happy glow in her eye. She leans forward and blows out...



Generative Video
approaches do not yet
understand physics, actions
or action consequences...

GenHowTo: Learning to Generate Actions and State Transformations from Instructional Videos



Tomáš Souček



Dima Damen



Michael Wray



Ivan Laptev



Josef Šivic



- Hands transform objects....

Input



peeled ♠ on chopping board



♠ in a blender



♠ smoothie in a blender



♠ = avocado

Input



GenHowTo



EF-DDPM



InstructPix2Pix



Prompt: a frosted cake with strawberries around the top



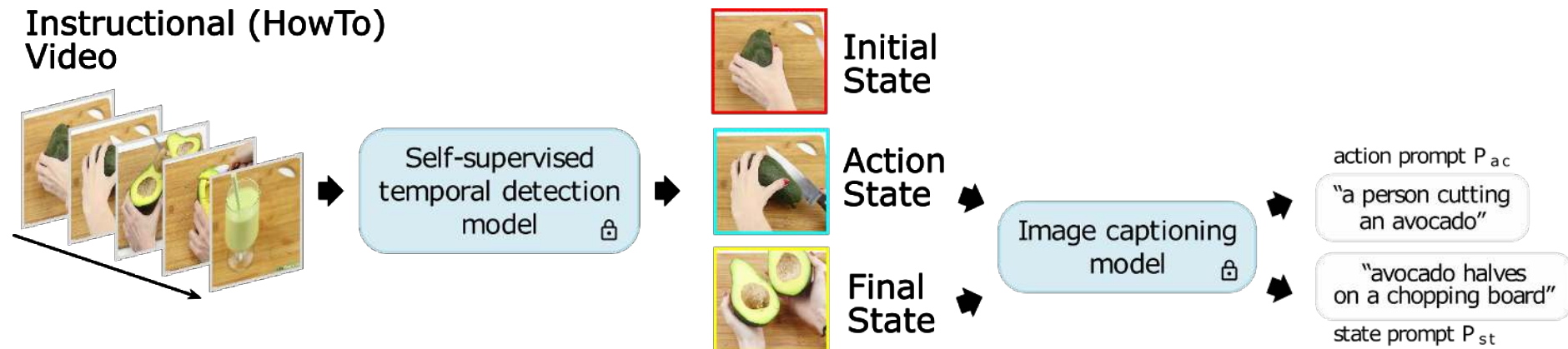
Prompt: a person kneading dough on a cutting board



Prompt: a person cutting a fish on a cutting board

- Two contributions.... Dataset & Method

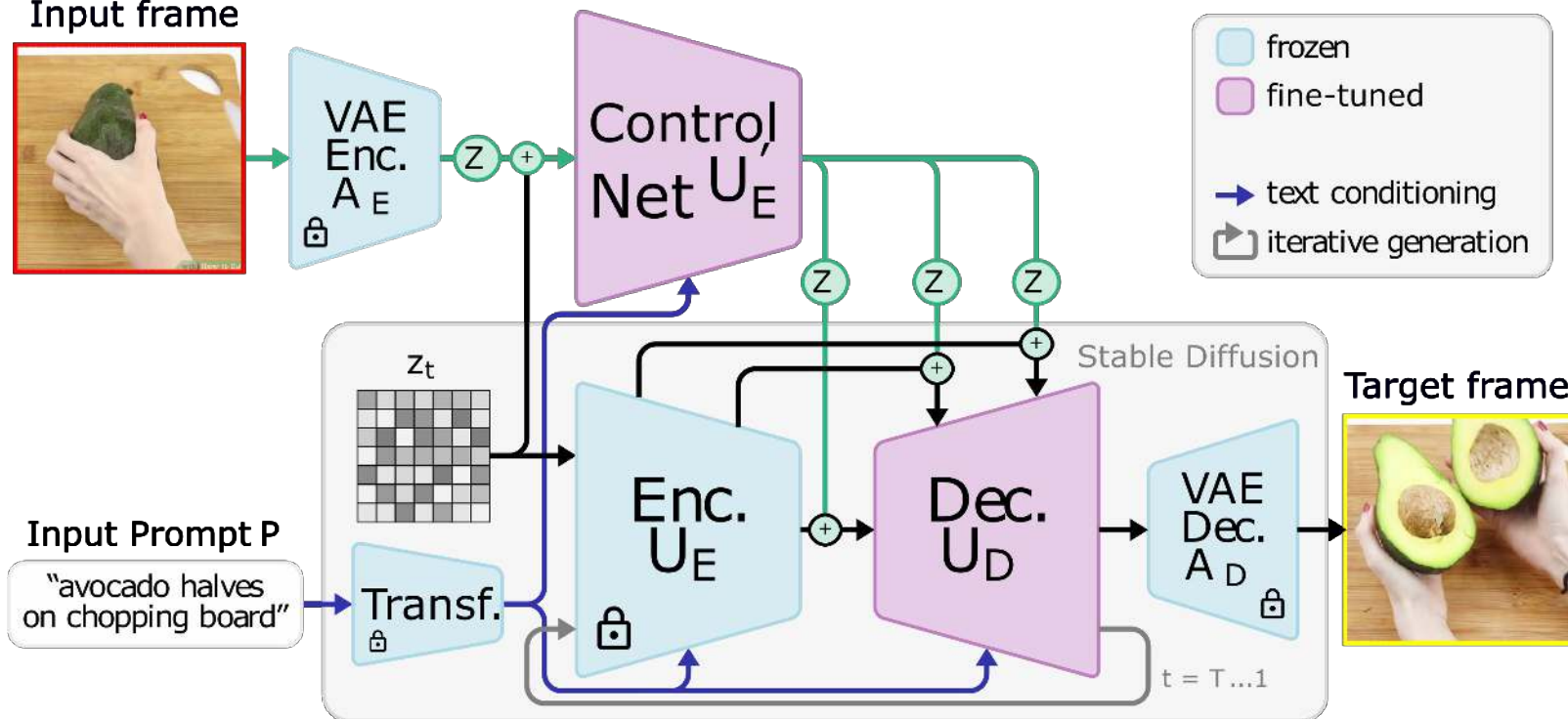
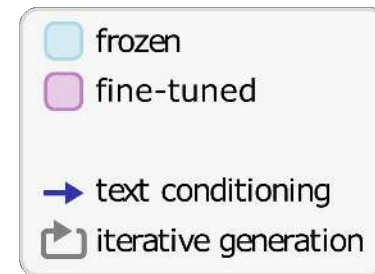
- Two contributions.... **Dataset** & Method



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 & Method

Input frame



GenHowTo...

with: Tomas Soucek
Ivan Laptev
Michael Wray
Josef Sivic

Input

less noise



more noise



- Qualitative Evaluation...

- Initial vs Final State
- Binary Classifier

Method	Acc _{ac} ↑	Acc _{st} ↑
<i>test set categories unseen during training</i>		
(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
<i>test set categories seen during training</i>		
(f) Edit Friendly DDPM [†]	0.69	0.80
(g) GenHowTo [†]	0.77	0.88
(h) <i>Real images</i>	0.96	0.97

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

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

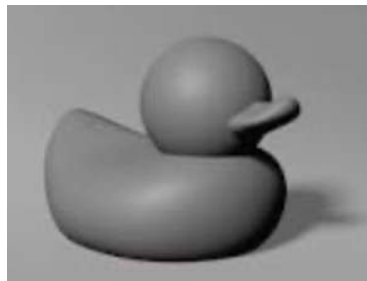
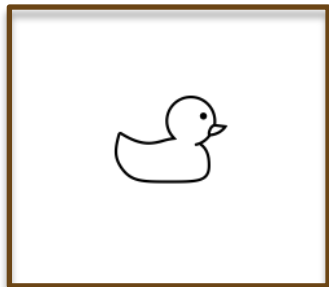


It is more important to
understand consequences of
actions that to generate
smooth motions

Analogous Tasks

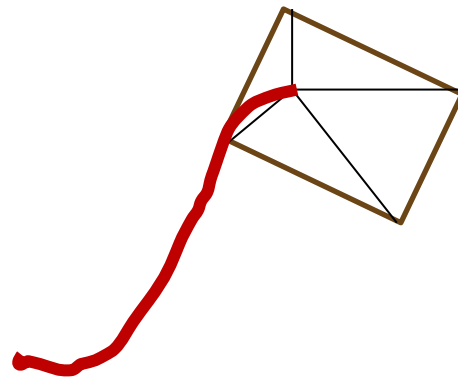
Image

- Image – to – 3D



Video

- Video – to - 3D

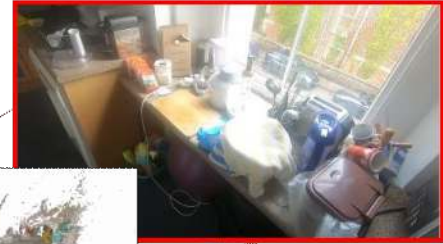




Almost all video-to-3D
focuses on a static scene

EPIC Fields

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





EPIC-KITCHENS

Spatial Cognition from Egocentric Video: Out of Sight, Not Out of Mind

Chiara Plizzari Shubham Goel Toby Perrett Jacob Chalk
Angjoo Kanazawa Dima Damen

<http://dimadamen.github.io/OSNOM>





3D Scene Mesh →

← Egocentric Image

↑ 3D Ego view w/ in-view objects

↑ Ego Camera in 3D

All active/moved objects in this video are represented by neon balls. Their initial positions are shown at the start of the video



← Egocentric Image



↑ 3D Ego view w/
in-view
objects



3D Scene
Mesh →

↖ Ego
Camera
in 3D

All active/moved objects in this video are represented by neon balls. Their initial positions are shown at the start of the video

Non-Analogous Tasks

Image



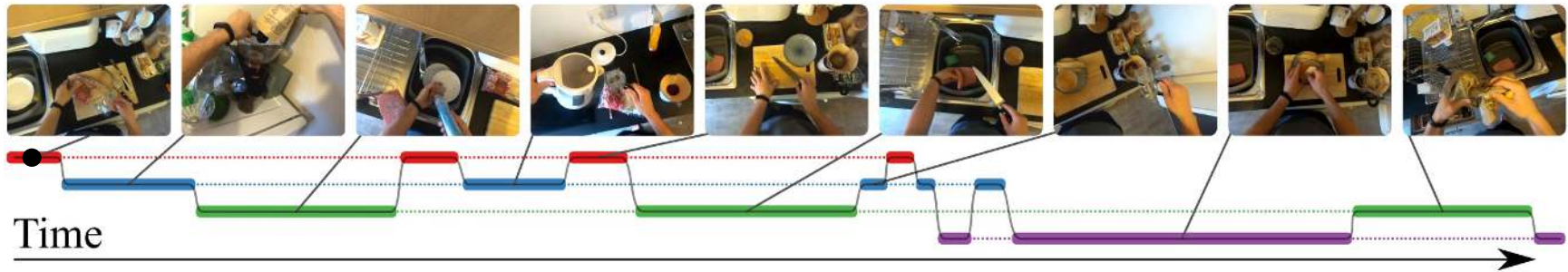
Video

- Understanding Goals in Long Videos



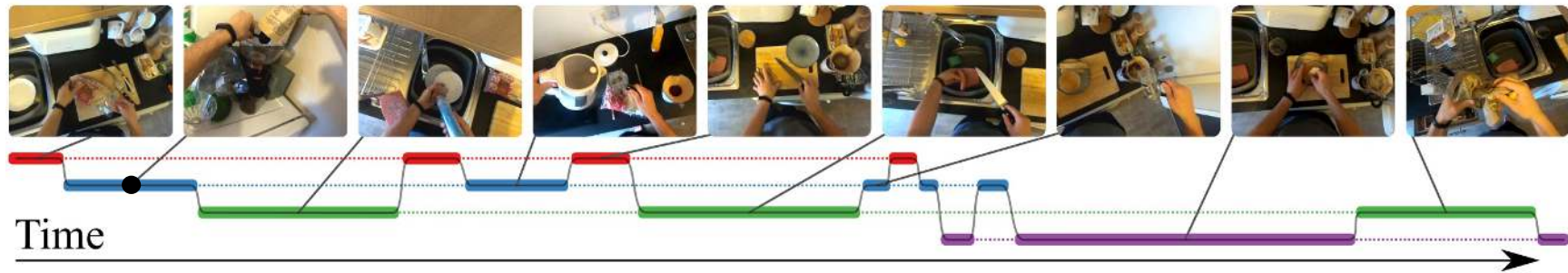
Goals...

with: Will Price
Carl Vondrick



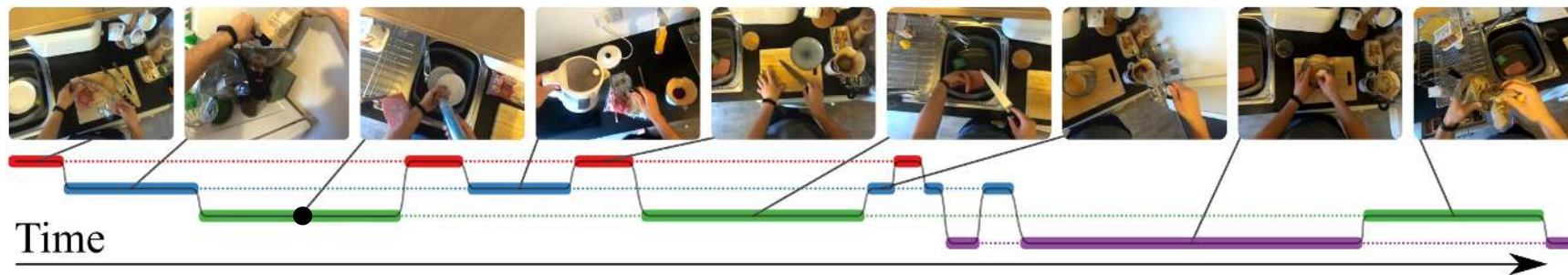
UnweaveNet

with: Will Price
Carl Vondrick



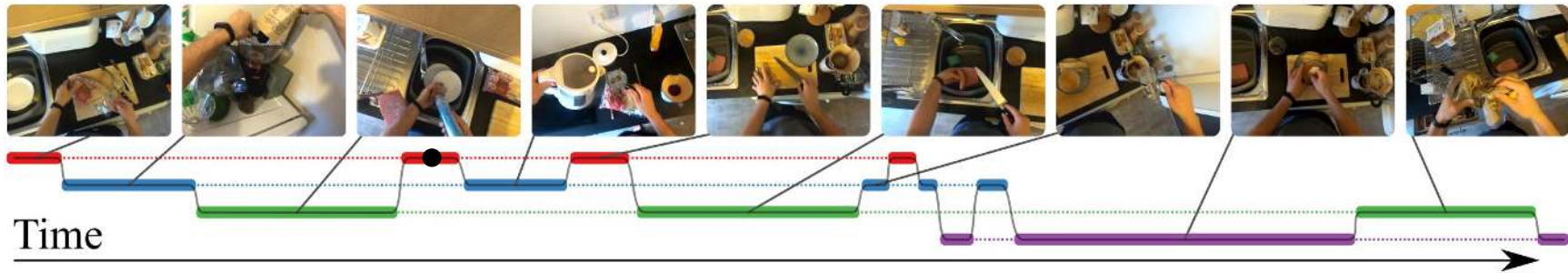
UnweaveNet

with: Will Price
Carl Vondrick



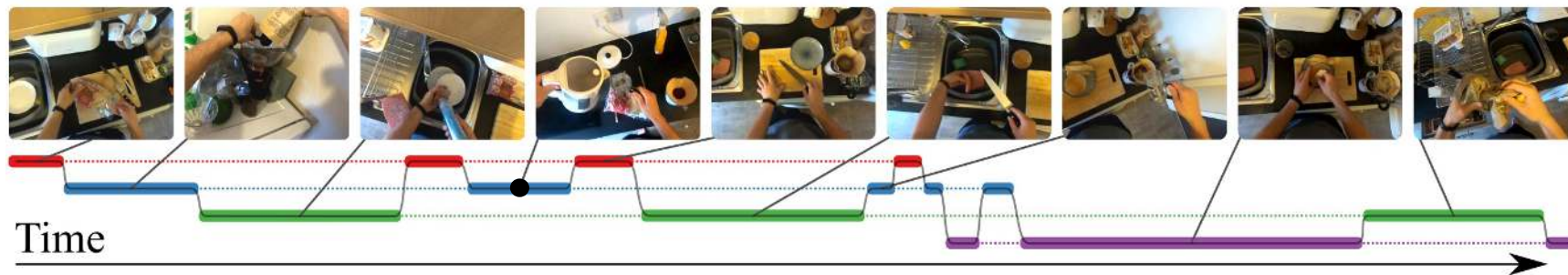
UnweaveNet

with: Will Price
Carl Vondrick



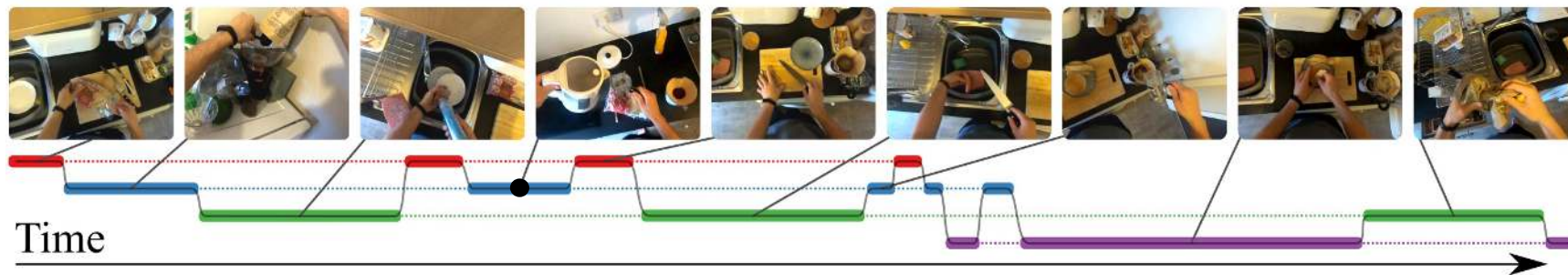
UnweaveNet

with: Will Price
Carl Vondrick



UnweaveNet

with: Will Price
Carl Vondrick



UnweaveNet

with: Will Price
Carl Vondrick

A blue rounded rectangle containing the text "UnweaveNet".

UnweaveNet

UnweaveNet

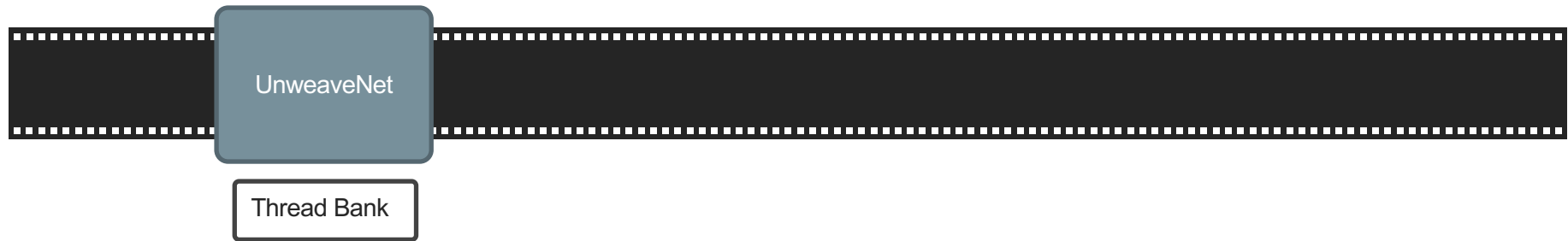
with: Will Price
Carl Vondrick

UnweaveNet

Thread Bank

UnweaveNet

with: Will Price
Carl Vondrick



UnweaveNet

with: Will Price
Carl Vondrick



UnweaveNet

with: Will Price
Carl Vondrick



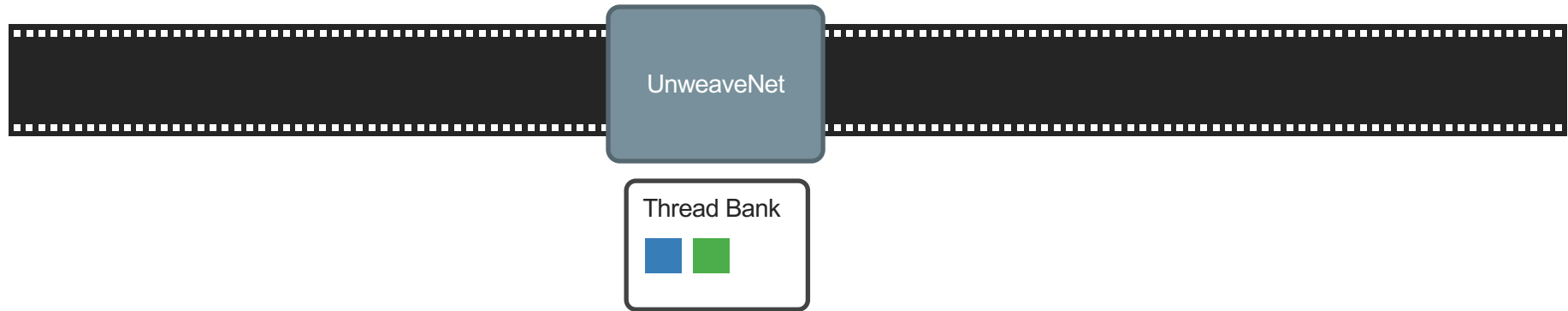
UnweaveNet

Thread Bank



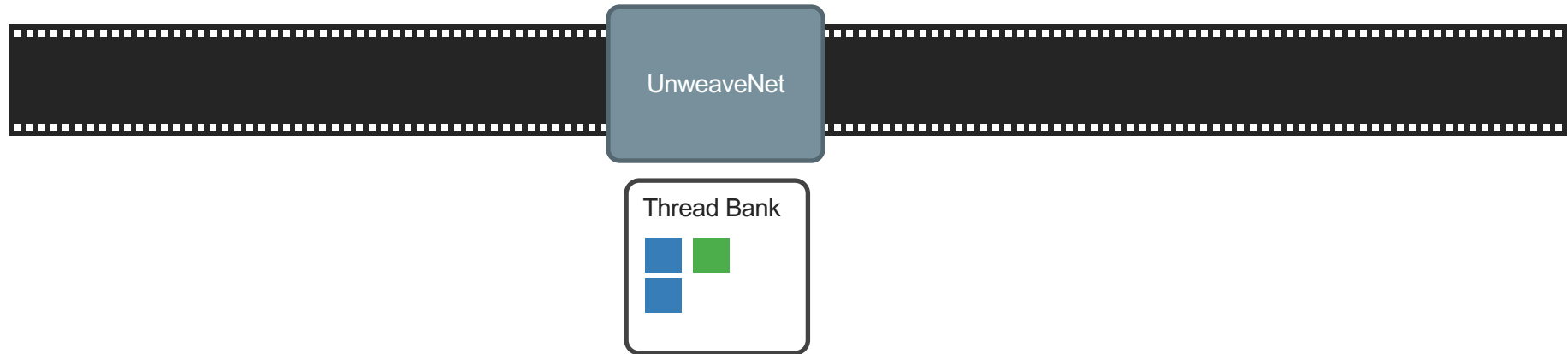
UnweaveNet

with: Will Price
Carl Vondrick



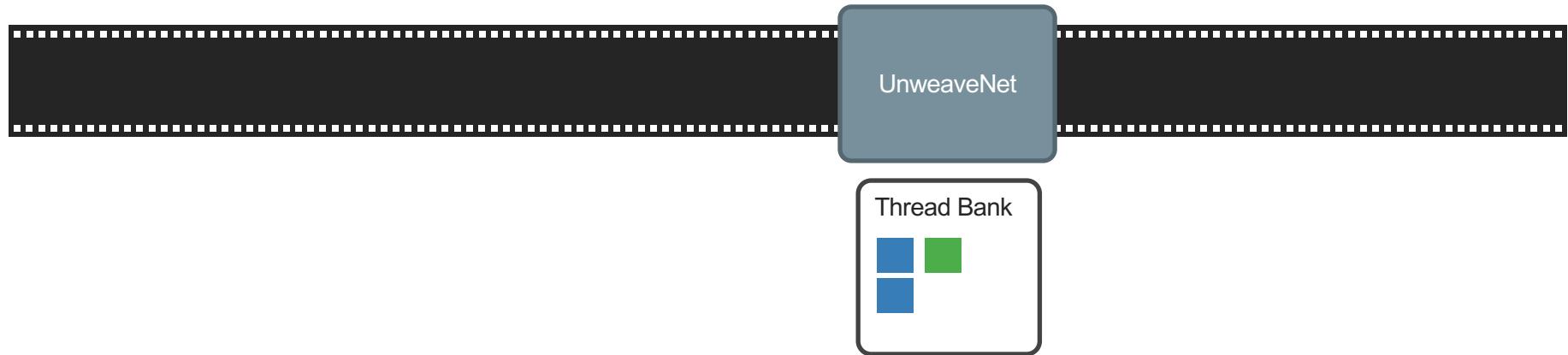
UnweaveNet

with: Will Price
Carl Vondrick



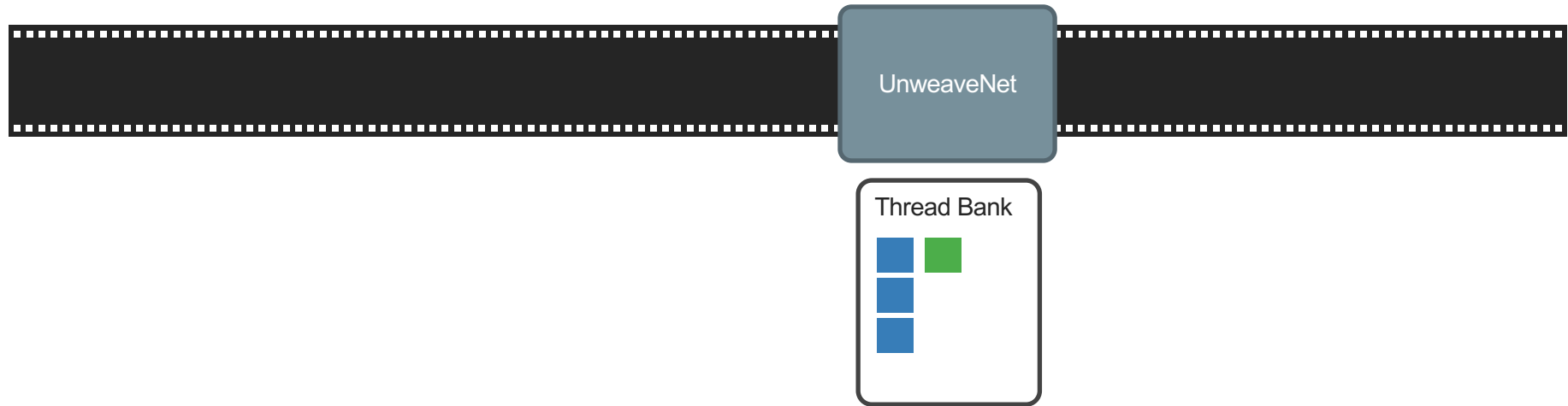
UnweaveNet

with: Will Price
Carl Vondrick



UnweaveNet

with: Will Price
Carl Vondrick



UnweaveNet

with: Will Price
Carl Vondrick



UnweaveNet

with: Will Price
Carl Vondrick

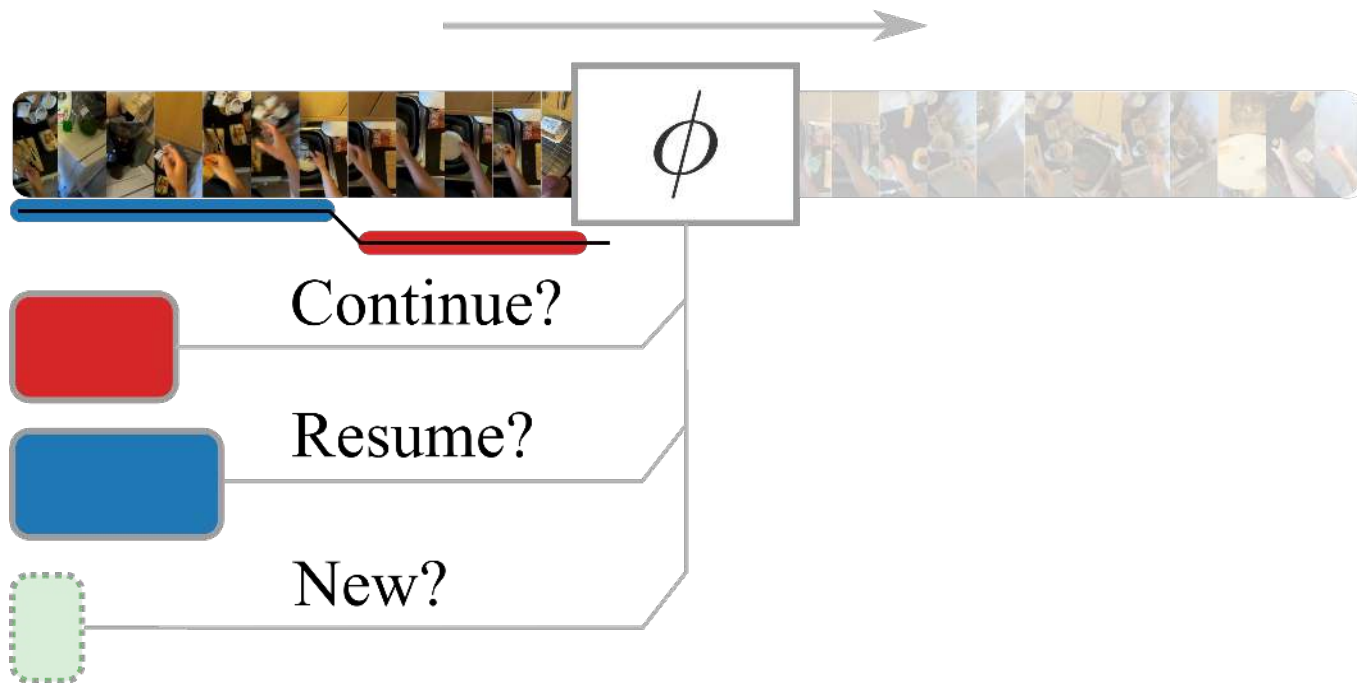


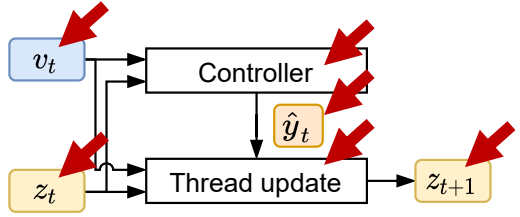
UnweaveNet

Thread Bank

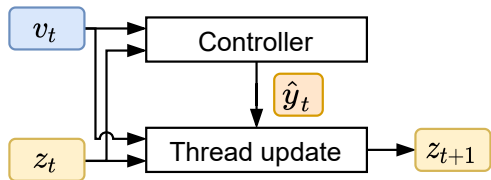


Unweaving

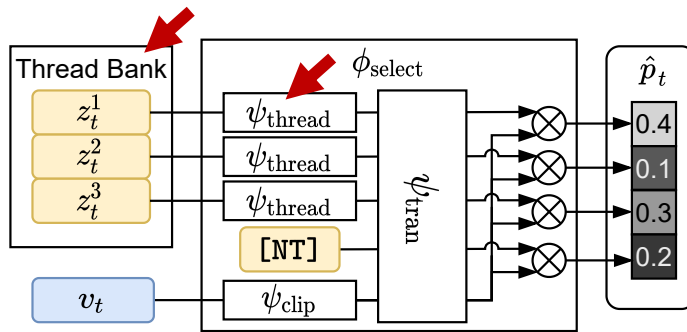




(a) UnweaveNet Overview



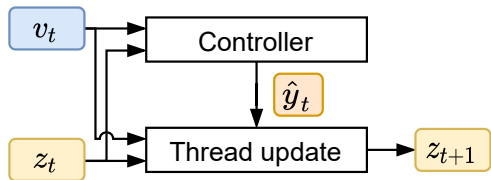
(a) UnweaveNet Overview



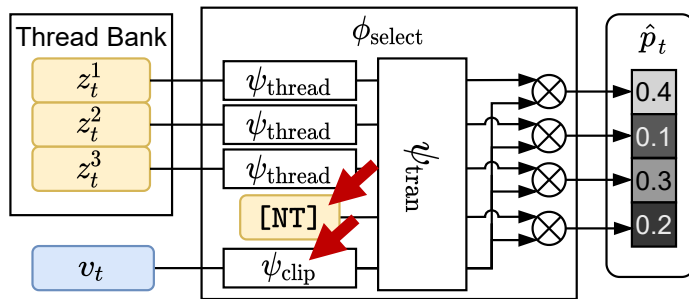
(b) Controller Architecture

Two learnt embeddings

$$\psi_{\text{thread}} : \mathbb{R}^D \rightarrow \mathbb{R}^E$$



(a) UnweaveNet Overview



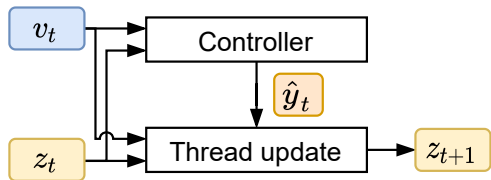
(b) Controller Architecture

Two learnt embeddings

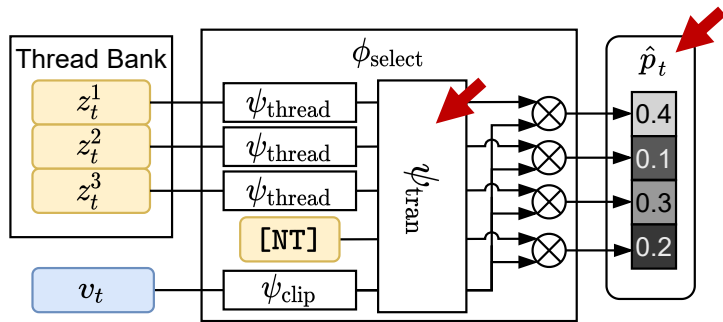
$$\psi_{\text{thread}} : \mathbb{R}^D \rightarrow \mathbb{R}^E$$

$$\psi_{\text{clip}} : \mathbb{R}^C \rightarrow \mathbb{R}^E$$

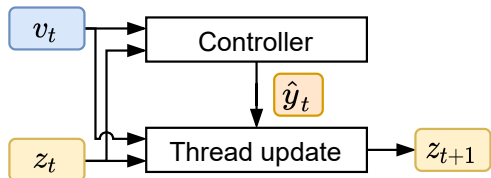
Learnt Encoding $[\text{NT}] \in \mathbb{R}^E$



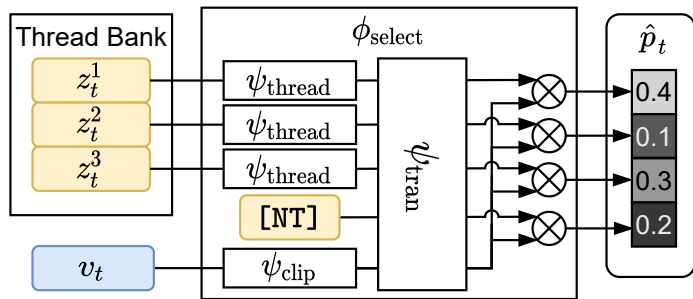
(a) UnweaveNet Overview



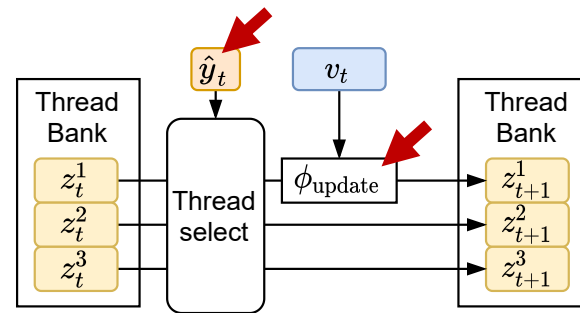
(b) Controller Architecture



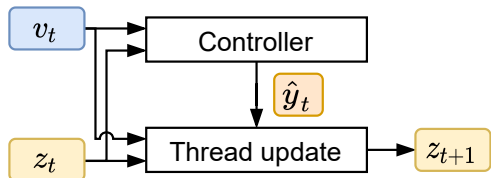
(a) UnweaveNet Overview



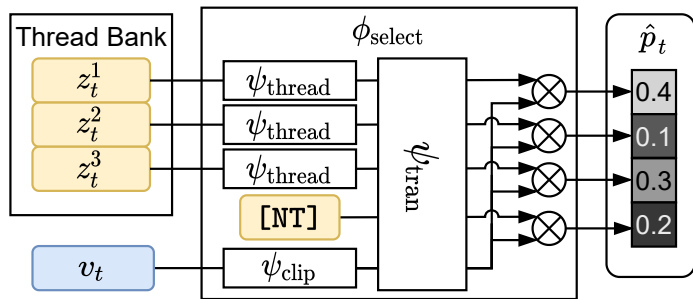
(b) Controller Architecture



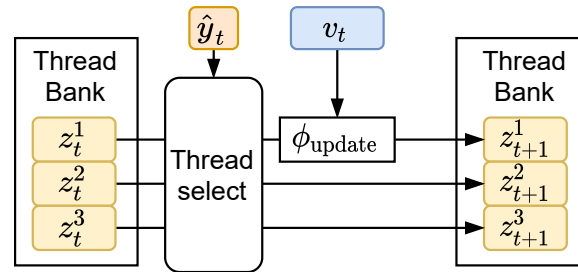
(c) Thread bank update



(a) UnweaveNet Overview



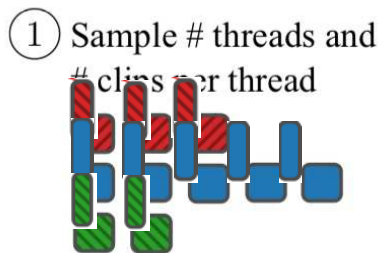
(b) Controller Architecture



(c) Thread bank update

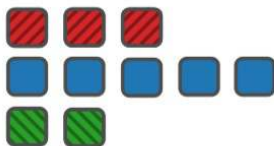
- Trained **end-to-end** including the backbone for clip features
- decisions made by ϕ_{select} are supervised using **teacher forcing**
 - at each time step, z_t is populated according to the ground-truth assignments $y_{1:t-1}$
 - A loss is then imposed on the output p_t given the correct decision y_t with focal hyperparameter γ due to the imbalance in decisions

- We propose self-supervised pretraining for UnweaveNet that samples threads from different parts of a long video and synthetically forms woven activity stories.



- We propose self-supervised pretraining for UnweaveNet that samples threads from different parts of a long video and synthetically forms woven activity stories.

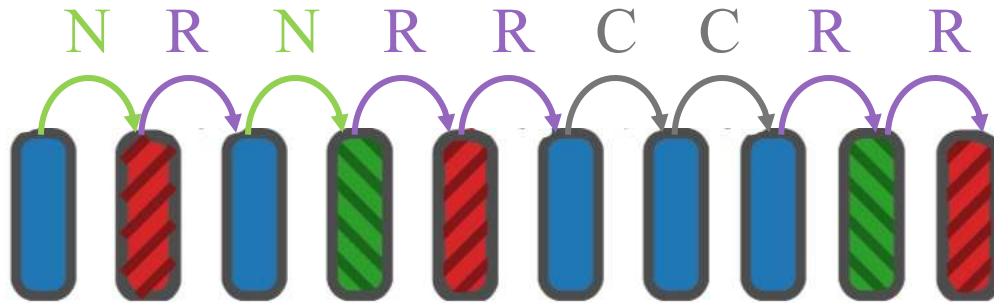
① Sample # threads and # clips per thread



② Position threads' clips within video



- 1000s of synthetic stories from training set.



- Labelled Sequences

Story Annotator

Story ID: c7c7f261-e5c7-4641-b0cf-b2b4e8c8ef3d

Video ID: P06_103

Start time: 7567

Splits: train+val

Author:

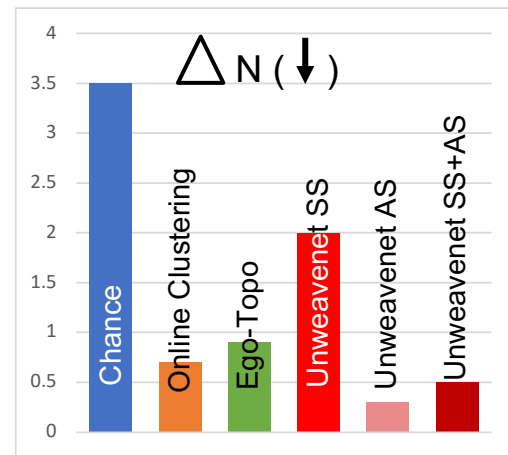
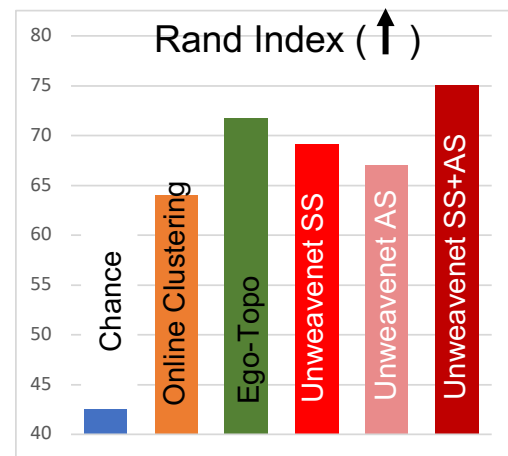
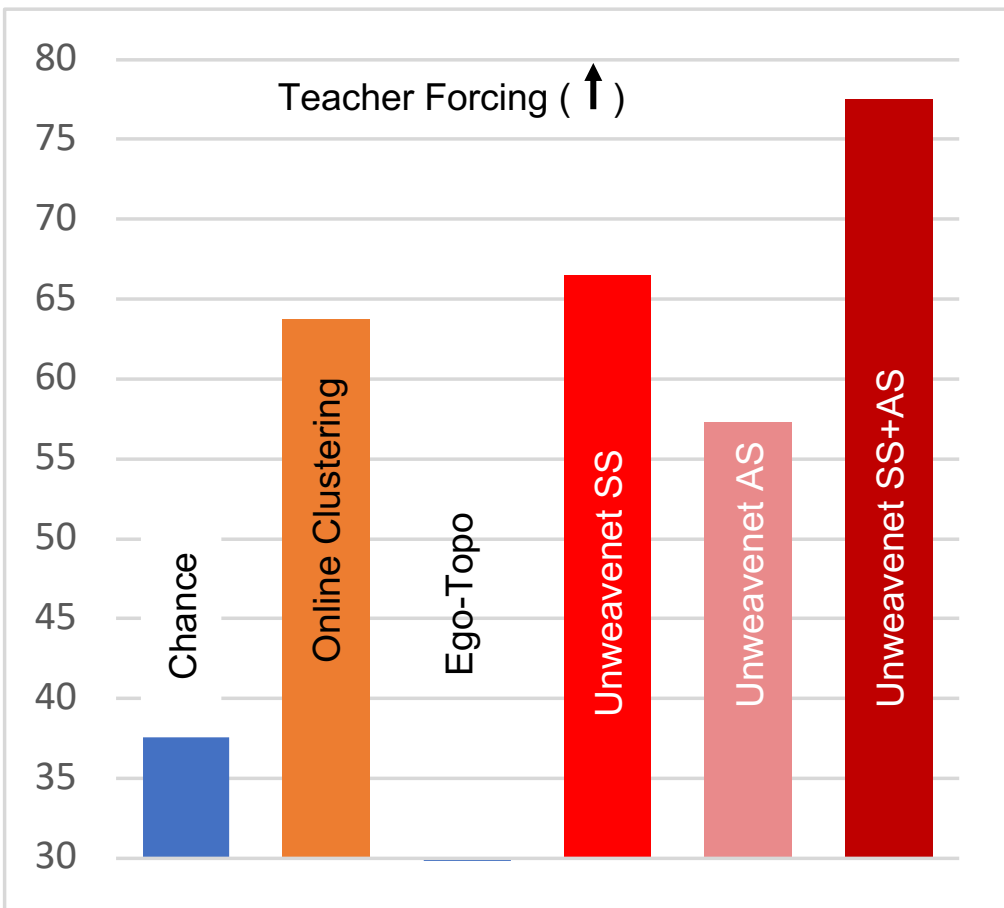


Split	# Threads		
	1	2	3
Train	718	201	32
Val	211	94	46
Test	50	50	50
Total	979	345	128

Table 1. EPIC-KITCHENS activity-story dataset by # of threads.

UnweaveNet

with: Will Price
Carl Vondrick



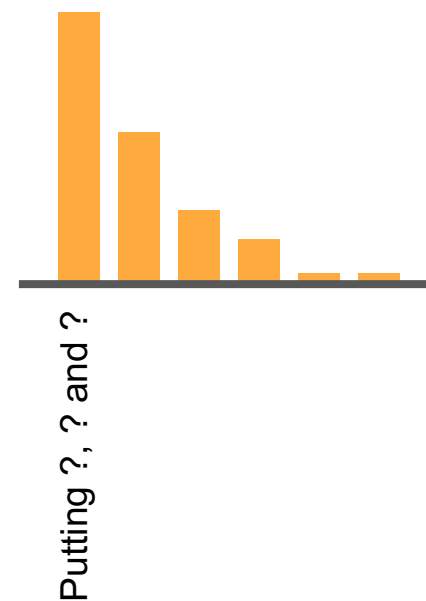
W Price, C Vondrick, D Damen (2022). UnweaveNet: Unweaving Activity Stories. CVPR

Carl Vondrick
July 26, 2024

Explainable?

Frame Attributions in Video Models

with: Will Price



Frame Attributions in Video Models

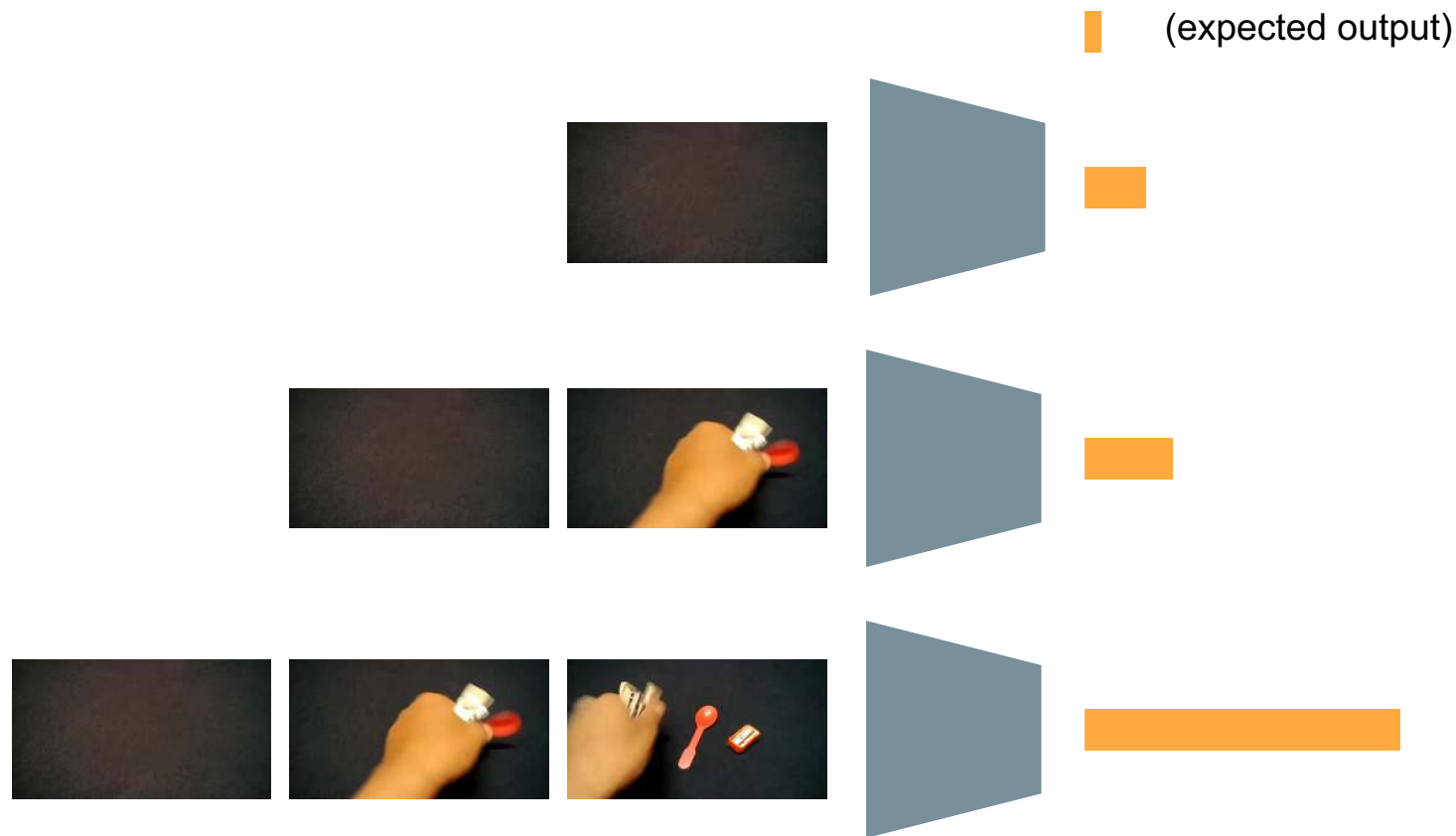
with: Will Price



Expected output
(Prior probability for
classification model)

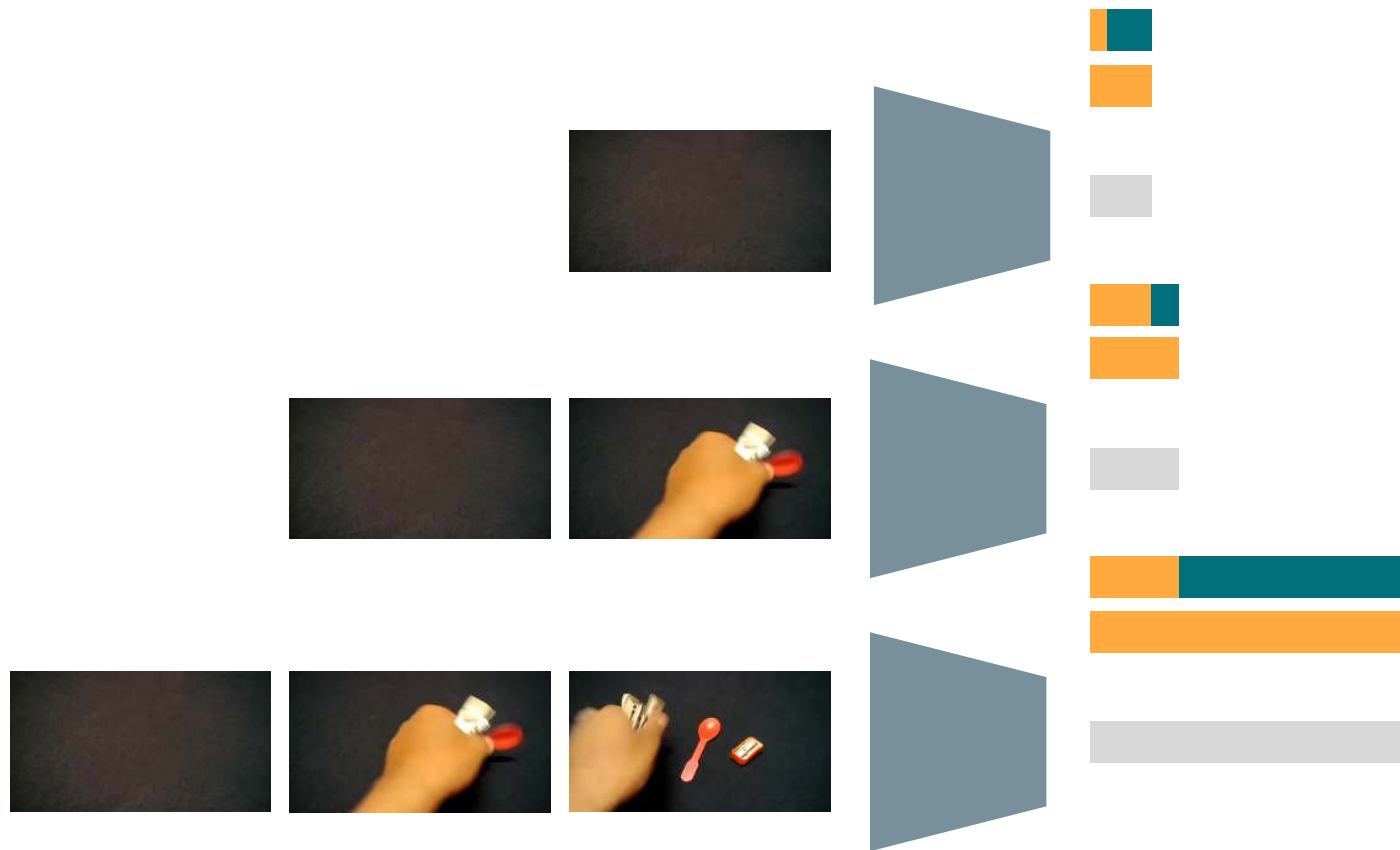
Frame Attributions in Video Models

with: Will Price



Frame Attributions in Video Models

with: Will Price



Frame Attributions in Video Models

with: Will Price



MODEL



MODEL



Frame Attributions in Video Models

with: Will Price



Frame Attributions in Video Models

with: Will Price



MODEL

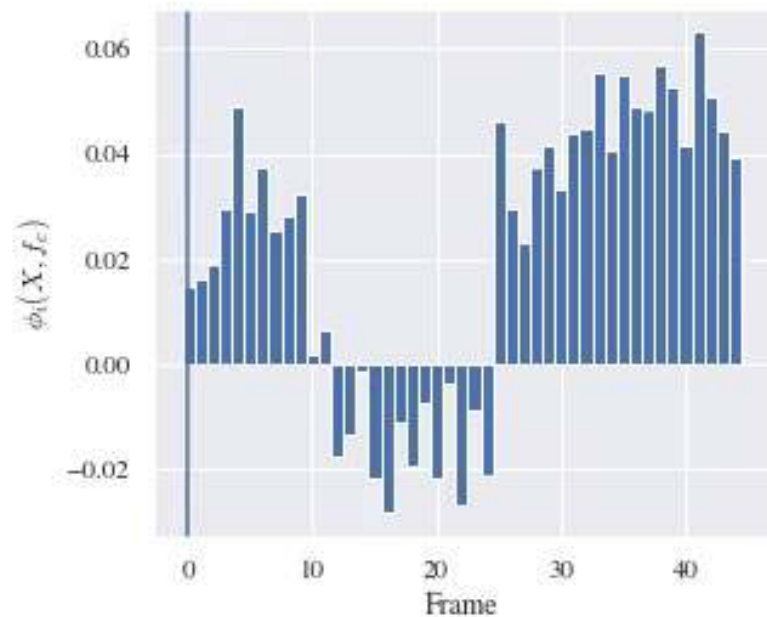
$$\Delta_3(\{1,2,4,5\}) = -.2$$



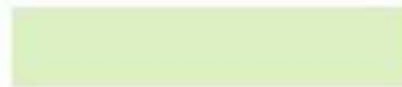
MODEL

Frame Attributions in Video Models

with: Will Price



Showing that something is empty



Frame Attributions in Video Models

with: Will Price
Tom Stark

ESVs Dashboard for Epic

Select a verb: open

Select a noun: drawer

Select a video: P01_103_84

Select number of frames: 1 2 3 4 5 6 7 8

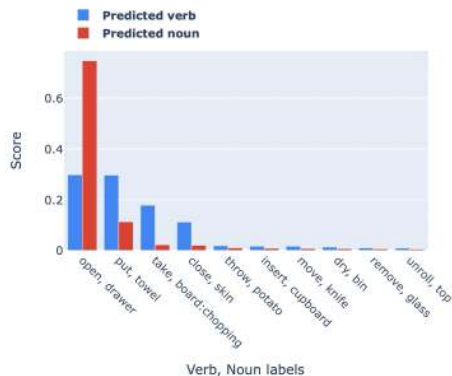
Original Video:

Selected frame: 2



Selected Verb: 3, Selected Noun: 8, Video P01_103_84

Model Predictions



ESV Predictions



Frame Attributions in Video Models

with: Will Price
Tom Stark

ESVs Dashboard for Epic

Select a verb

cut

Select a noun

tomato

Select a video

P01_17_126

Select number of frames



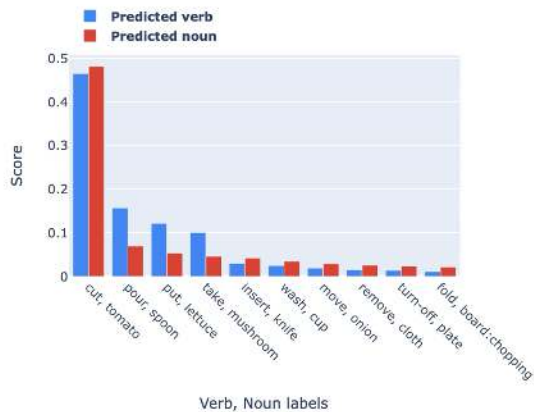
Original Video:

Selected frame: 529

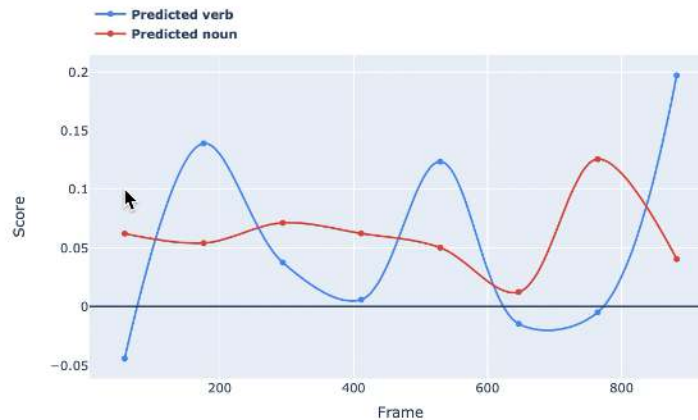


Selected Verb: 7, Selected Noun: 43, Video P01_17_126

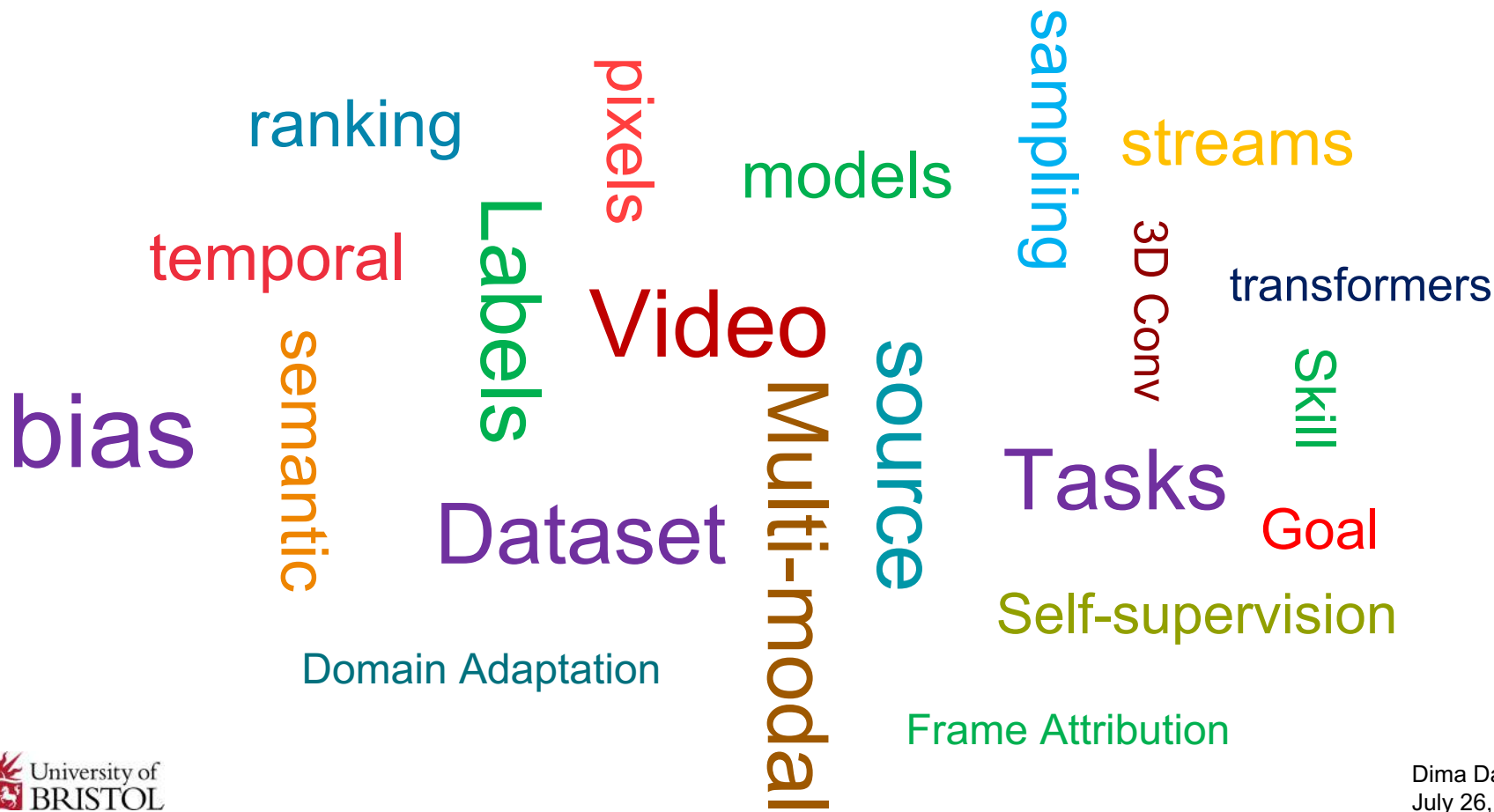
Model Predictions



ESV Predictions



Summary Wordle



and many more...

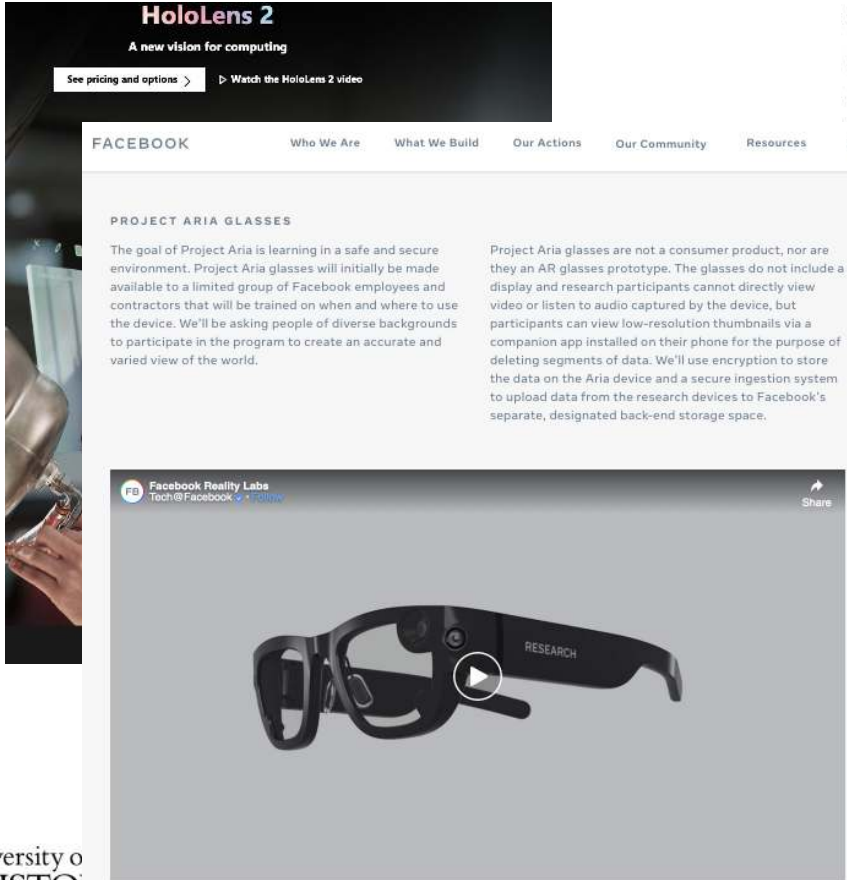




Video Understanding

An Egocentric Perspective

The future is here...



HoloLens 2
A new vision for computing
See pricing and options > Watch the HoloLens 2 video


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



Samsung patent application reveals augmented reality headset design

It comes as the Gear VR slowly fades away

by Jon Porter



The future is here...



Let's start with a show of hands...



Hands-Up if you are ready to wear a head-mounted or glass-mounted camera...



Hands-Up if this is NOT the future...



A world of isolated individuals....



Dangerous for crossing the road...



Mind-altering...

45 years ago...



Wa

Walkman ban ok

WOODBRIDGE, N.J. (AP) — Trians and cyclists who tune in to tape players with lightweight also tuning out traffic hazards, sa

Personal St

'Let'



Richard Butler, 24, rides the No. 30 (Jack-

Isolation is a result of plugging in

By John Jenks

Plug in to solitary entertainment. America. The Sony Walkman and its clones are here to stay.

Sales of the lightweight cassette tape recorders and radios continue to brisk in the Milwaukee area and "plugged in" people are becoming more and more visible.

The Walkman is a valuable weapon in the fight against boredom when walking, jogging and traveling. It's a good way to listen to music also for workers in noisy jobs — shutting out the noise of machinery and bringing in favorite music or other entertainment.

It also can be a source of solitude in the midst of the busy-busy of 1980s America — just plug in and tune out.

"As society gets more and more cluttered you will see an increase in



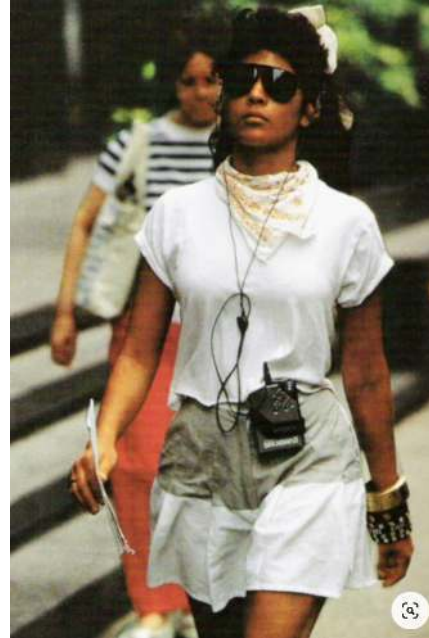
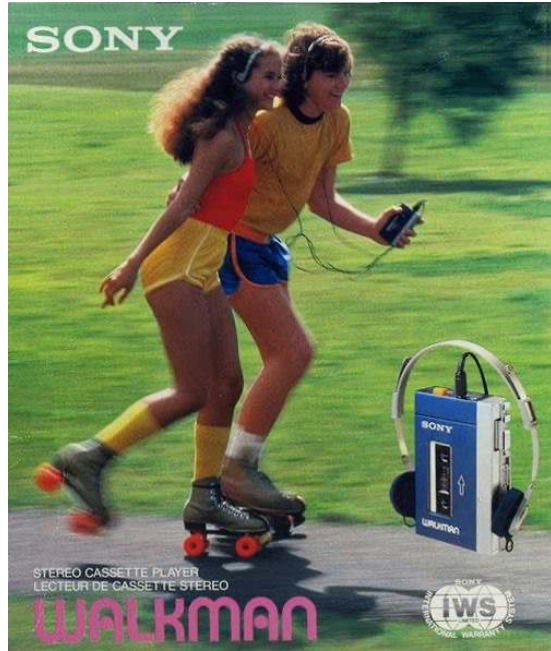
TUESDAY, 19 MARCH, 1966 9



Alone in the crowd: as more people opt for Walkman-type isolation, we are slowly peeling away from our fellow citizens, according to the latest American theory

Are we lost in a world of our own?

45 years ago...



How did we get here?



1907



1970s

1987



2004
GoPro



How did we get here?



How did we get here?



1907



1970s



1999
Eyetrapp



2013
Google
Glass

1987



1995



2004
GoPro



2024



Dima Damen
July 26, 2024

How did we get here?



2/15/1996

28 years ago...

How did we get here?



1907



70s



1999
Eyetrapp



2013
Google
Glass

14 years

1987

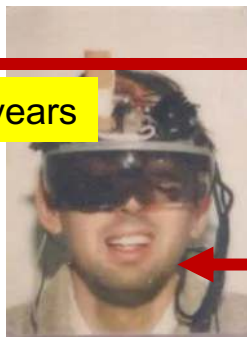
1995

2004
GoPro

2024



17 years



29 years



Dima Damen
July 26, 2024

How did we get here?

Motorcyclist fights fine for GoPro-style camera on helmet

A MOTORCYCLIST who was fined for having a GoPro-style camera mounted on his helmet. **Rebekah Cavanagh**
2 min read September 16, 2015 - 1:15

How Much Could Google's New Camera-Embedded Eyeglasses Change the Way People Use Entertainment, raise privacy questions

By **Kevin Kelleher**

Can Apple Rescue the Vision Pro?

The \$3,500 "spatial computing" device has gathered dust on my shelf. Can tweaks and upgrades save it from obsolescence?

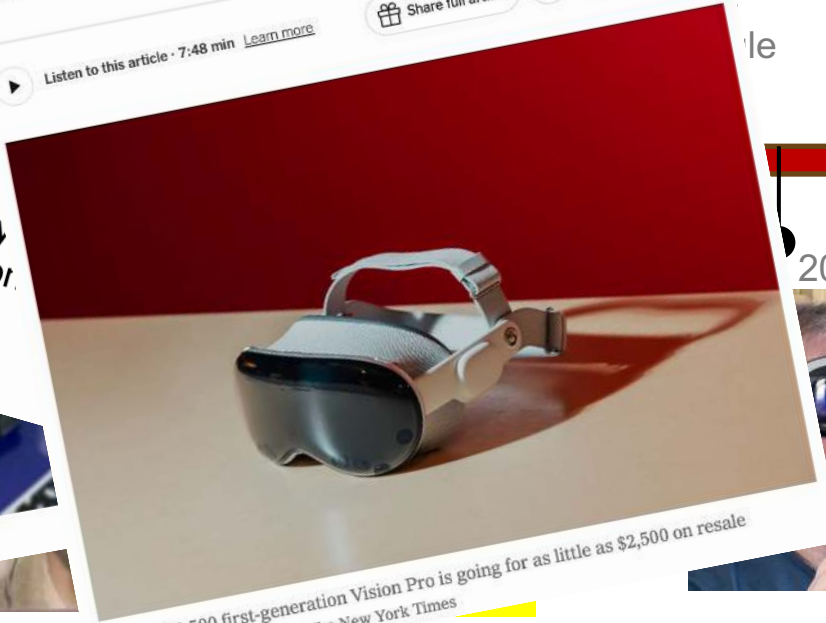
Listen to this article - 7:48 min [Learn more](#)

Share full article

106



Undated : GoPro HD Motorsport Hero camera - motorbike / motorcycle helmet mounted



Apple's \$3,500 first-generation Vision Pro is going for as little as \$2,500 on resale websites. **Clara Mokri for The New York Times**



2024



An Outlook into the Future of Egocentric Vision

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



Politecnico
di Torino



University of
BRISTOL

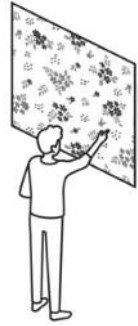


UNIVERSITÀ
degli STUDI
di CATANIA

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

How did we do this?

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



EGO-Designer



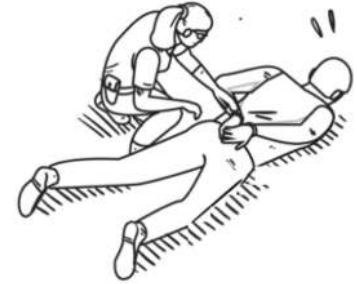
EGO-Tourist



EGO-Worker

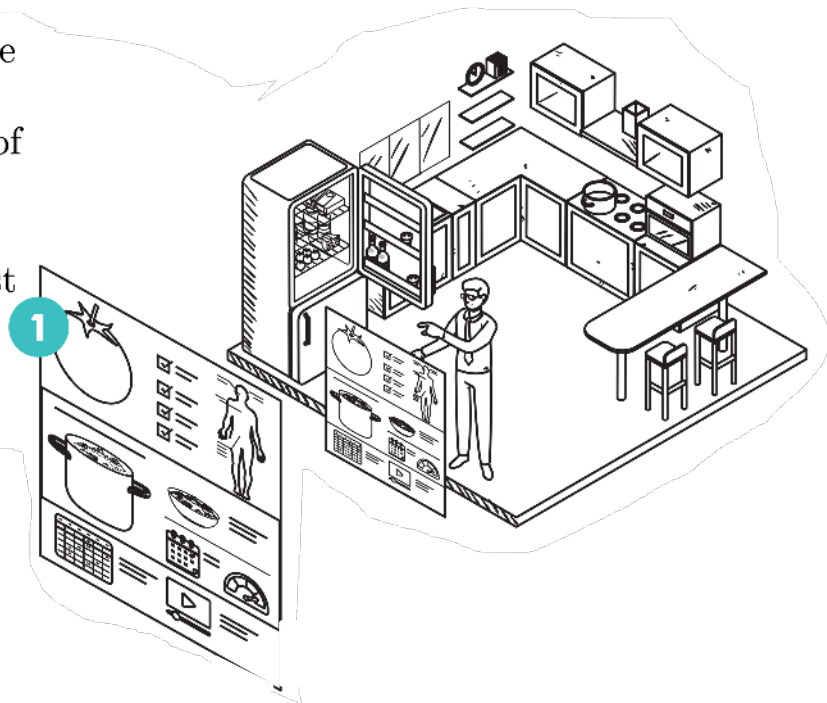


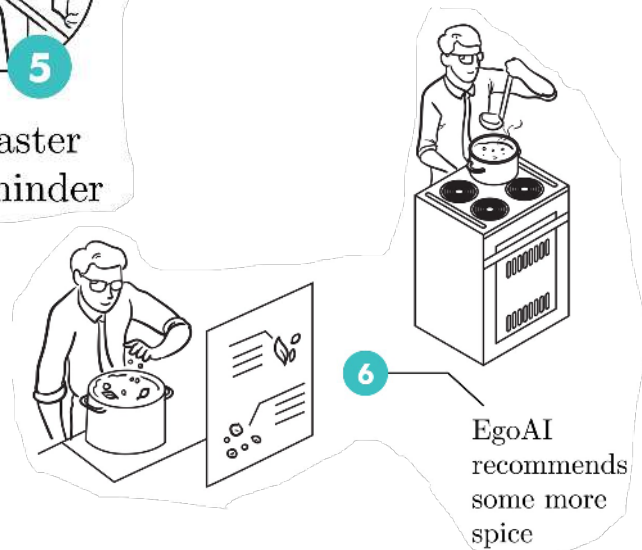
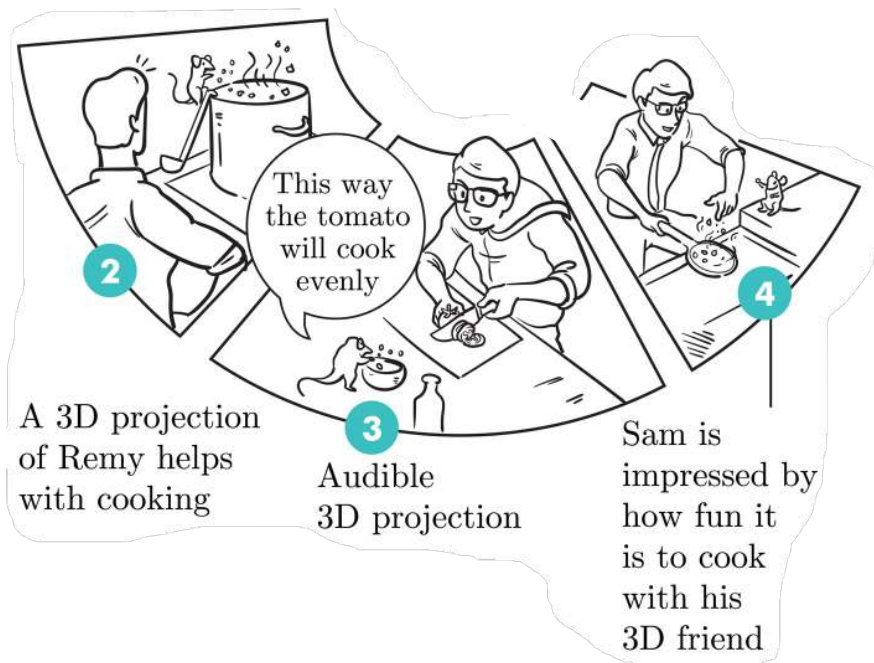
EGO-Home

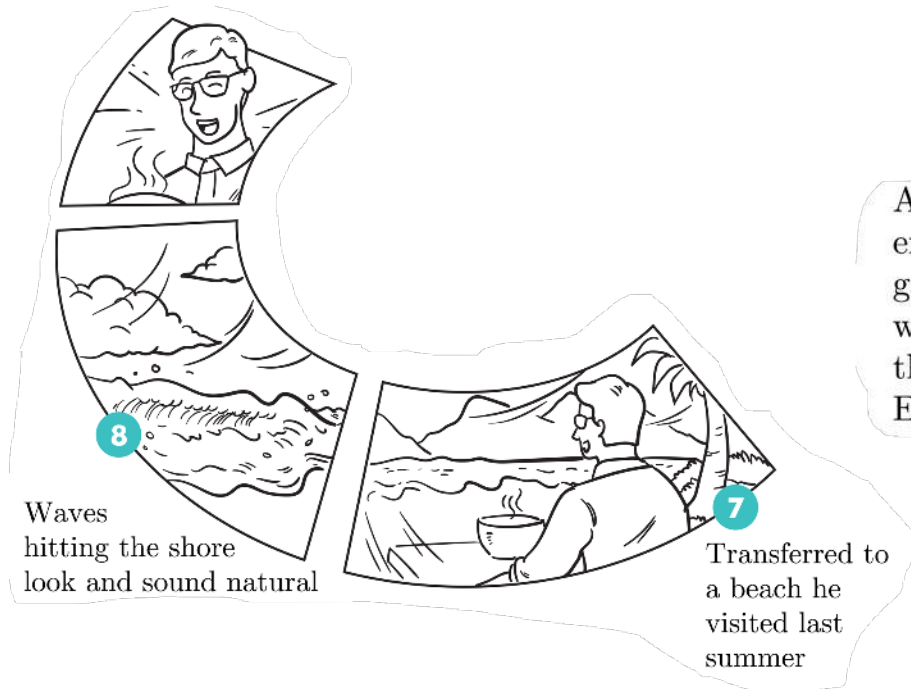


Ego-Police

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

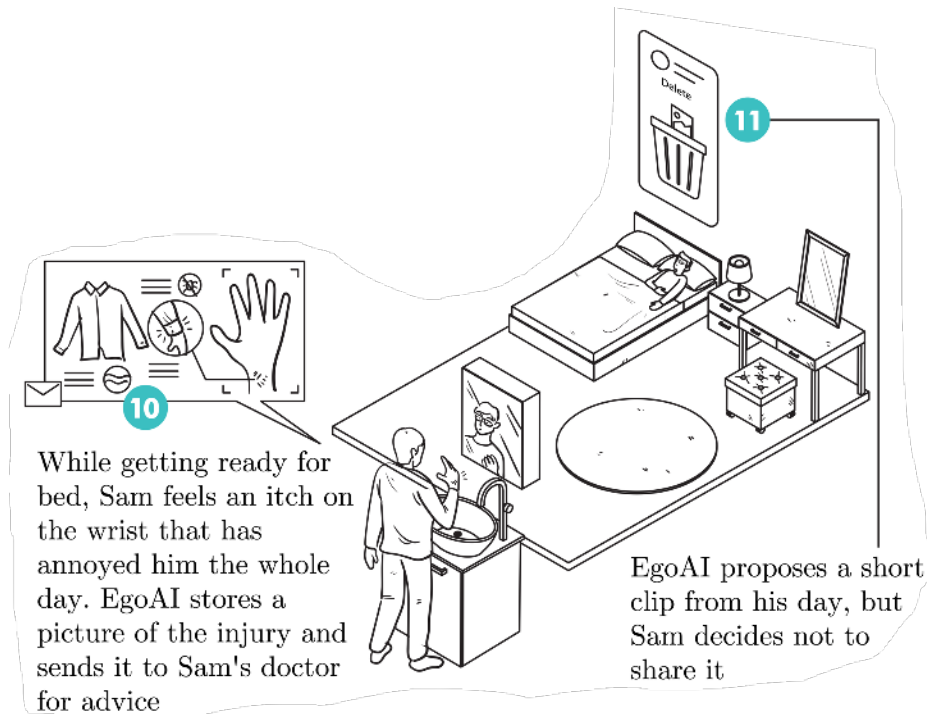






After dinner, Sam enjoys a group card game with his friends, who are connected through their own EgoAI

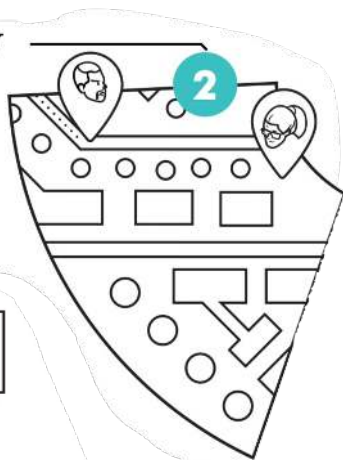




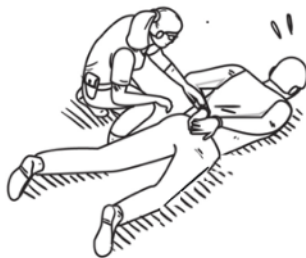
From Stories to Tasks

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



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



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

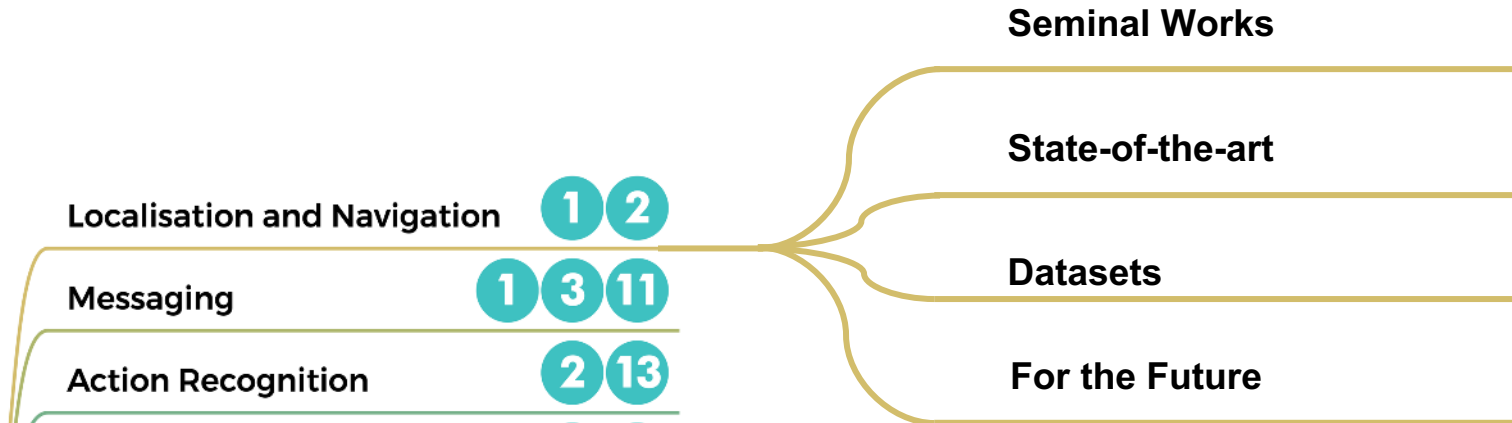
13

Privacy

14

The Survey Part

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

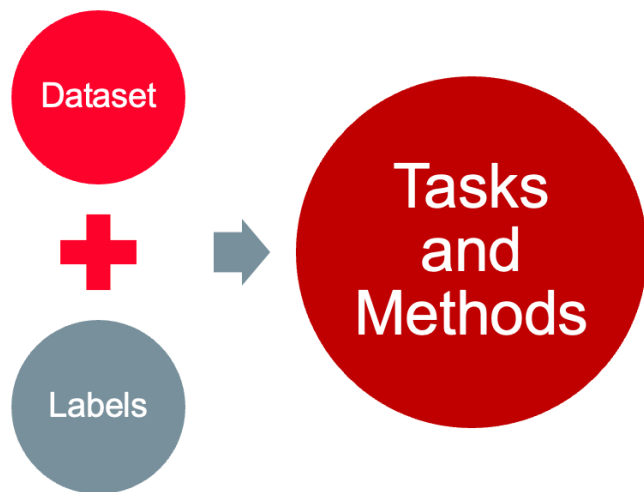
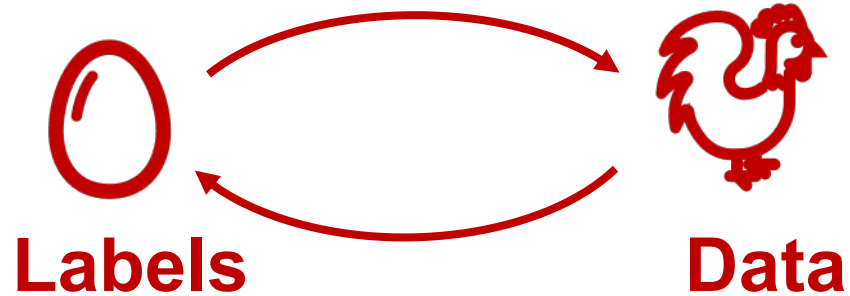


The Survey Part

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

- 12 tasks
- 46 pages (excluding references)
- 462 references

In this talk...



Thank you

For further info, datasets, code, publications...

<http://dimadamen.github.io>



@dimadamen



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

Q&A