



# A fine-grained Perspective onto Object Interactions

# Natural interactions

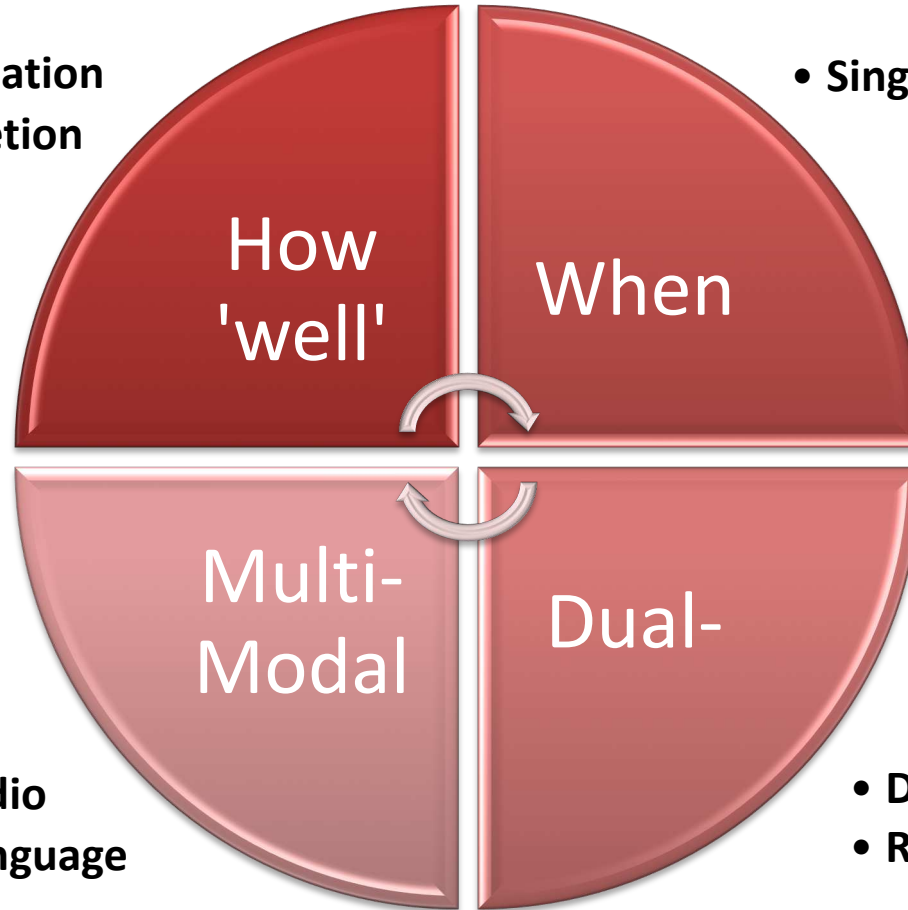
---



# Fine-Grained Object Interactions

---

- Skill Determination
- Action Completion



- Single-timestamp

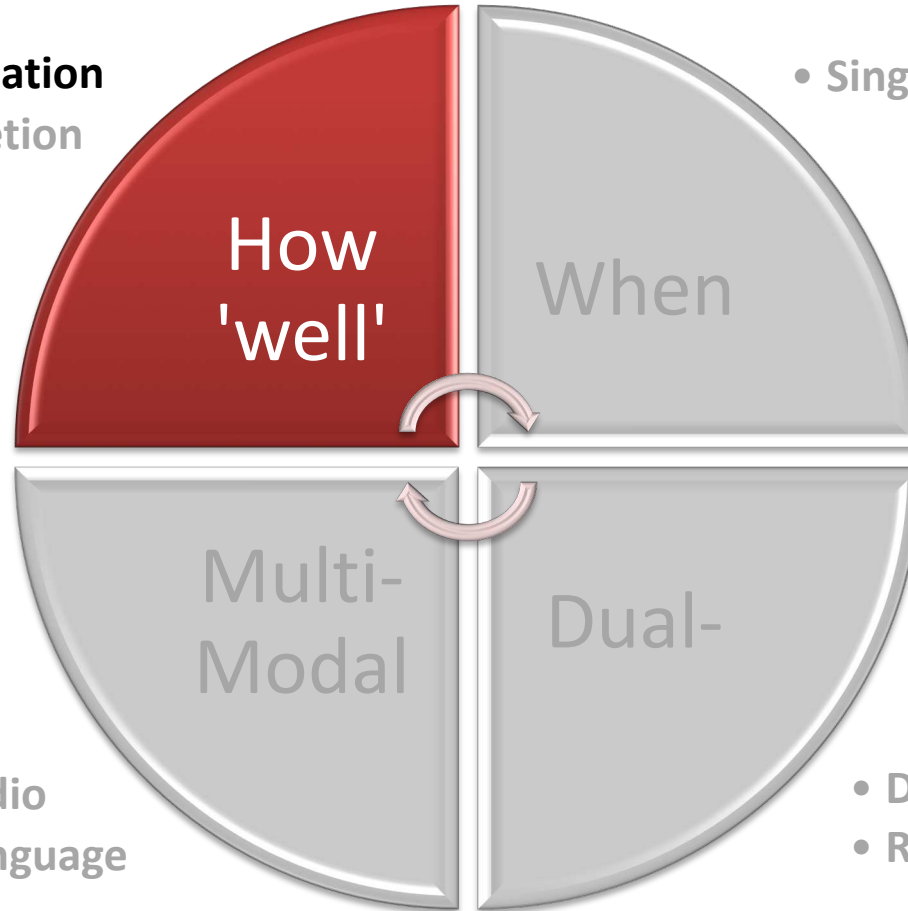
- Vision+Audio
- Vision+Language

- DDLSTM
- Retro-Actions

# Fine-Grained Object Interactions

---

- **Skill Determination**
- Action Completion



- Single-timestamp

- Vision+Audio
- Vision+Language

- DDLSTM
- Retro-Actions



# Who's Better? Who's Best? Skill Determination in Video using Deep Ranking

with: Hazel Doughty  
Walterio Mayol-Cuevas

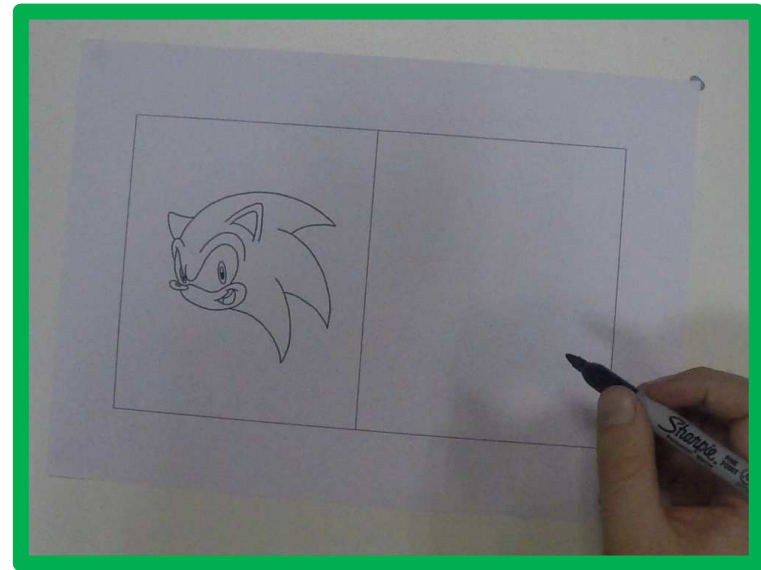
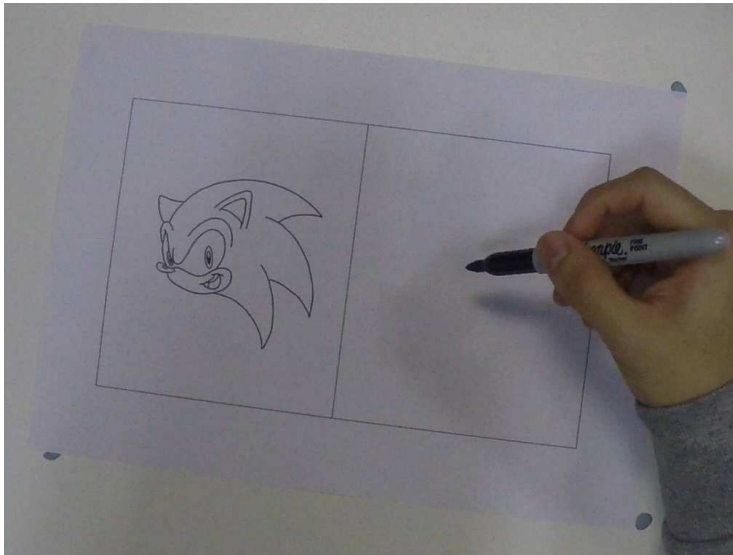


Assess relative skill for a collection of video sequences, applicable to a variety of tasks.

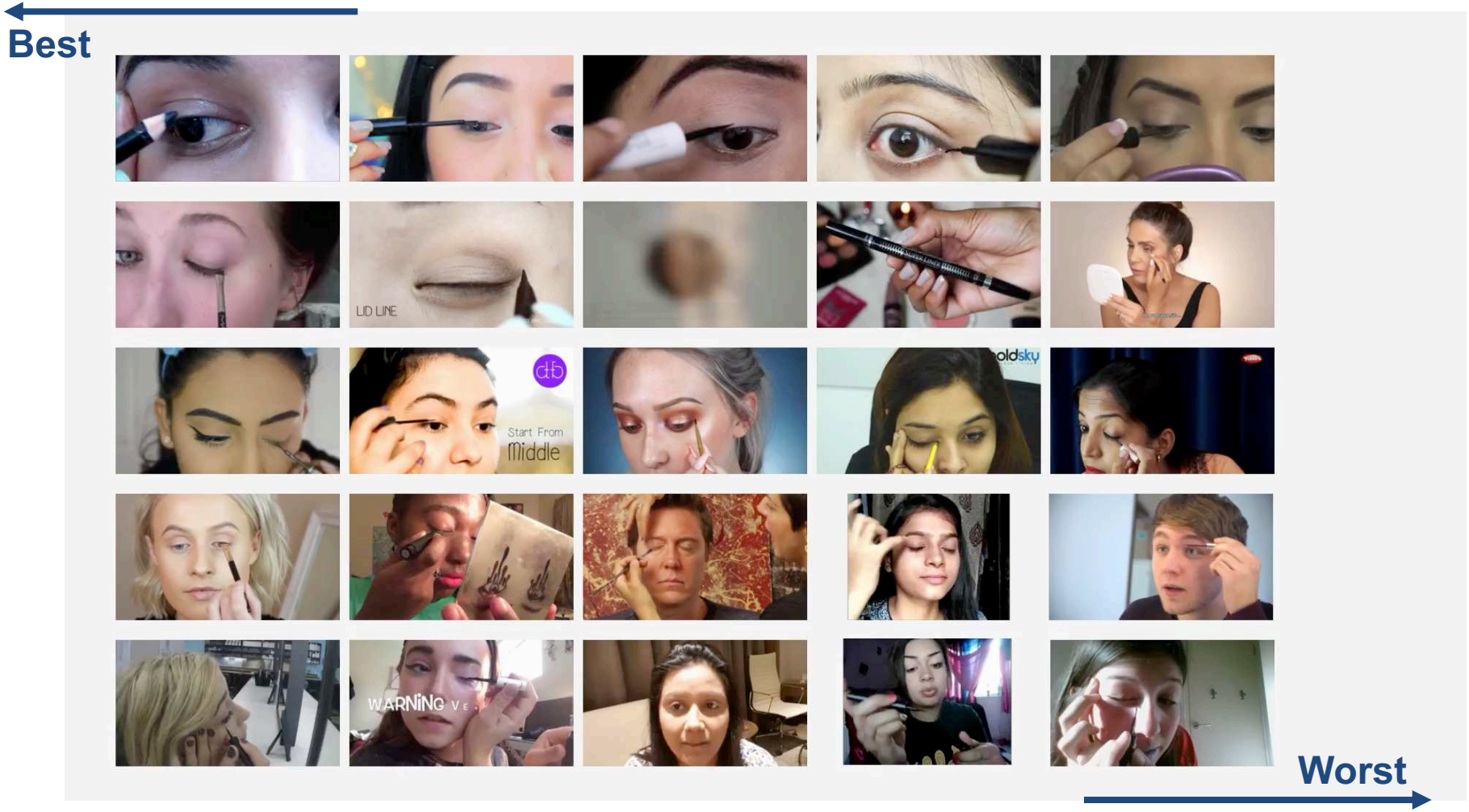
# Skill Determination from Video

---

**Input:** Pairwise annotations of videos, indicating higher skill or no skill preference



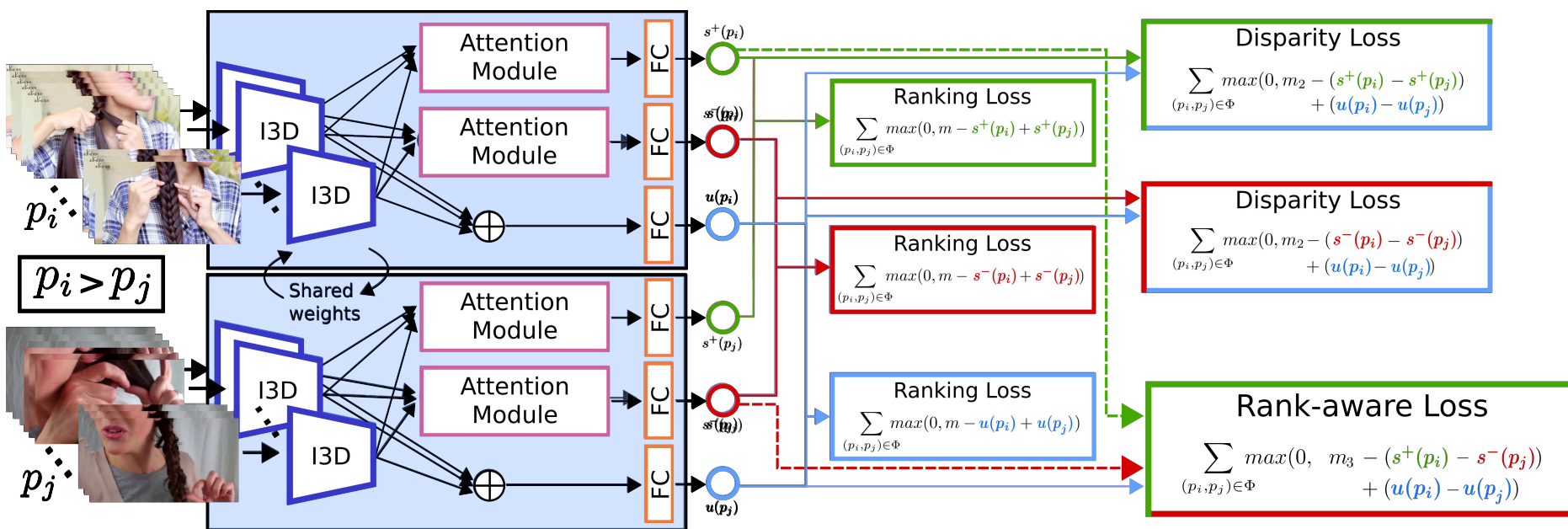
# Skill Determination in Video





# The Pros and Cons: Rank-Aware Temporal Attention

with: Hazel Doughty  
Walterio Mayol-Cuevas



# The Pros and Cons: Rank-Aware Temporal Attention

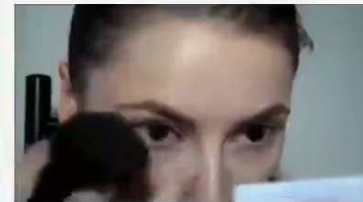
with: Hazel Doughty  
Walterio Mayol-Cuevas

## Low-skill Attention Module

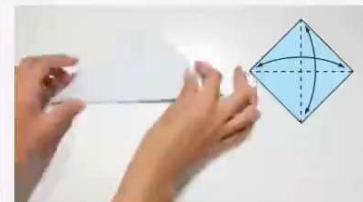
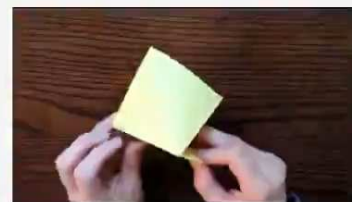
Surgery



Apply Eyeliner



Origami





# The Pros and Cons: Rank-Aware Temporal Attention

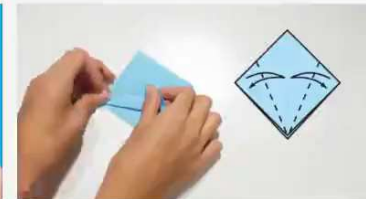
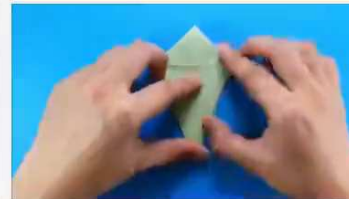
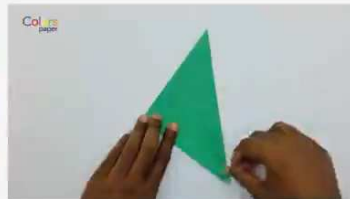
with: Hazel Doughty  
Walterio Mayol-Cuevas

## High-skill Attention Module

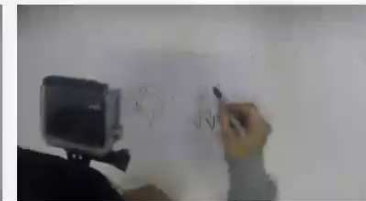
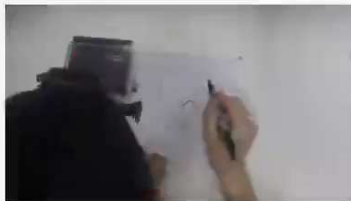
Dough  
Rolling



Origami



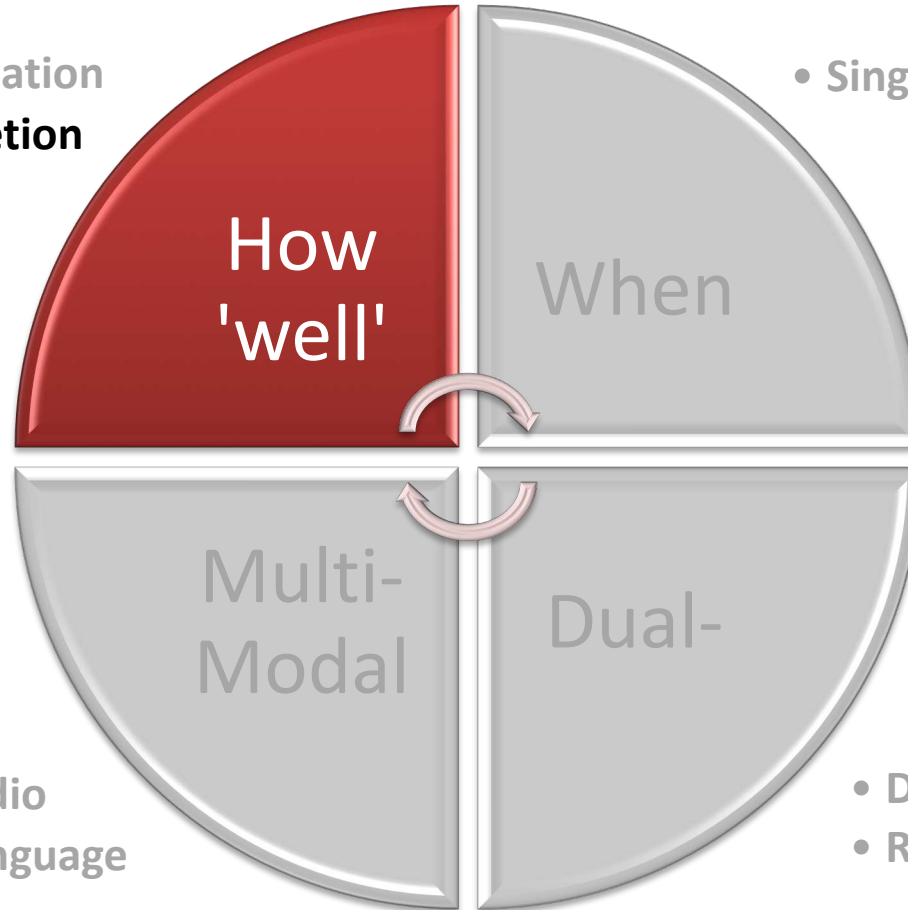
Drawing



# Fine-Grained Object Interactions

---

- Skill Determination
- **Action Completion**



- Single-timestamp

- Vision+Audio
- Vision+Language

- DDLSTM
- Retro-Actions

# Action Completion Detection

---



# Action Completion Detection

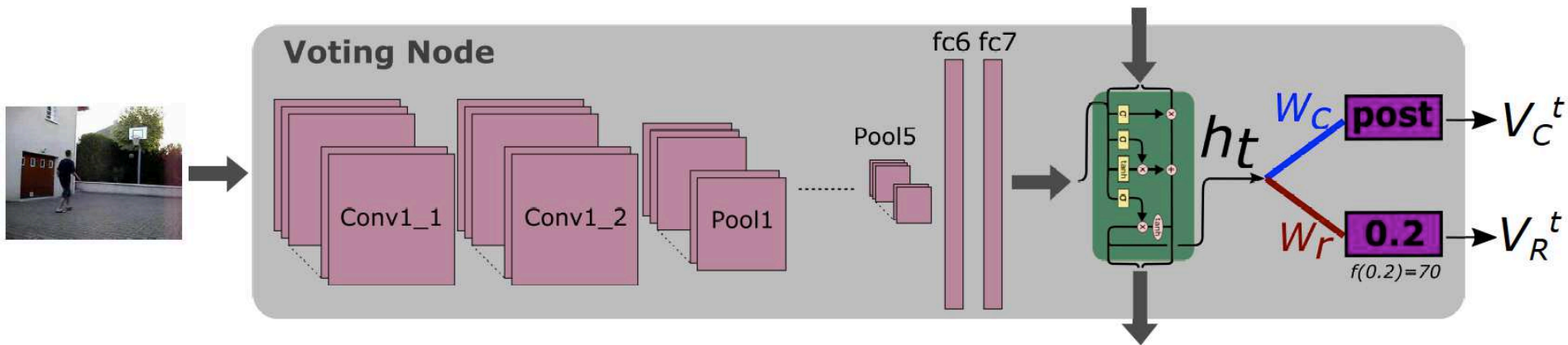


Pre-V ←  
V<sub>R</sub><sup>T</sup> ←  
C-C ←  
R-R ←  
R-C ←  
C-R ←  
Ground truth ←



# Action Completion Detection

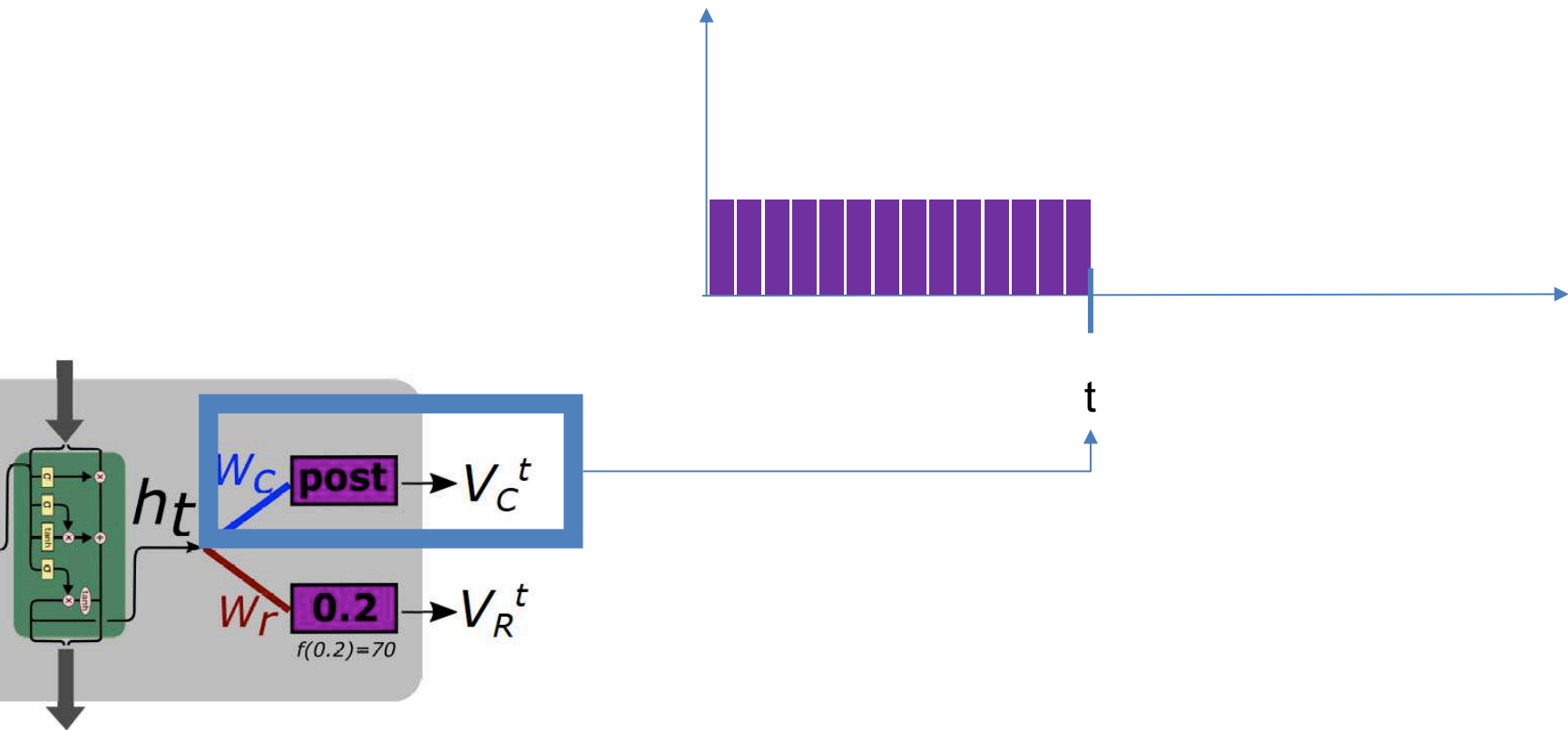
- Each frame in the sequence, contributes to the completion moment detection via 'voting'





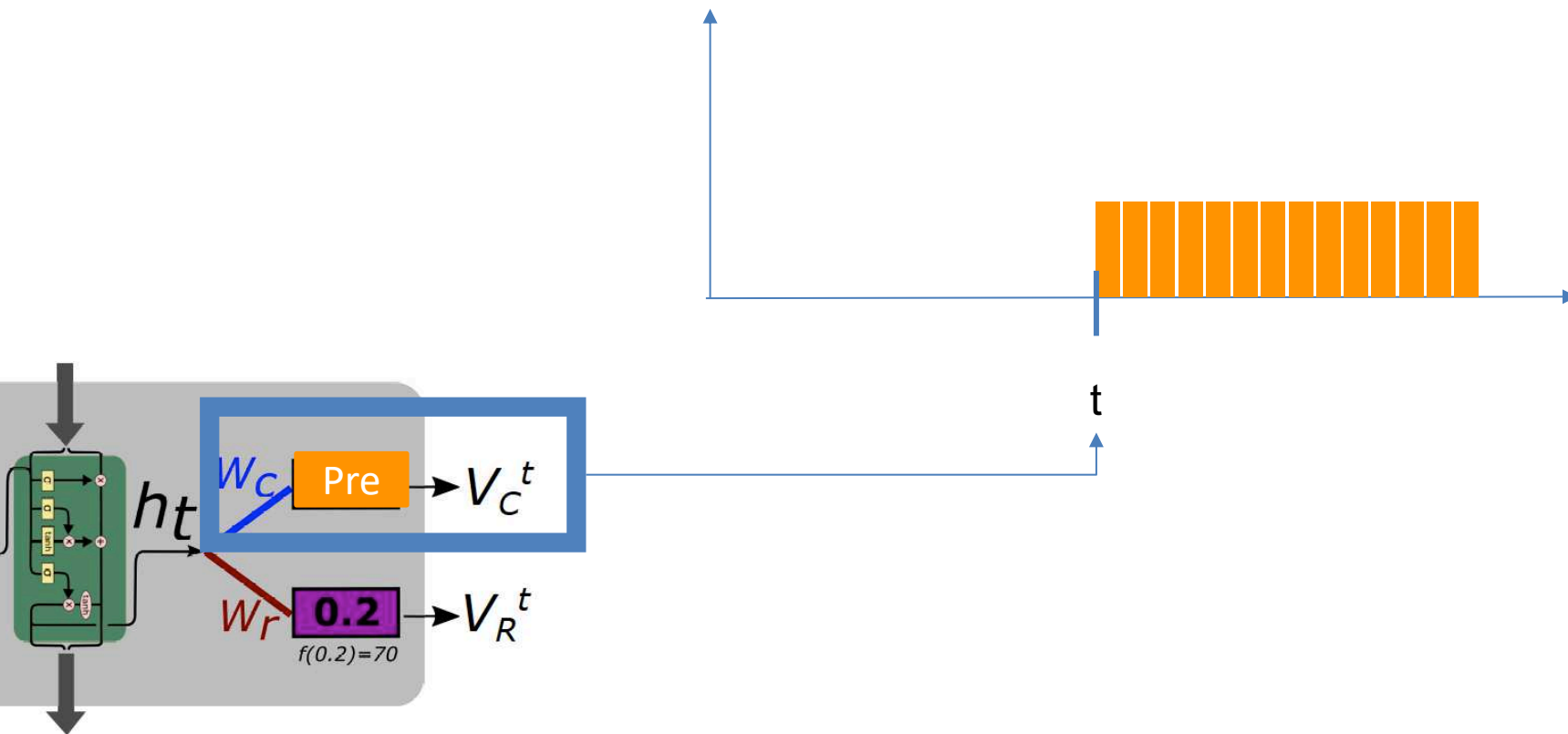
# Action Completion Detection

## 1. Classification-Based Voting



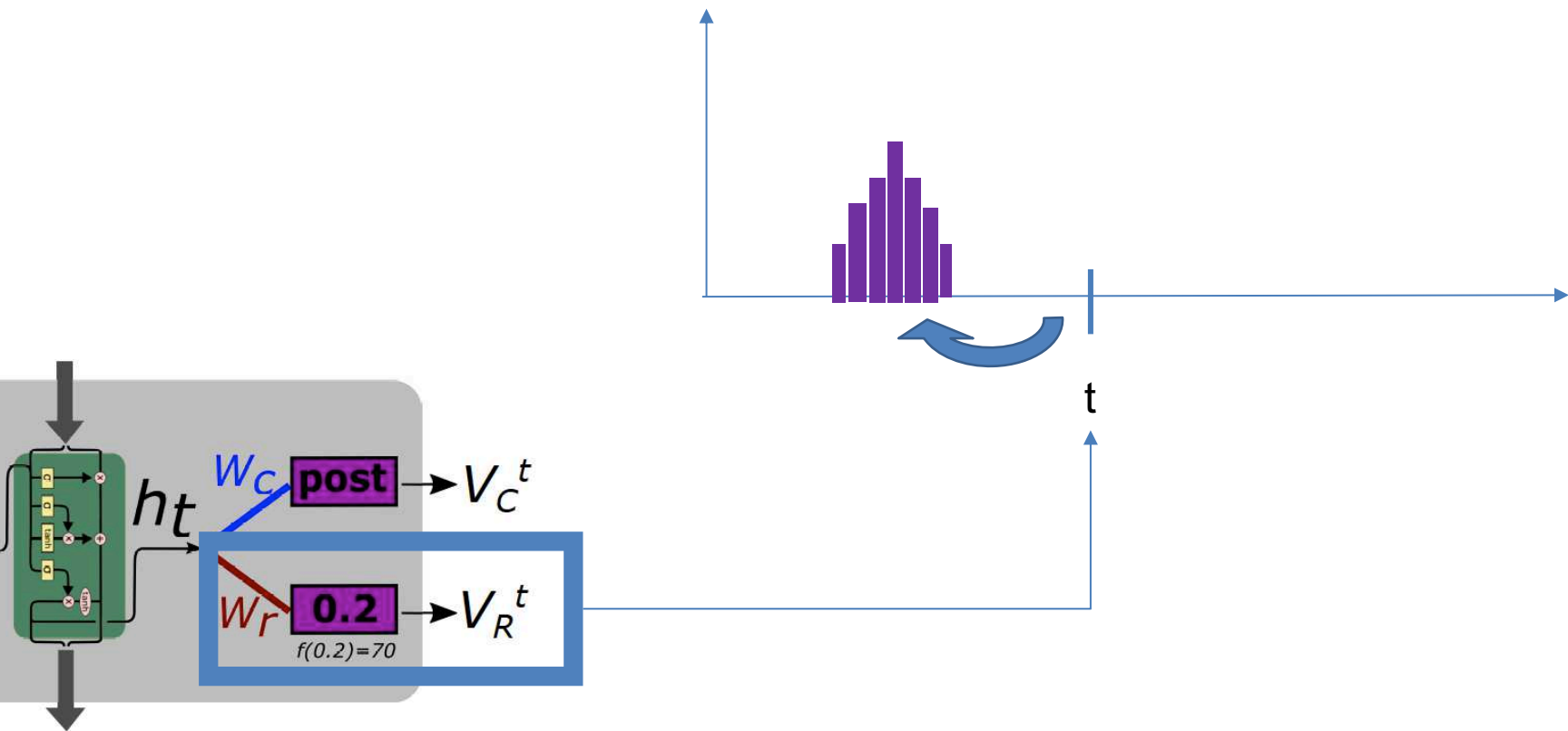
# Action Completion Detection

## 1. Classification-Based Voting



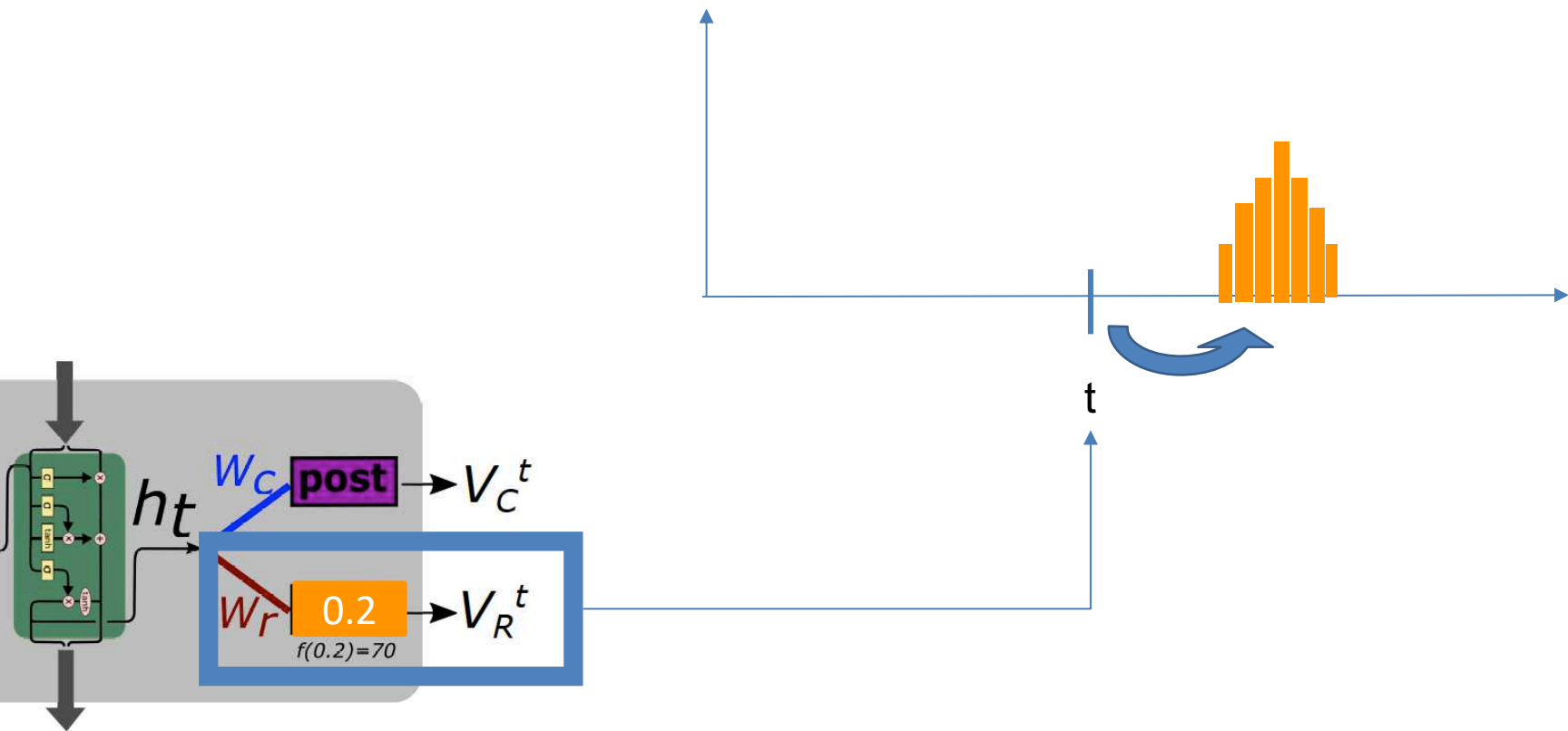
# Action Completion Detection

## 2. Regression-Based Voting

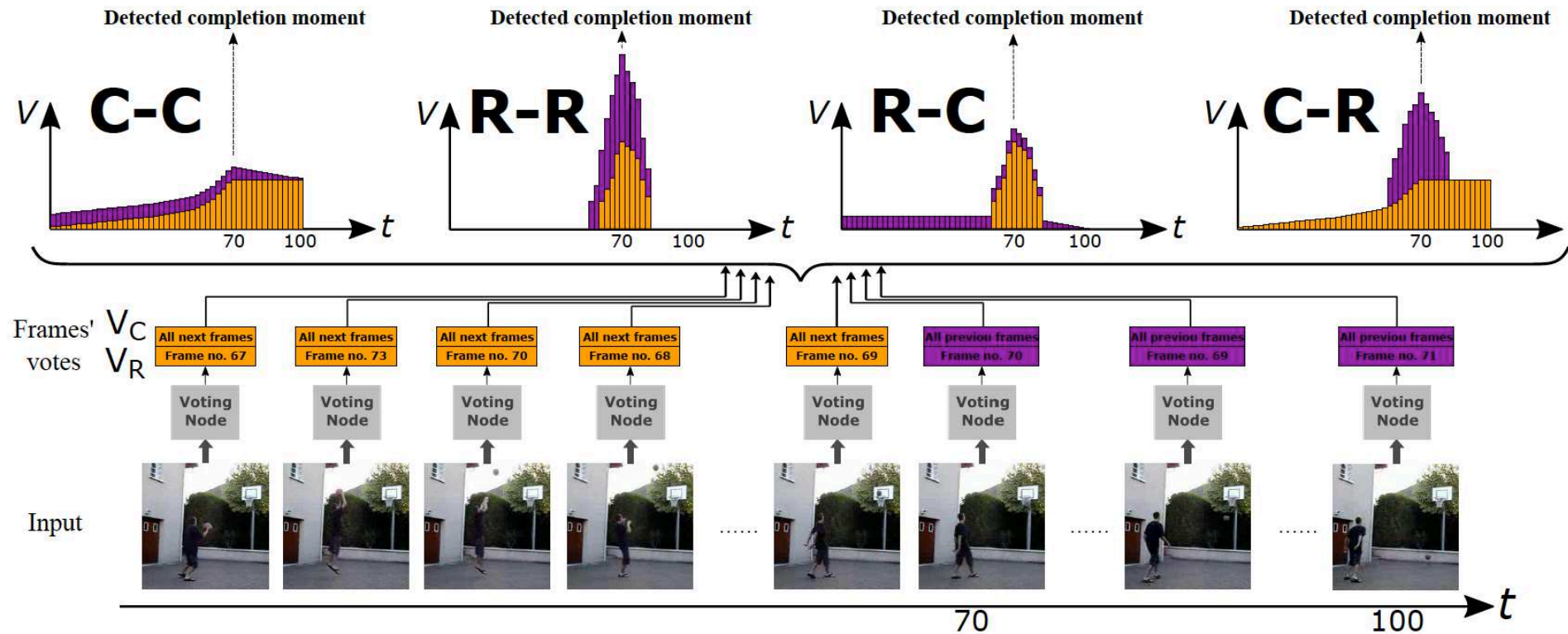


# Action Completion Detection

## 2. Regression-Based Voting



# Action Completion Detection



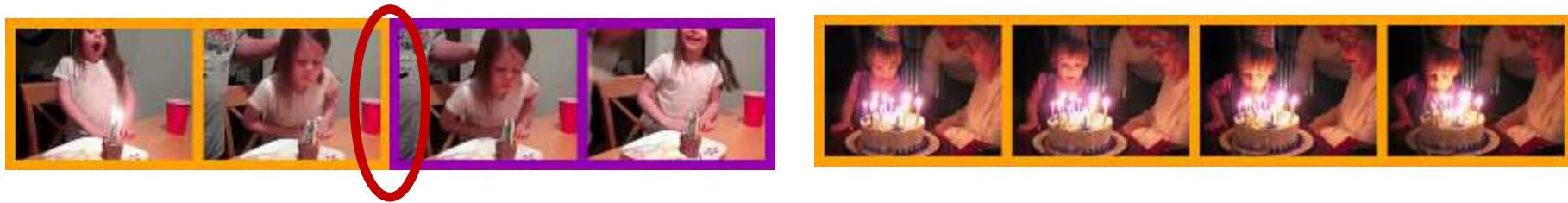


# Action Completion Detection



# Action Completion Detection

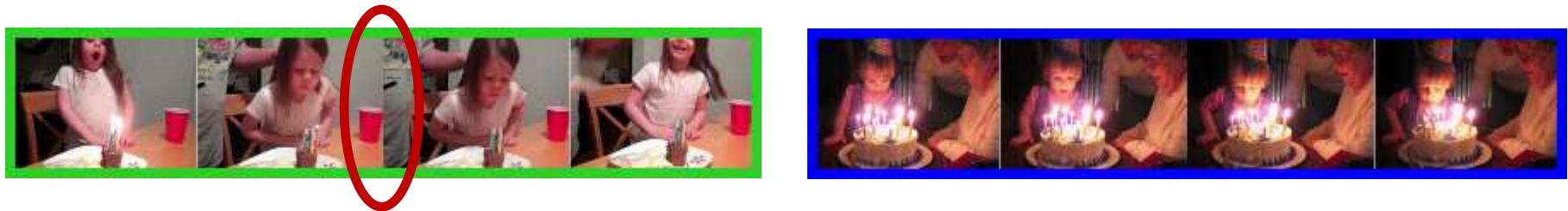
**Frame-level labels:** annotations are expensive, subjective and noisy.



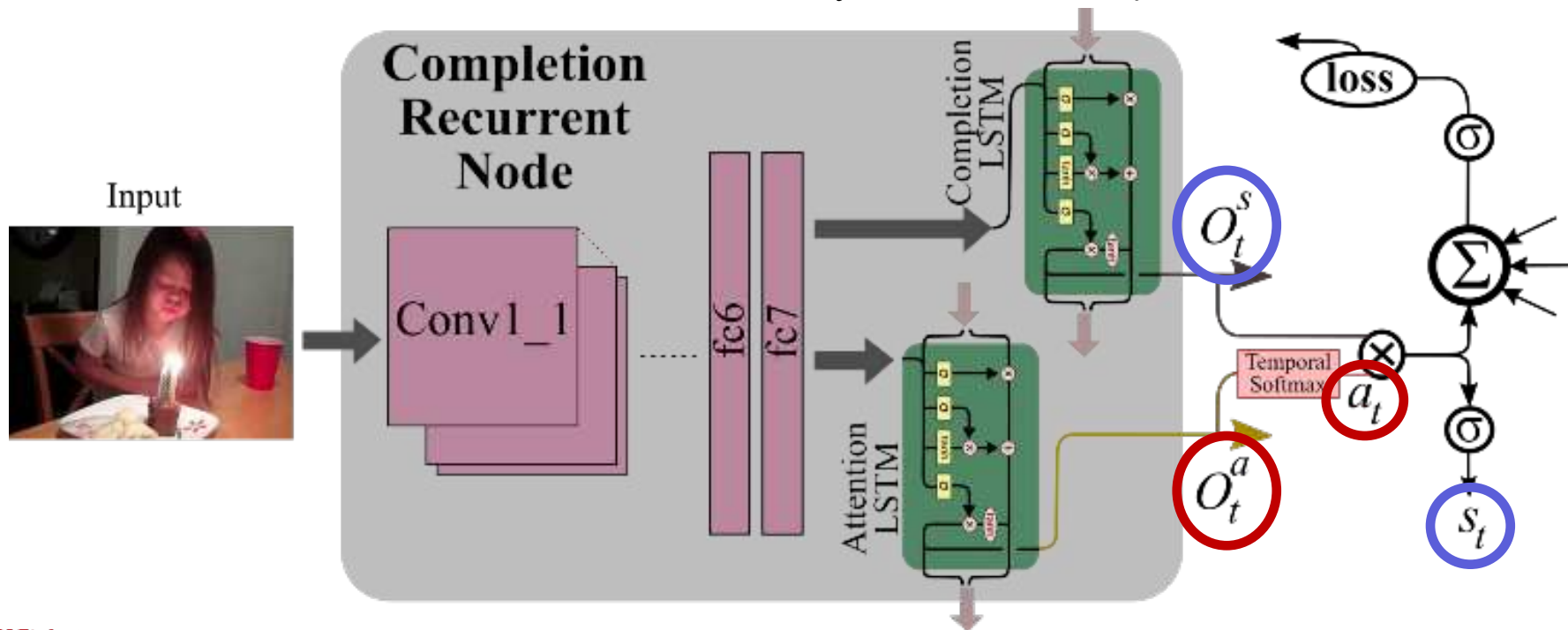
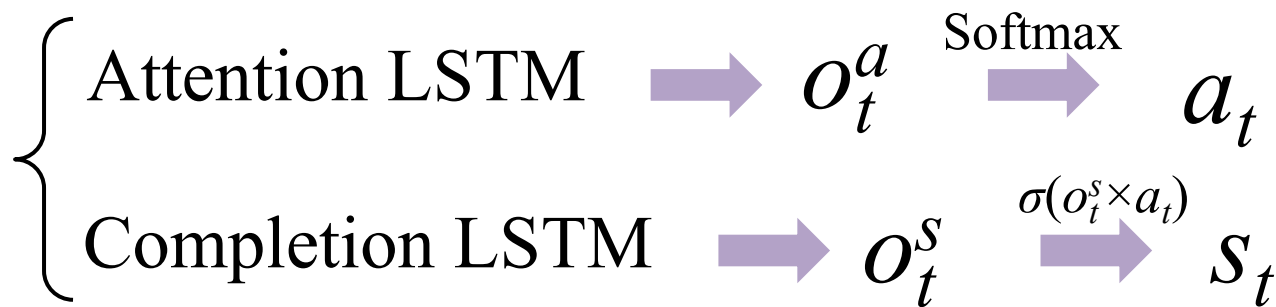
We detect completion using only weak labels during training.



sequence-level *complete* and *incomplete* labels



# Action Completion Detection



# Action Completion Detection

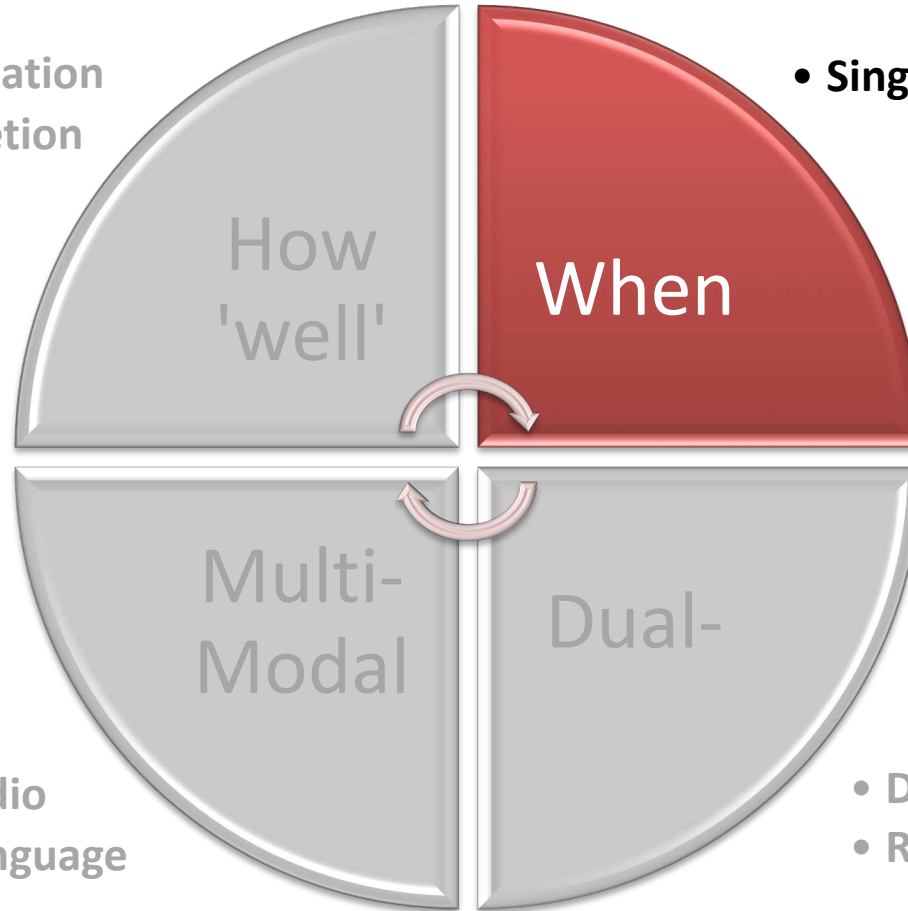


Completion scores ←  
Attention scores ←  
WS-U ←  
WS-Att ←  
GT ←

# Fine-Grained Object Interactions

---

- Skill Determination
- Action Completion



- **Single-timestamp**

- Vision+Audio
- Vision+Language

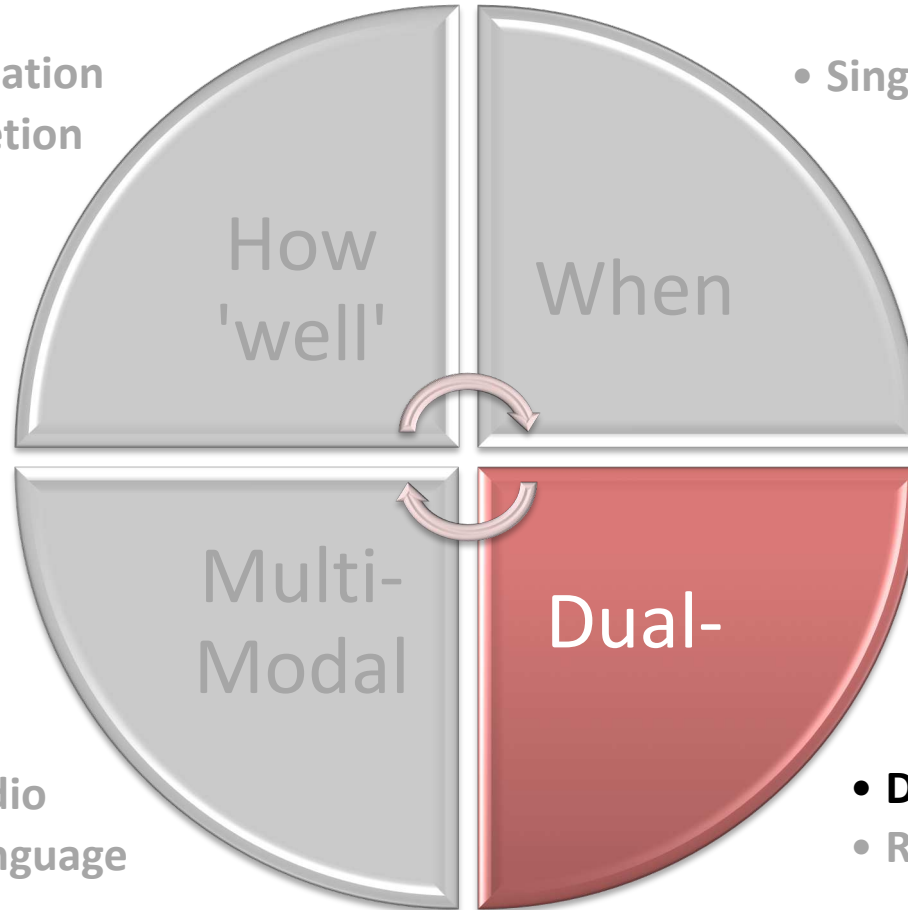
- DDLSTM
- Retro-Actions



# Fine-Grained Object Interactions

---

- Skill Determination
- Action Completion



- Single-timestamp

- Vision+Audio
- Vision+Language

- **DDLSTM**
- Retro-Actions

# Dual-Domain LSTM for Cross-Dataset Action Recognition

---

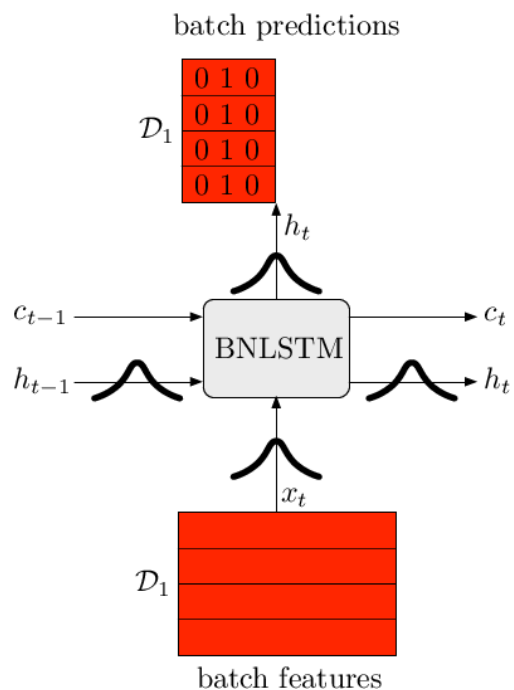
with: Toby Perrett



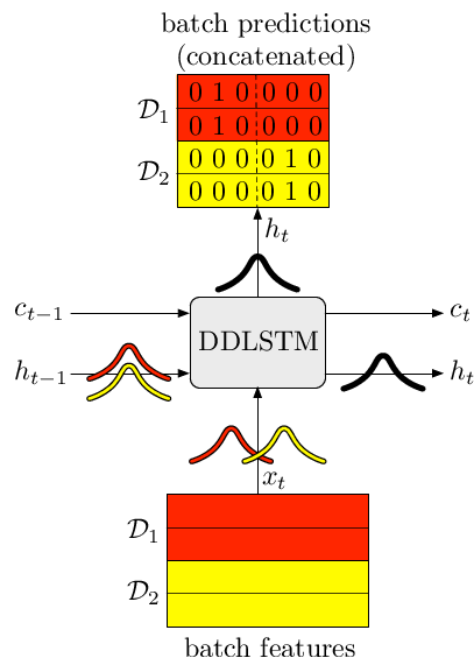
# Dual-Domain LSTM for Cross-Dataset Action Recognition

with: Toby Perrett

## BNLSTM 1 dataset

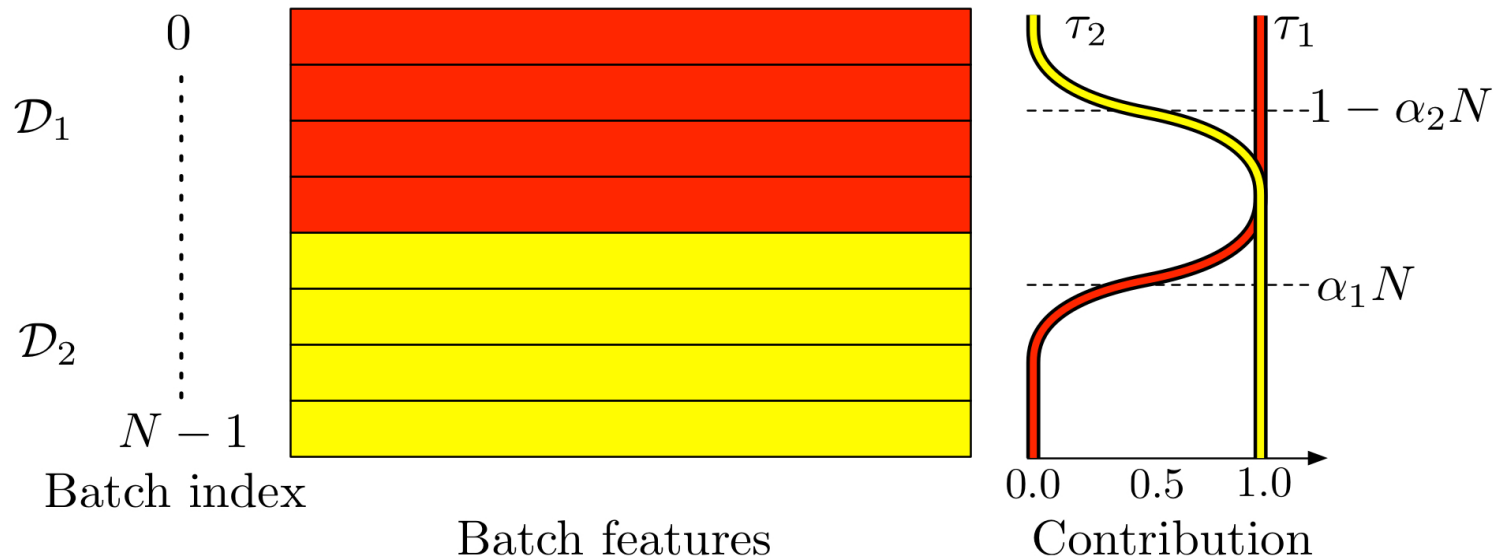


## DDLSTM 2 datasets



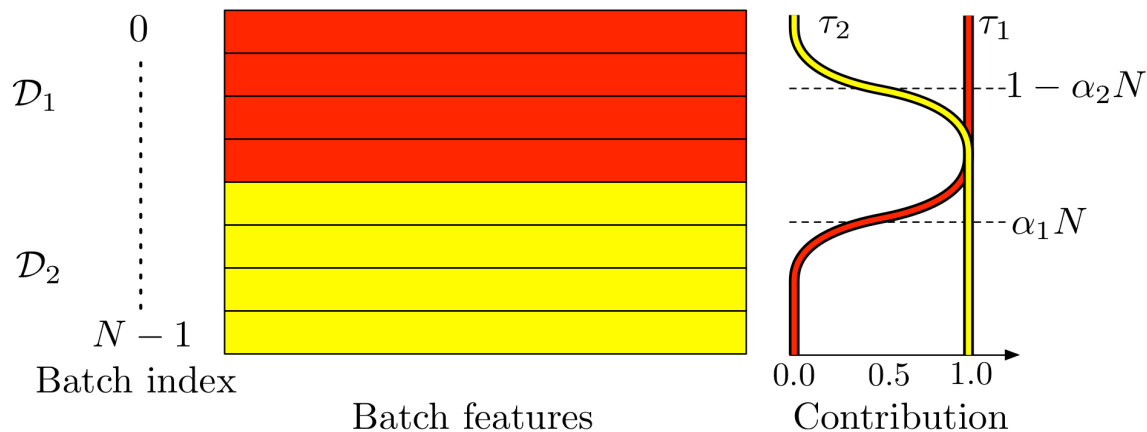
# Dual-Domain LSTM for Cross-Dataset Action Recognition

with: Toby Perrett



# Dual-Domain LSTM for Cross-Dataset Action Recognition

with: Toby Perrett

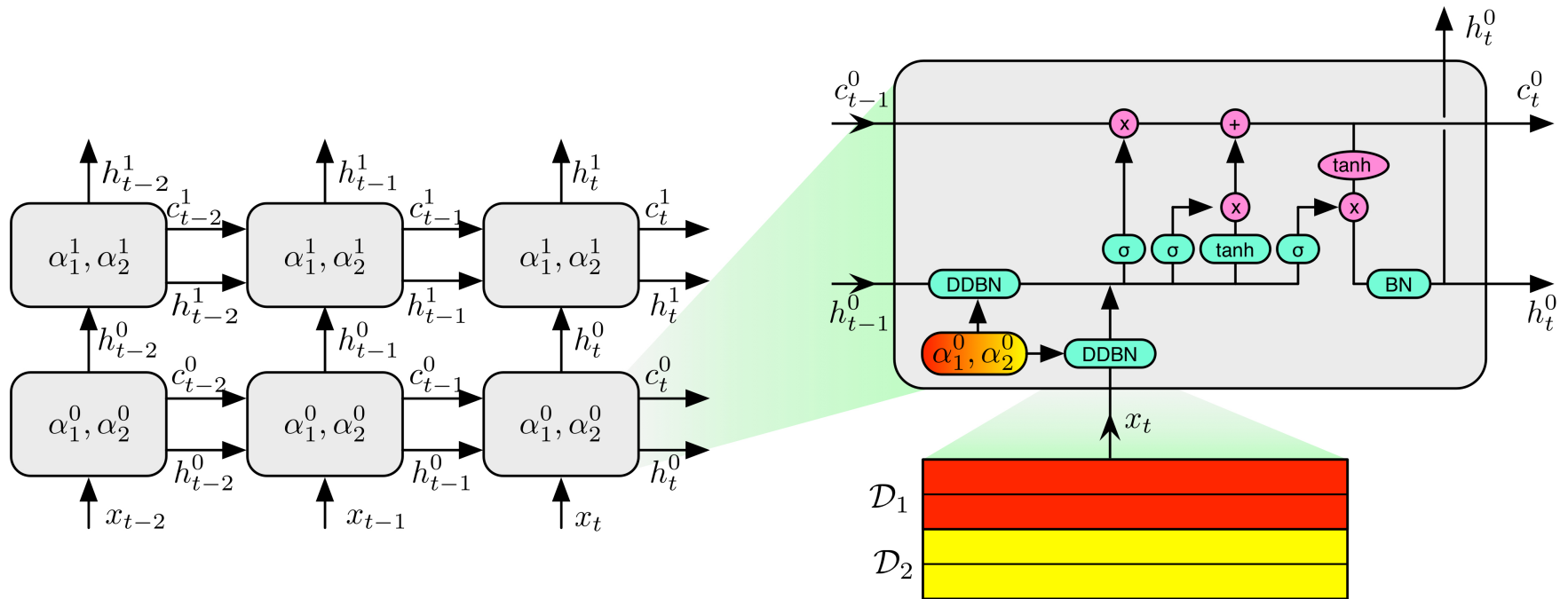


$$\tau_1(\alpha_1, j) = \frac{1 - \tanh(j - \alpha_1 N)}{2}$$

$$\tau_2(\alpha_2, j) = \frac{1 + \tanh(j - \alpha_2 N)}{2}$$

# Dual-Domain LSTM for Cross-Dataset Action Recognition

with: Toby Perrett

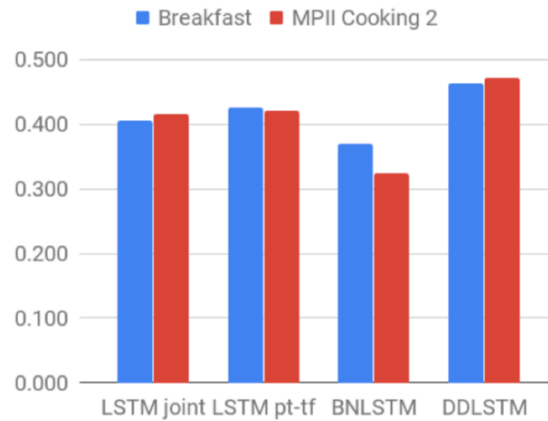


# Dual-Domain LSTM for Cross-Dataset Action Recognition

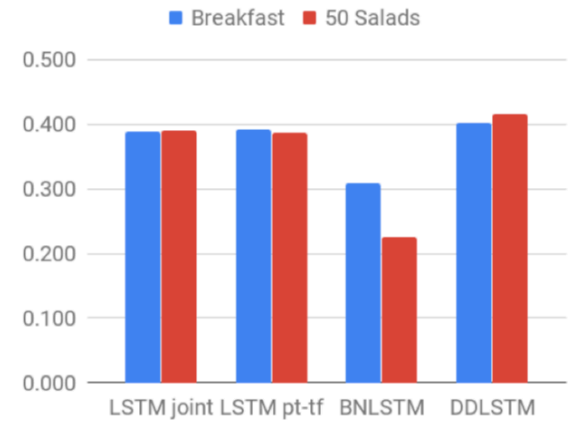
with: Toby Perrett



(a) Breakfast



(b) 50 Salads



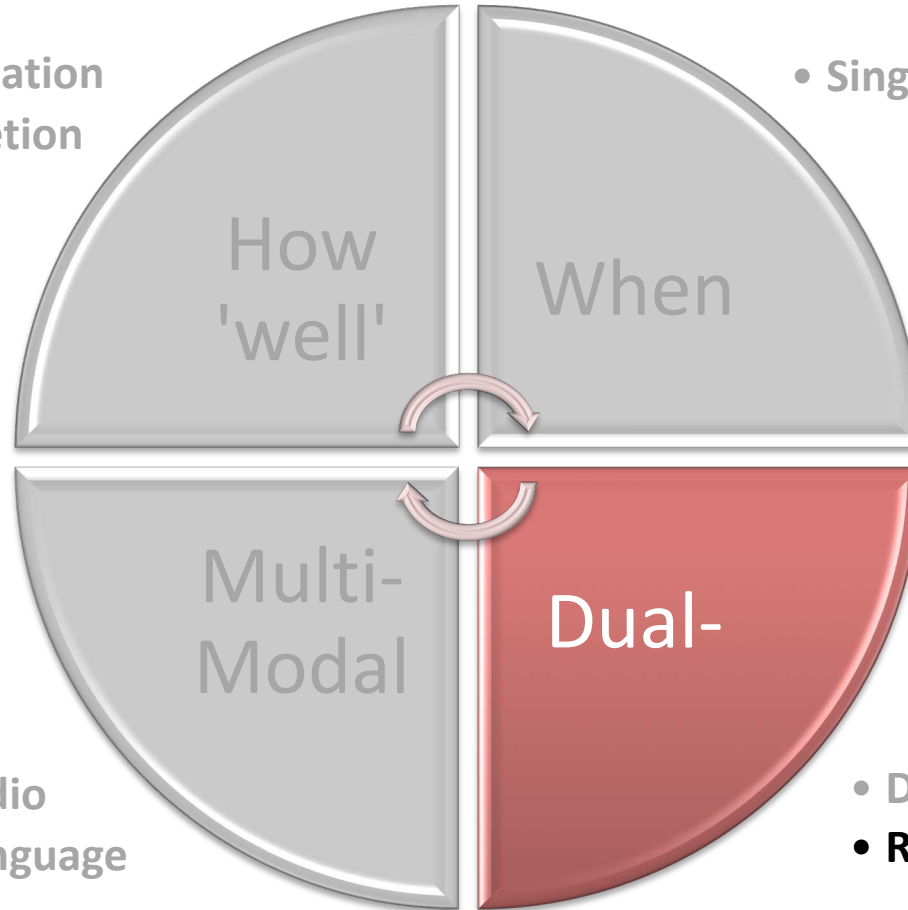
(c) MPII Cooking 2



# Fine-Grained Object Interactions

---

- Skill Determination
- Action Completion



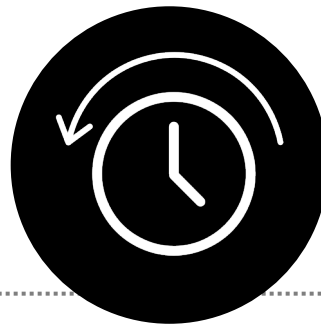
- Single-timestamp

- Vision+Audio
- Vision+Language




- DDLSTM
- **Retro-Actions**

# Retro-actions

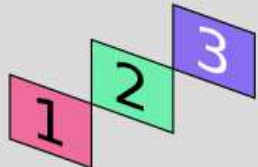





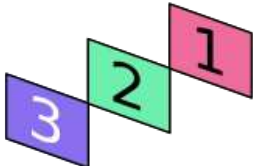


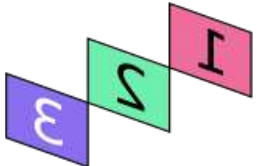


---



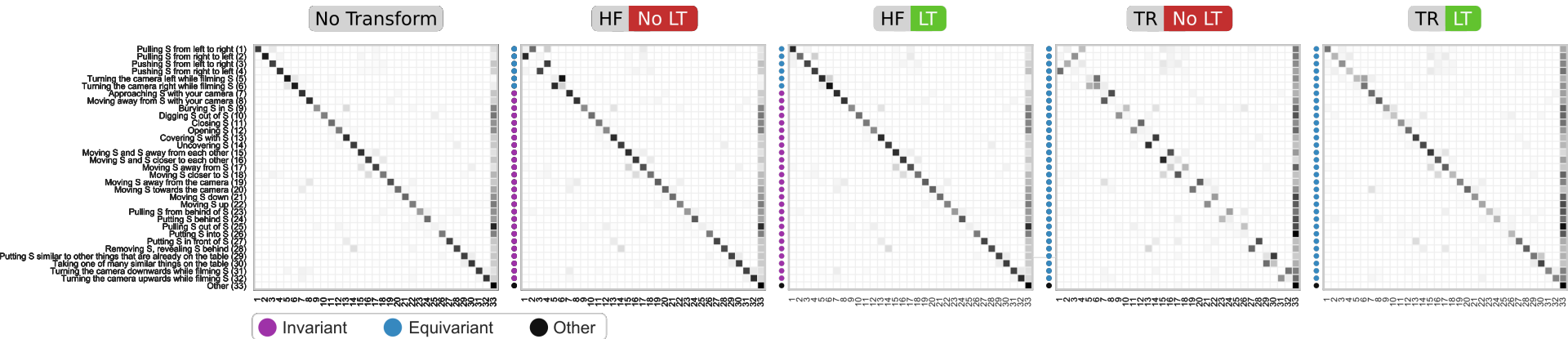
# Retro-actions

INVARIANT	<p>moving [part] of [something]</p>  <p>moving [part] of [something]</p>
EQUIVARIANT	<p>removing [something], revealing [something] behind</p>  <p>putting [something] in front of [something]</p>
IRREVERSIBLE	<p>poking a stack of [something] so the stack collapses</p>  <p>irreversible</p>

# Retro-actions

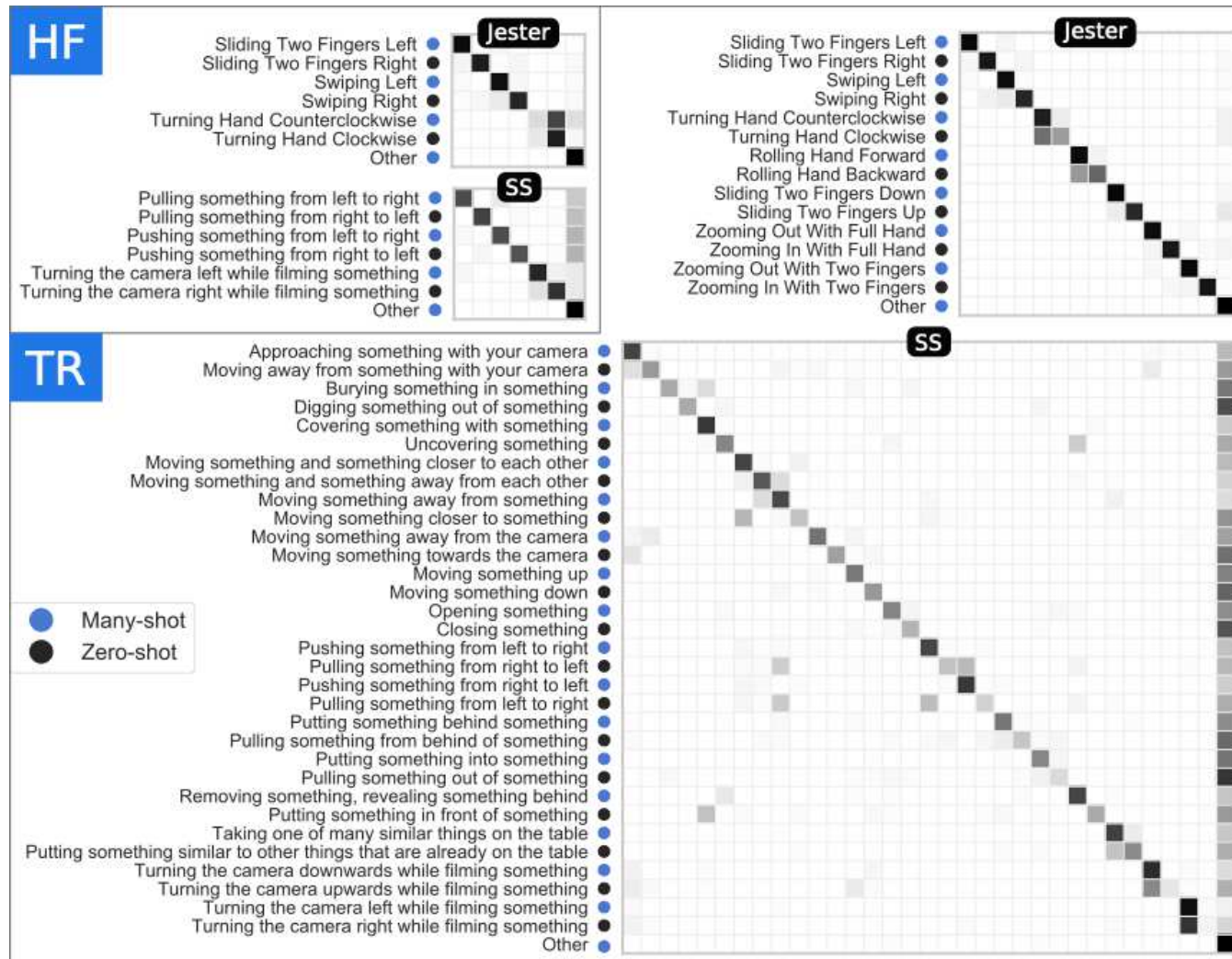
ORIG		<p>Opening</p> 	<p>Pulling Left to Right</p> 
HF		<p>Opening</p> 	<p>Pulling <b>Right to Left</b></p> 
TR		<p>Closing</p> 	<p>Pushing <b>Right to Left</b></p> 
HF+TR		<p>Closing</p> 	<p>Pushing <b>Left to Right</b></p> 

# Retro-actions





# Retro-actions – Zero-Shot Learning

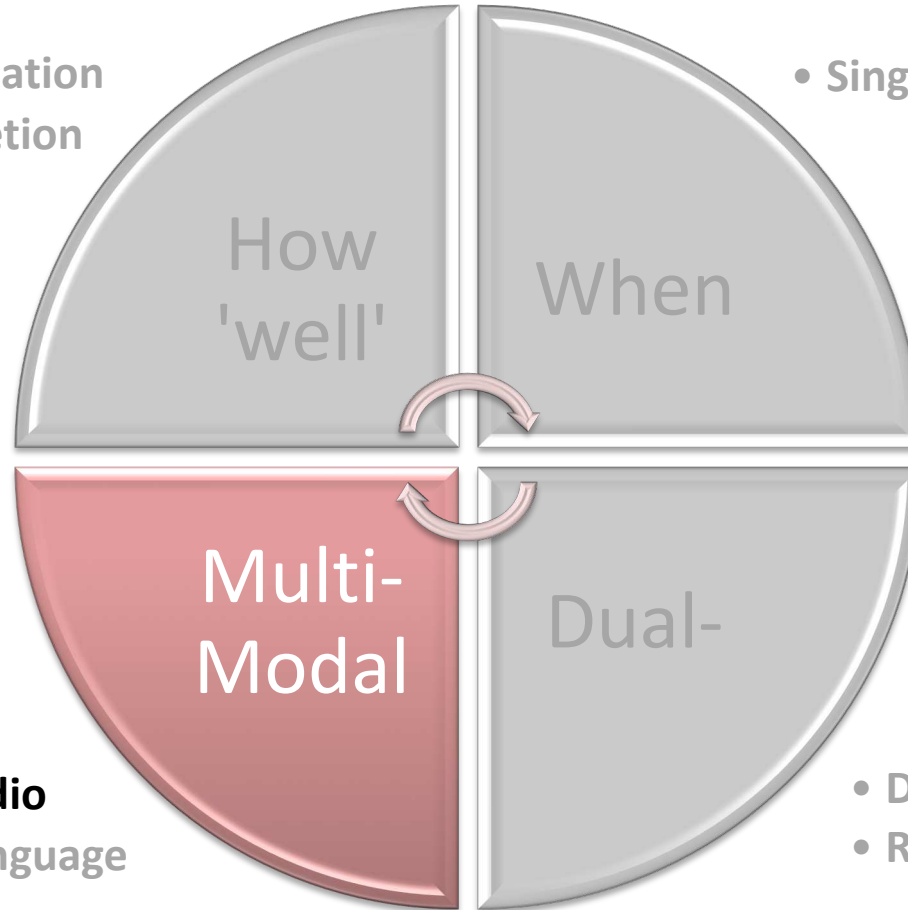




# Fine-Grained Object Interactions

---

- Skill Determination
- Action Completion



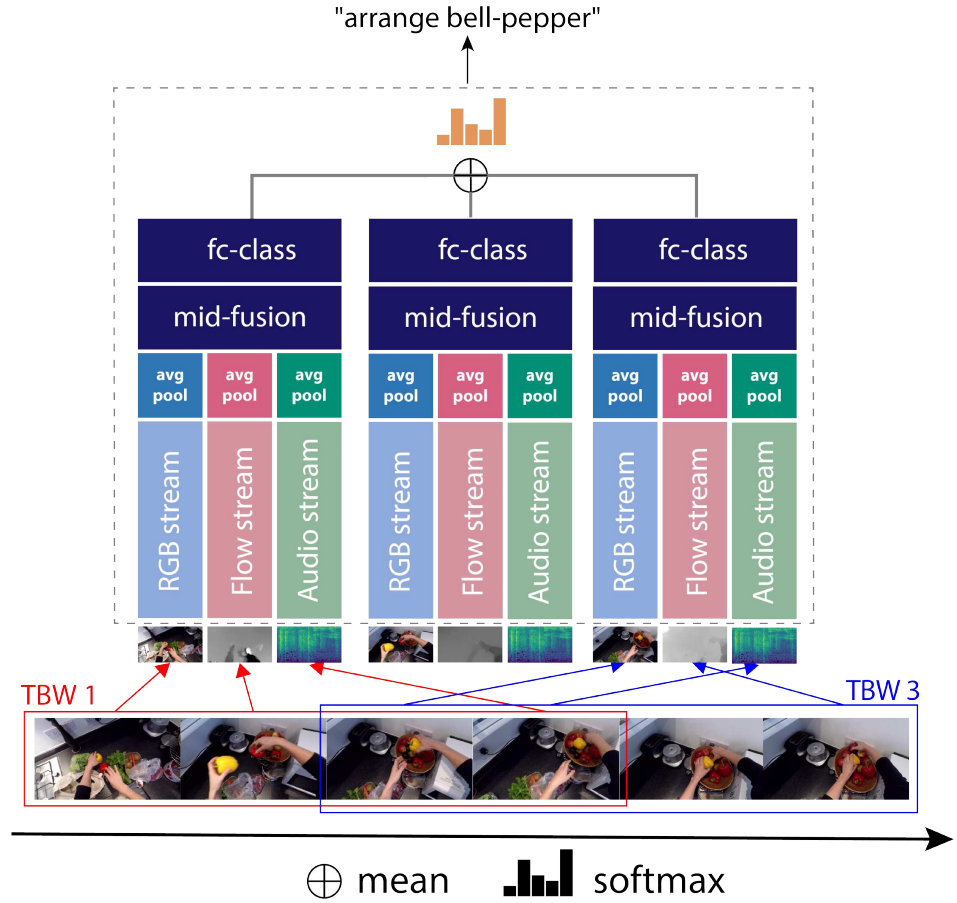
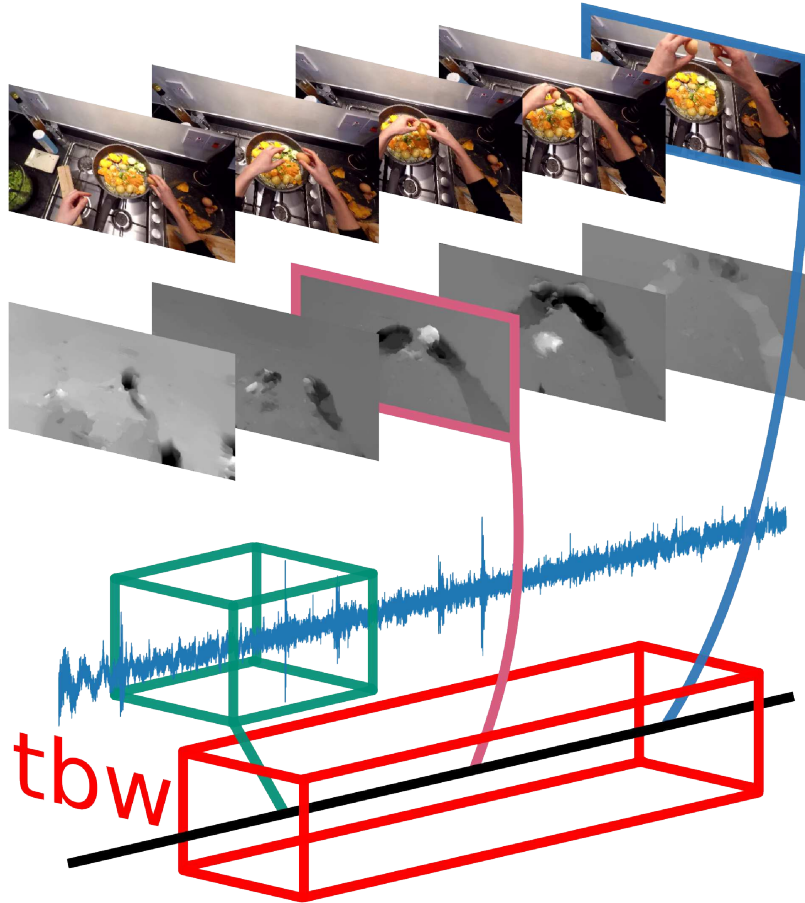
- Single-timestamp

- **Vision+Audio**
- Vision+Language

- DDLSTM
- Retro-Actions

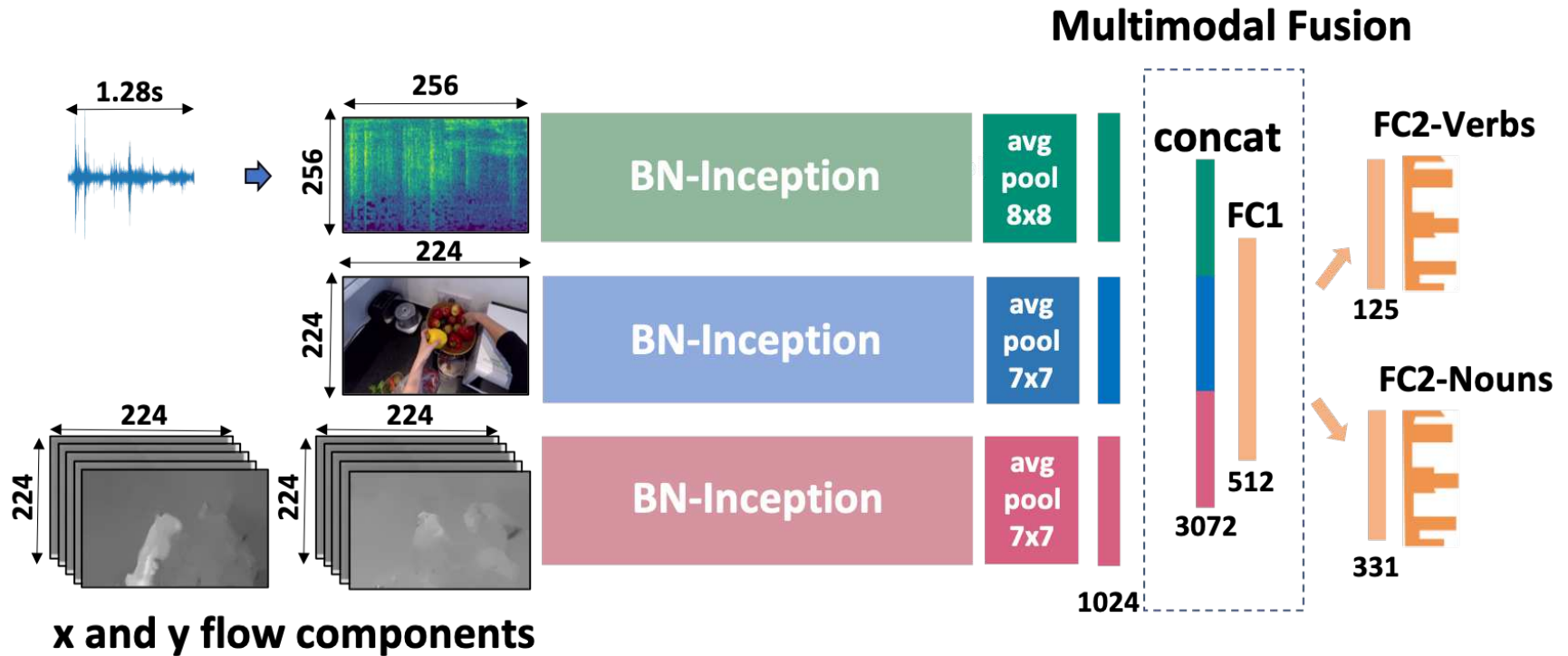
# Audio-Visual Temporal Binding for Egocentric Action Recognition

with: Vangelis Kazakos  
Arsha Nagrani  
Andrew Zisserman



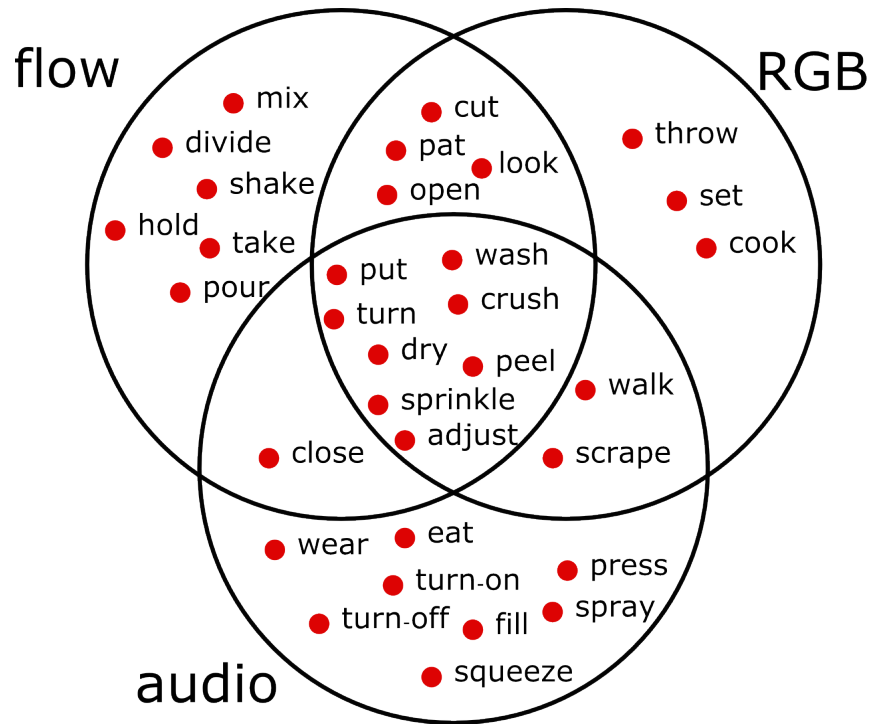
# Audio-Visual Temporal Binding for Egocentric Action Recognition

with: Vangelis Kazakos  
Arsha Nagrani  
Andrew Zisserman



# Audio-Visual Temporal Binding for Egocentric Action Recognition

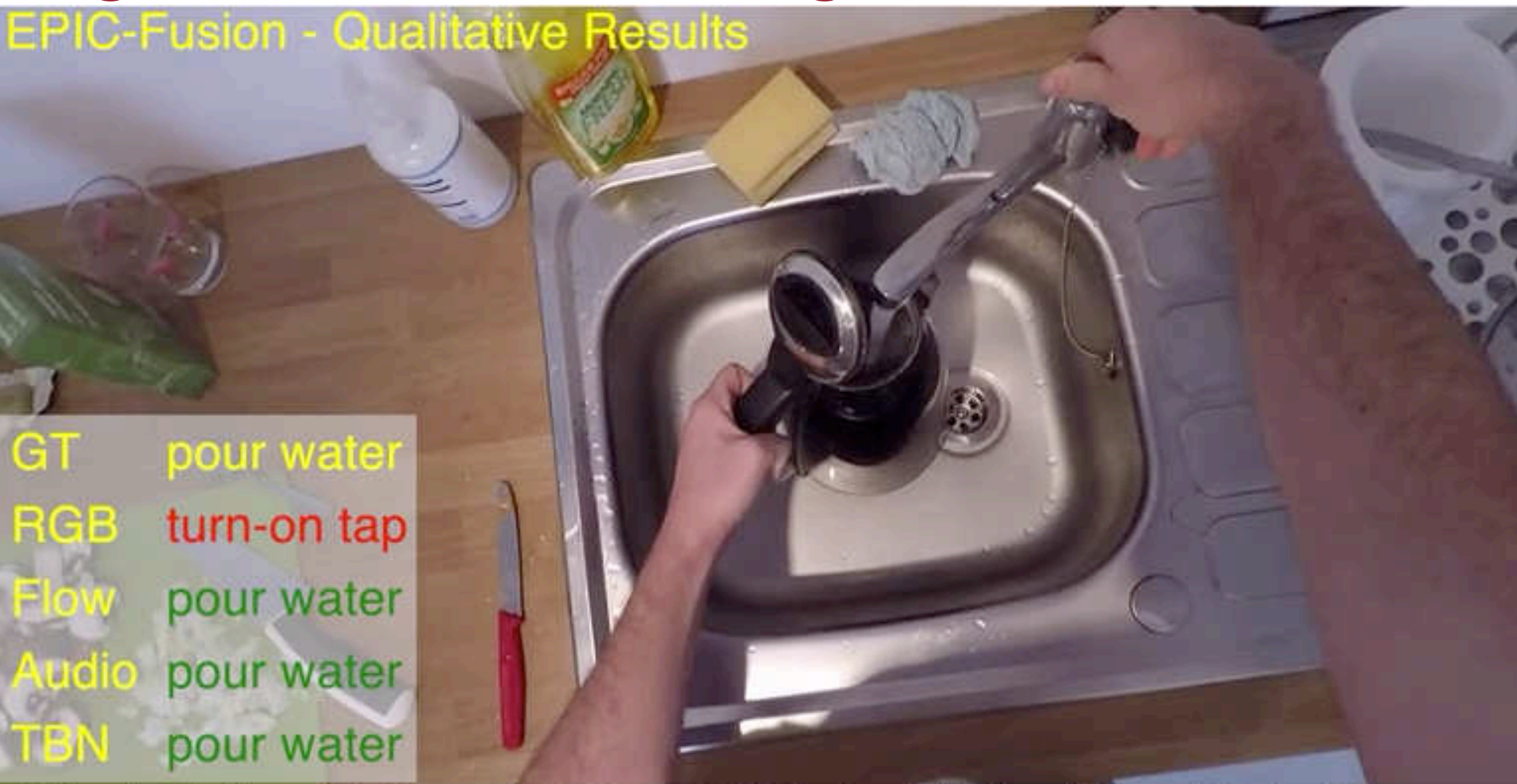
with: Vangelis Kazakos  
Arsha Nagrani  
Andrew Zisserman



# Audio-Visual Temporal Binding for Egocentric Action Recognition

with: Vangelis Kazakos  
Arsha Nagrani  
Andrew Zisserman

## EPIC-Fusion - Qualitative Results

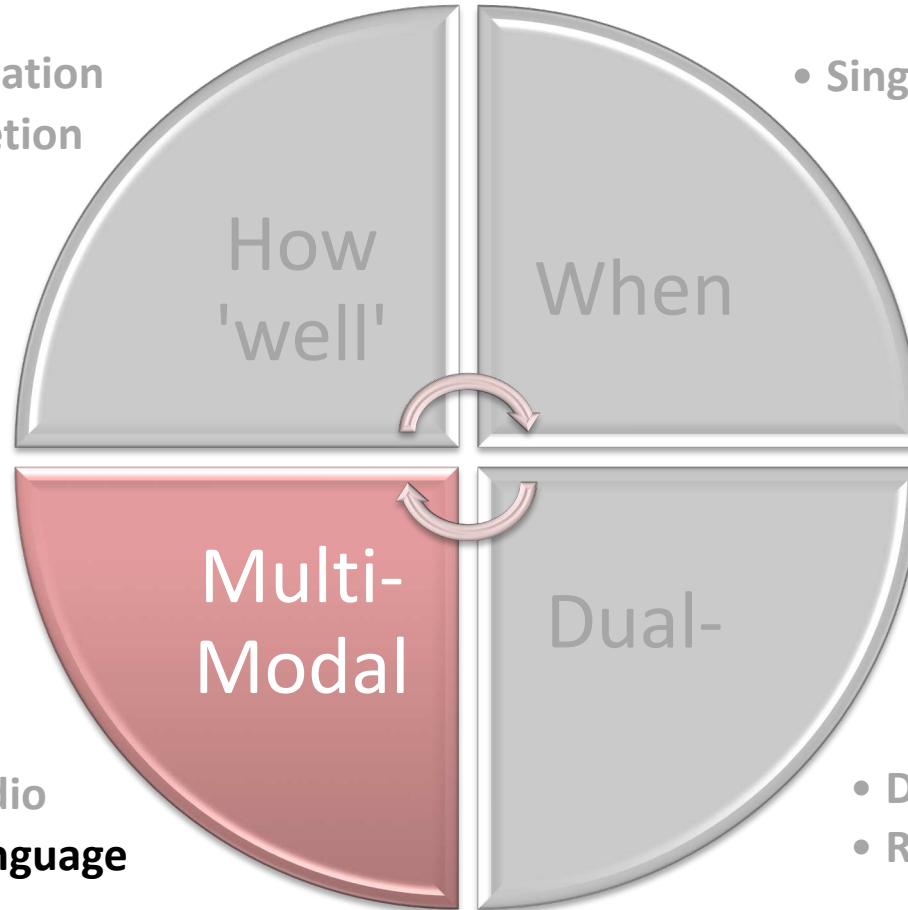


E. Kazakos, A. Nagrani, A. Zisserman, D. Damen, EPIC-Fusion: Audio-Visual Temporal Binding for Egocentric Action Recognition, ICCV 2019

# Fine-Grained Object Interactions

---

- Skill Determination
- Action Completion



- Single-timestamp

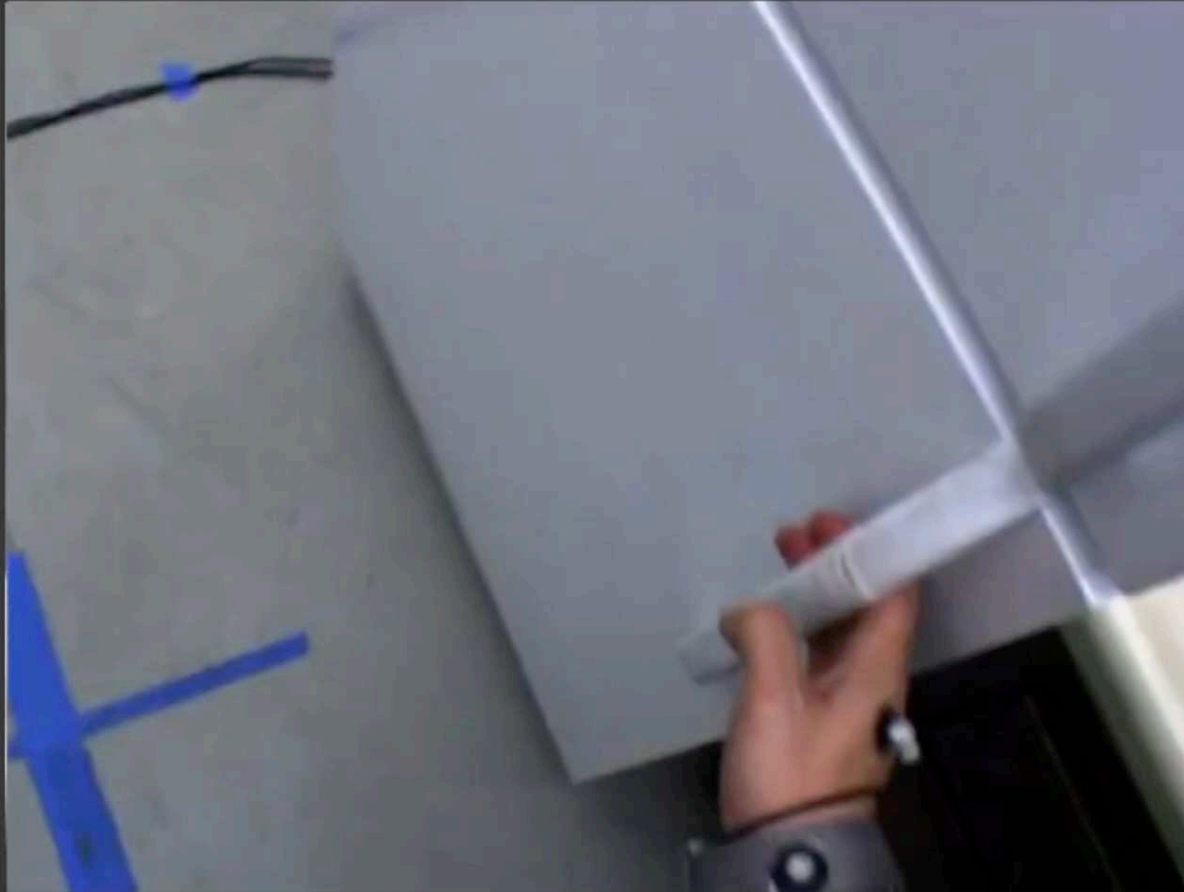
- Vision+Audio
- **Vision+Language**

- DDLSTM
- Retro-Actions



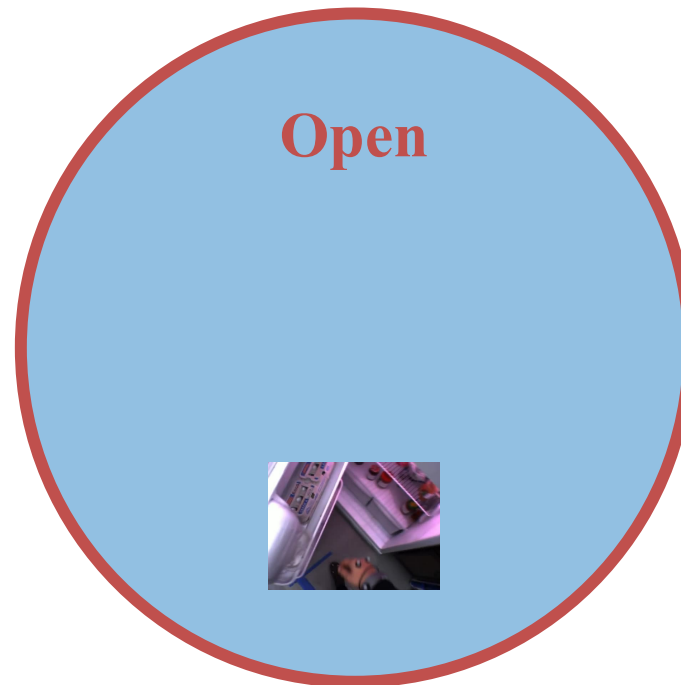
# The *Verbs* Dilemma

---



# The *Verbs* Dilemma

---



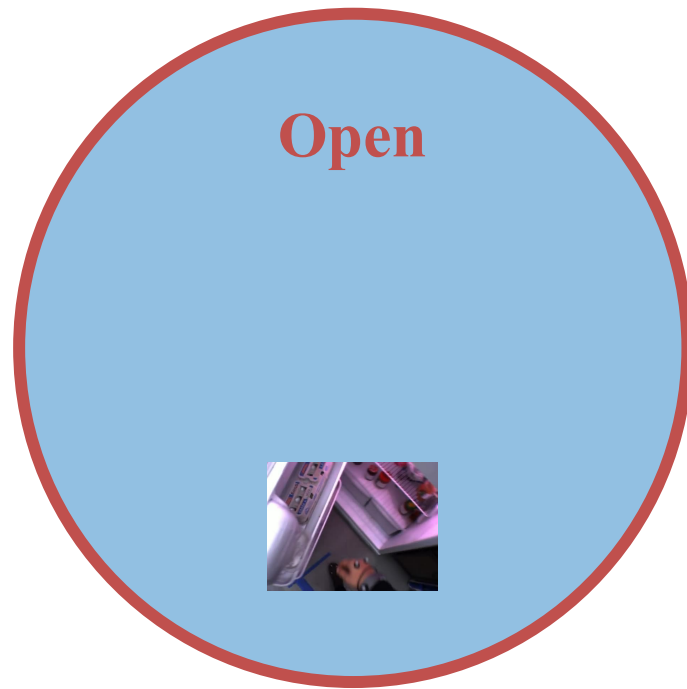
# The Verbs Dilemma

---



# The *Verbs* Dilemma

---



# The Verbs Dilemma

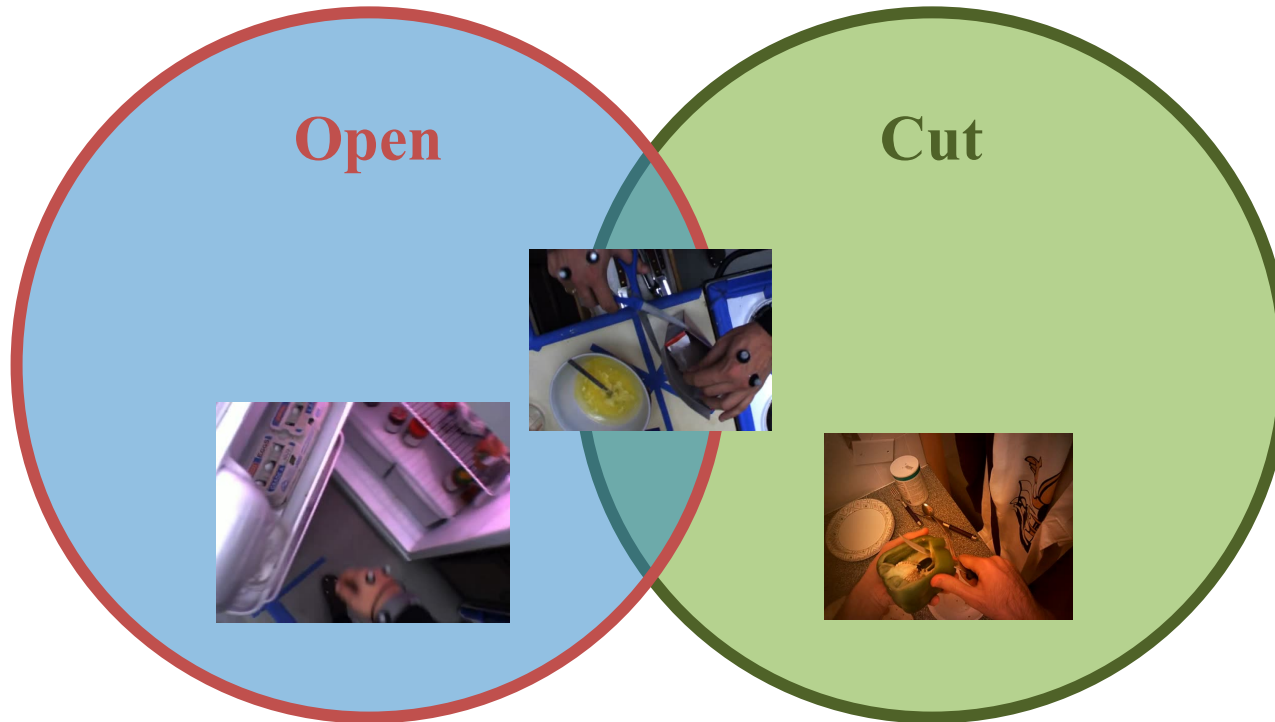
---





# The *Verbs* Dilemma

---



# The *Verbs* Dilemma

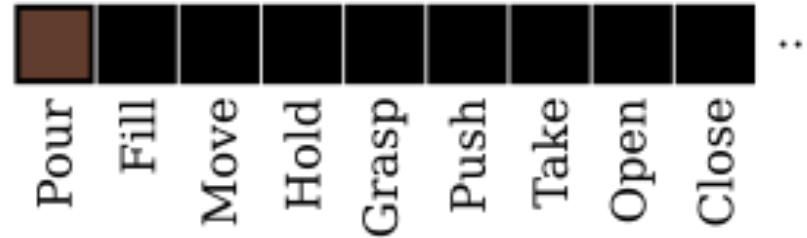
---

- 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

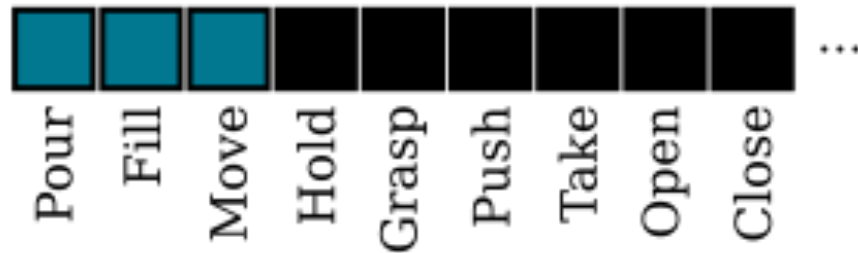
# The Verbs Dilemma



## Single Verb



## Multi Verb



## Soft Assigned Multi Verb



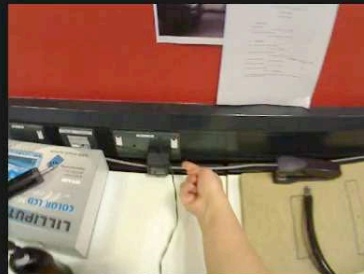
# The Verbs Dilemma

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.

# Fine-Grained Action Retrieval

with: Michael Wray  
Gabriela Csurka  
Diane Larlus

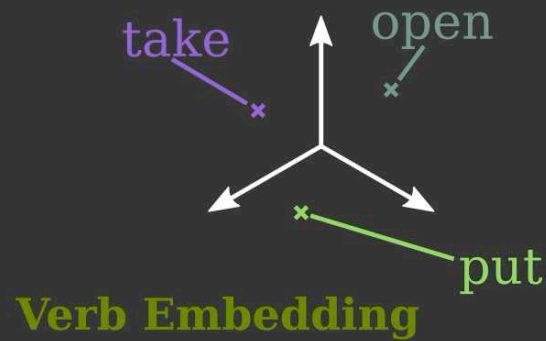
In this work we focus on  
**Fine-Grained Action Retrieval**

I put meat on a  
ball of dough



# Fine-Grained Action Retrieval

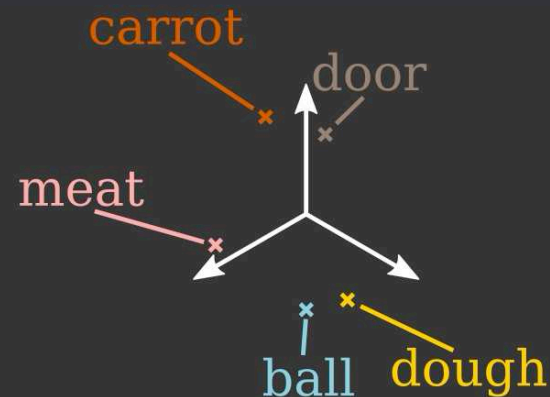
We embed the video  
and representations



**Verb Embedding**

[put]

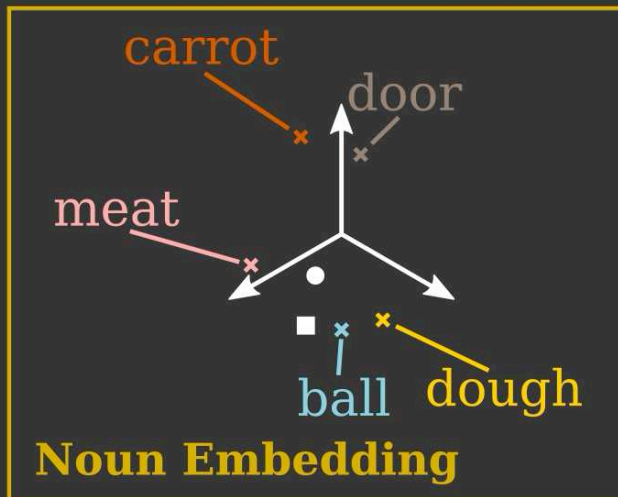
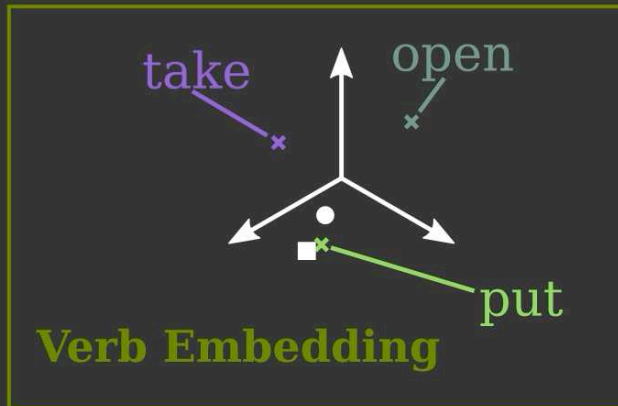
[meat, ball, dough]



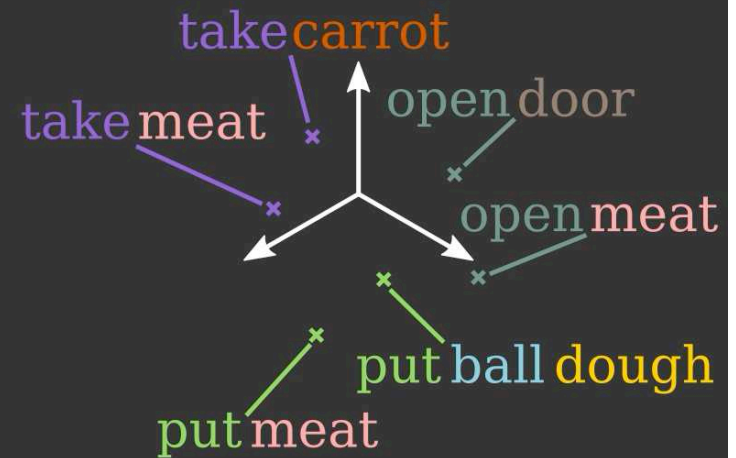
**Noun Embedding**



# Fine-Grained Action Retrieval



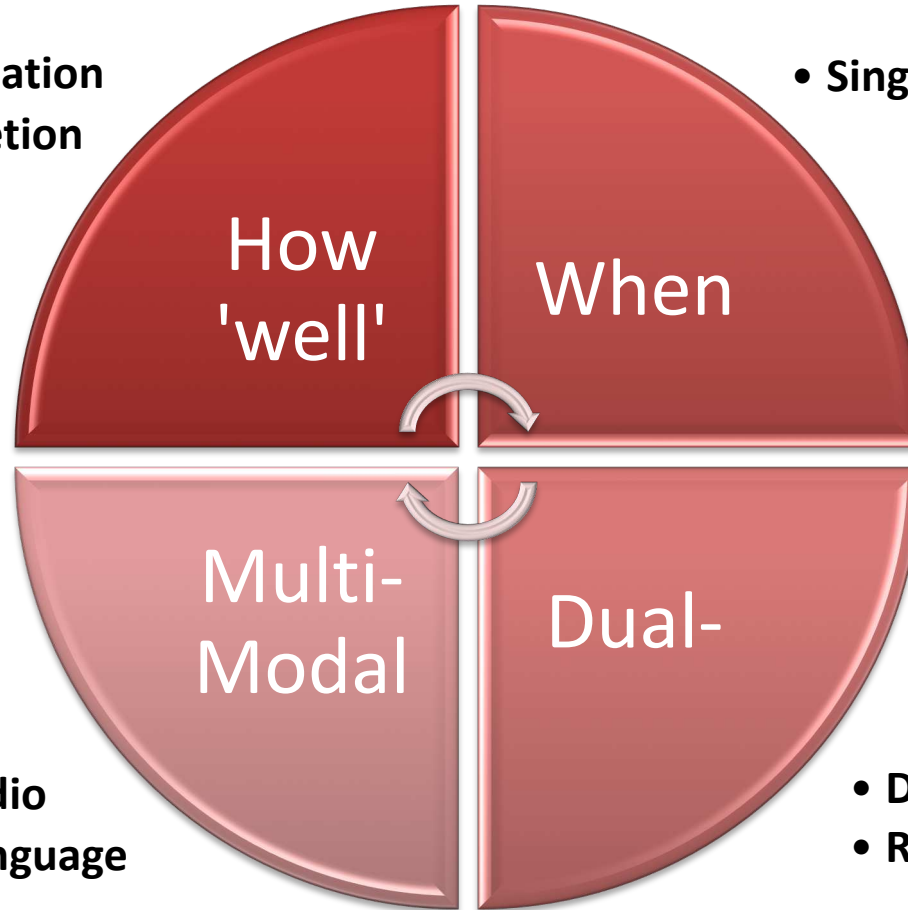
Finally, we combine the outputs and embed these into an action space



# Fine-Grained Object Interactions

---

- Skill Determination
- Action Completion



- Single-timestamp

- Vision+Audio
- Vision+Language

- DDLSTM
- Retro-Actions

# Bristol and University of Bristol

---



# Thank you...

---

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

<http://dimadamen.github.io>



@dimadamen



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



# Scaling Egocentric Vision: The **EPIC-KITCHENS** Dataset



Dima Damen



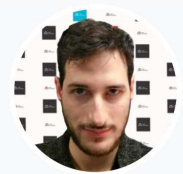
Hazel Doughty



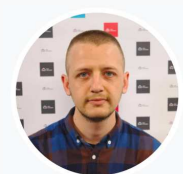
Giovanni M. Farinella



Sanja Fidler



Antonino Furnari



Evangelos Kazakos



Davide Moltisanti



Jonathan Munro



Toby Perrett



Will Price



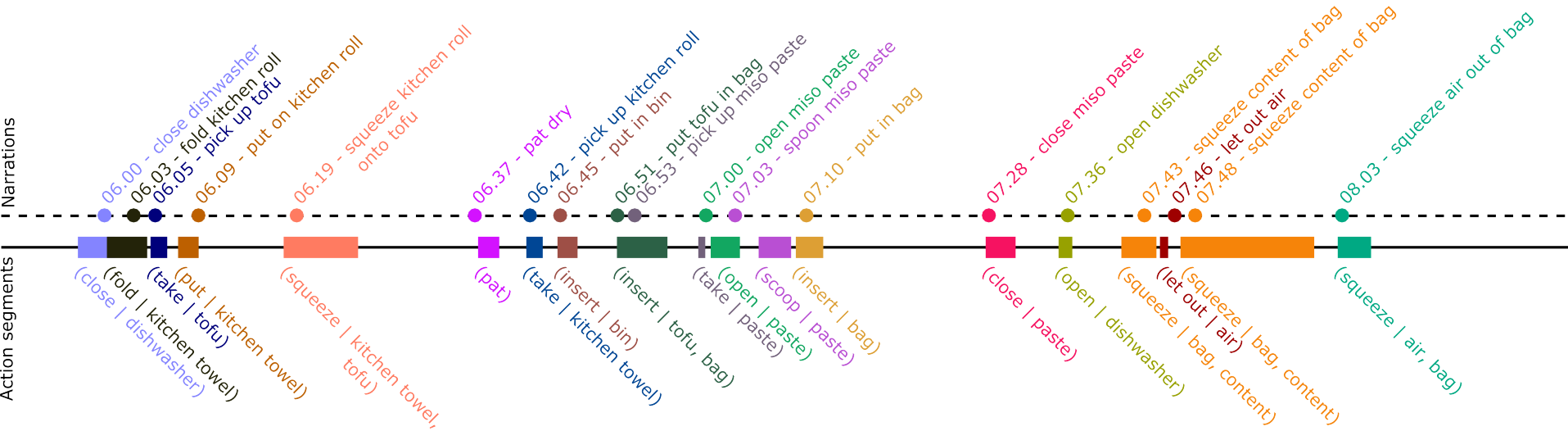
Michael Wray

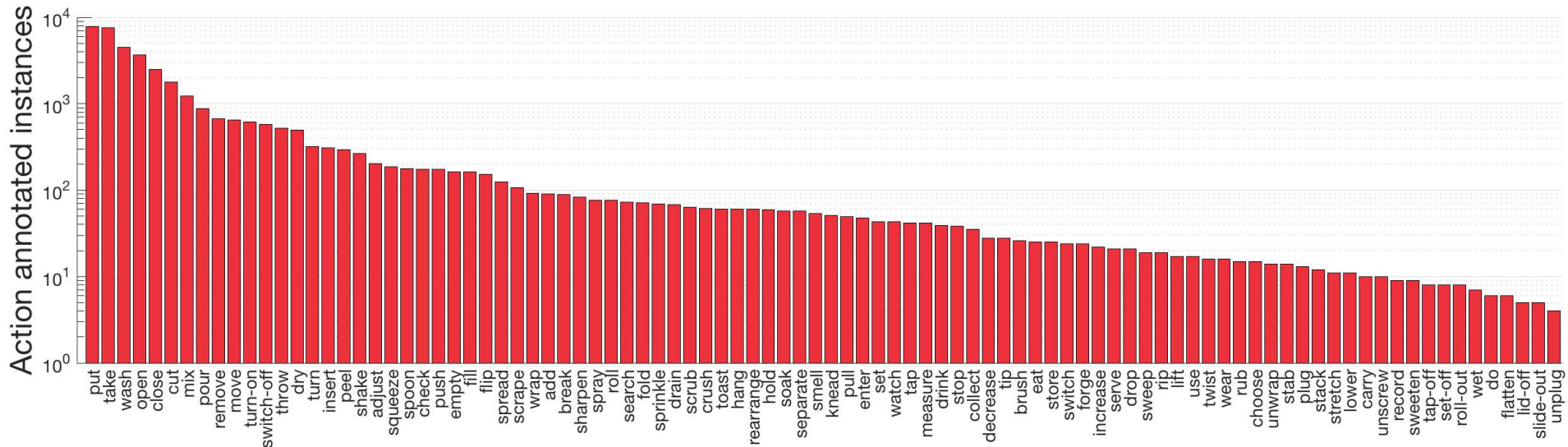


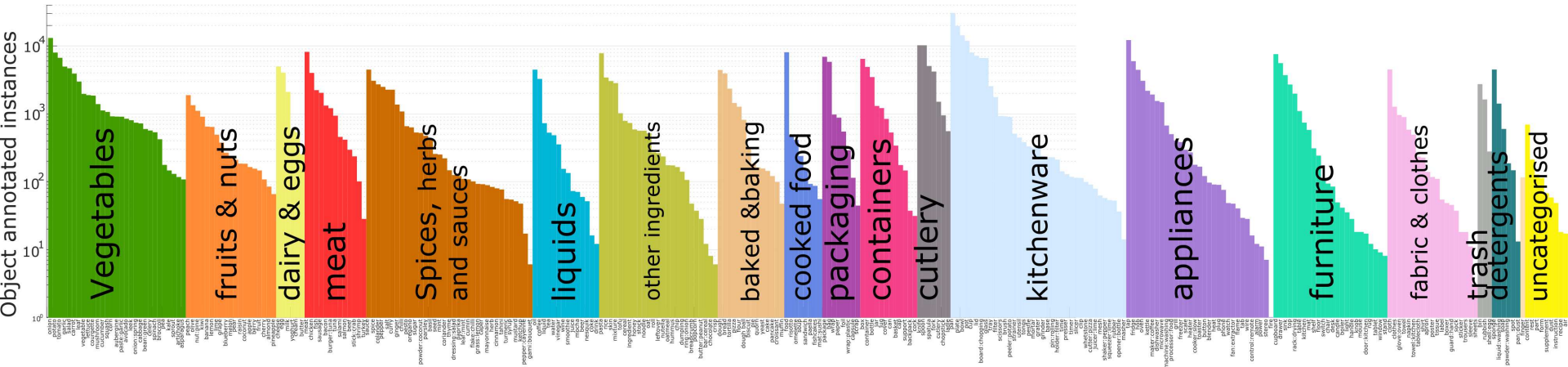




# Narrations to Action Segments









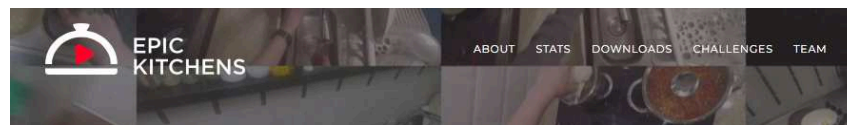
39 000  
ACTION SEGMENTS



454 200  
OBJECT ANNOTATIONS



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



## NEWS

- EPIC-KITCHENS accepted for oral presentation at ECCV 2018 in Munich this September
- News coverage: [UoB](#), [The Spoon](#), [Il Sole 24 Ore](#), [La Sicilia](#), [Elpais](#)
- EPIC-Kitchens Released: 9th of April 2018!!!
- Watch [YouTube Release Trailer here](#)

### What is EPIC-Kitchens?

The largest dataset in first-person (egocentric) vision; multi-faceted non-scripted recordings in native environments - i.e. the wearers' homes, capturing all daily activities in the kitchen over multiple days. Annotations are collected using a novel 'live' audio commentary approach.

### Characteristics

- 32 kitchens - 4 cities
- Head-mounted camera
- 55 hours of recording - Full HD, 60fps
- 11.5M frames
- Multi-language narrations
- 39,594 action segments
- 454,158 object bounding boxes
- 125 verb classes, 352 noun classes

### Updates

Stay tuned with updates on epic-kitchens2018, as well as EPIC workshop series by joining the [epic-community mailing list](#) send an email to: [sympa@sympa.bristol.ac.uk](mailto:sympa@sympa.bristol.ac.uk) with the subject *subscribe epic-community* and a *blank* message body.

