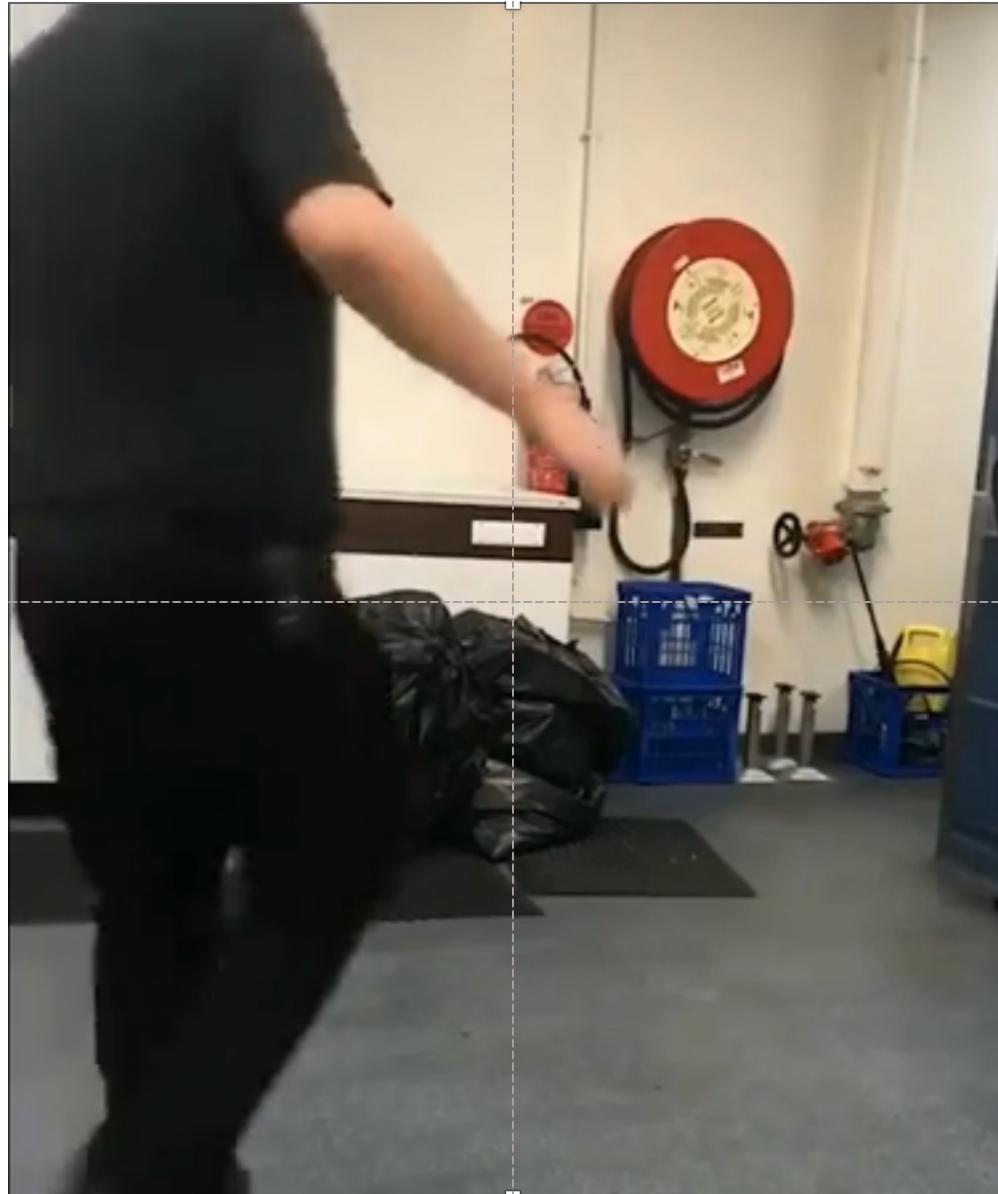


Video Understanding ❤️

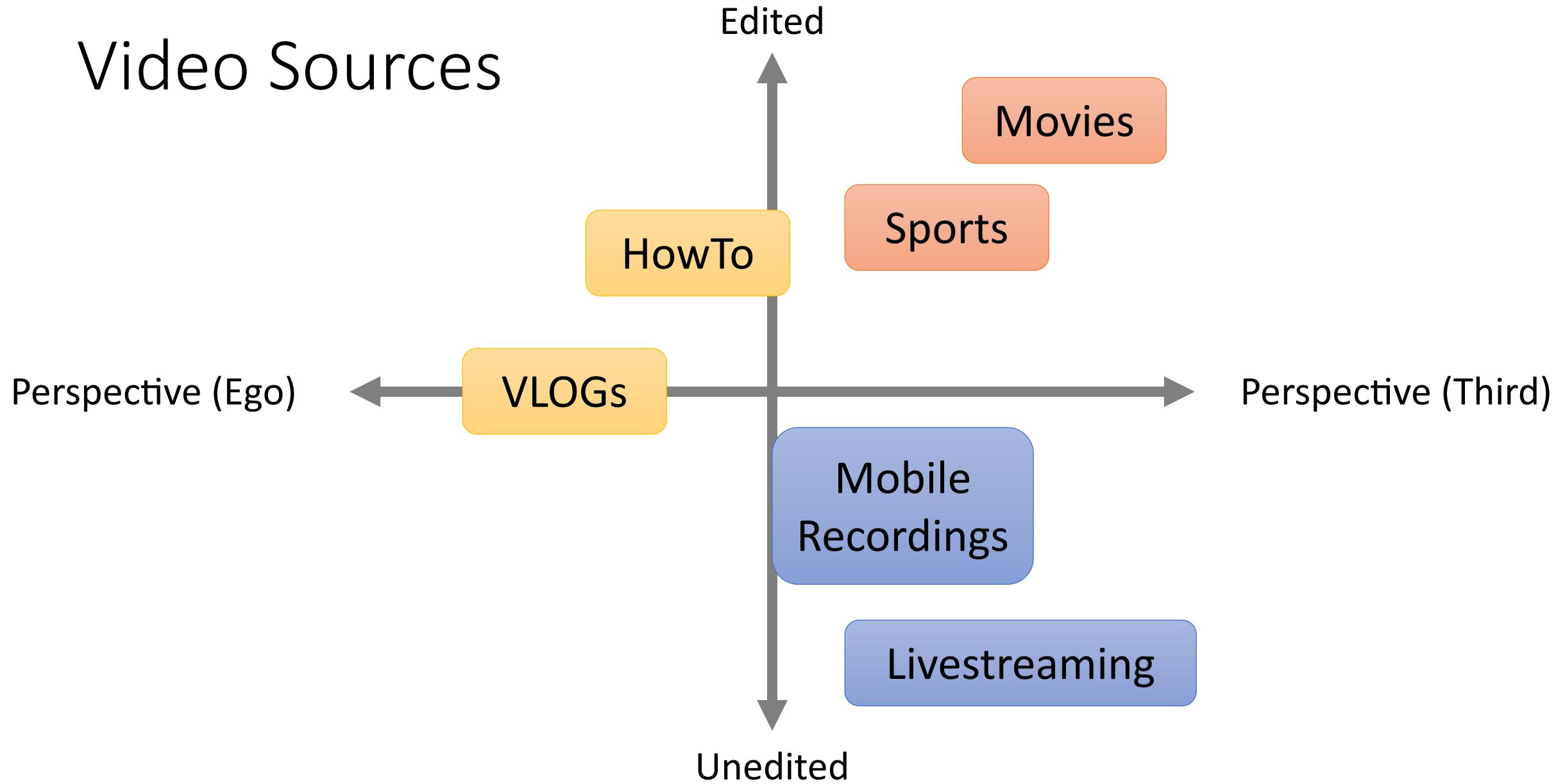
Makarand Tapaswi

CS7.505 Spring 2024

13th April 2024



Video Sources



Outline

- Part 1: Popular video understanding tasks
- Part 2: Video backbones
 - (break)
- Part 3: Video-language models
- Part 4: Open challenges

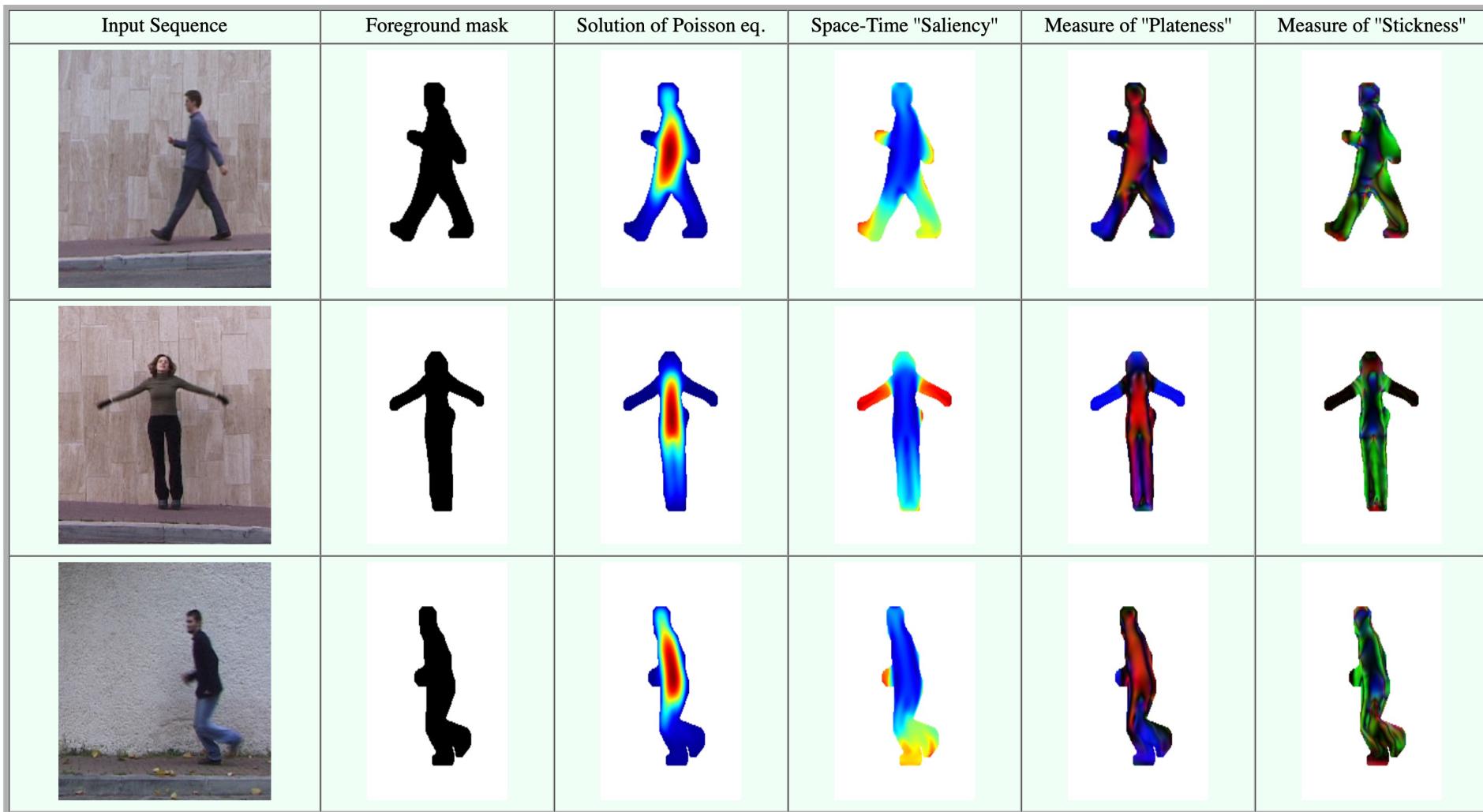
Part 1

Video Understanding

Golden years of actions: KTH Dataset

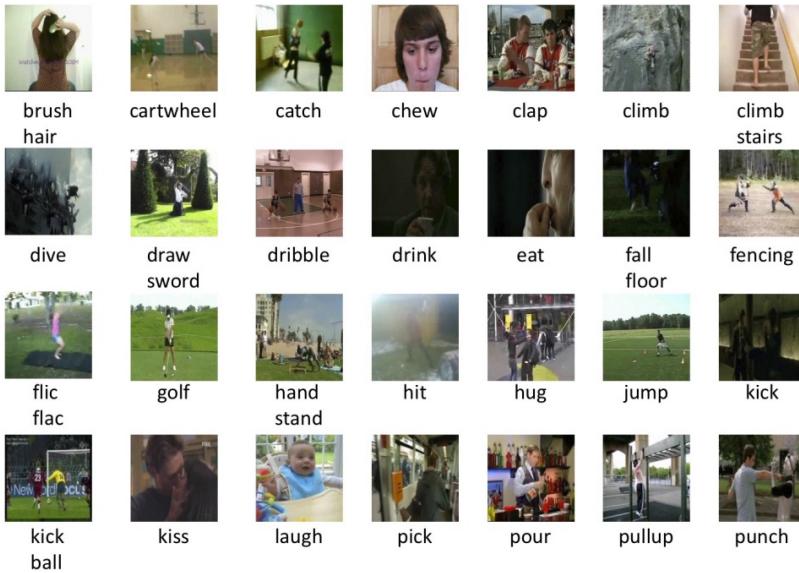


Golden years of actions: Weizmann dataset

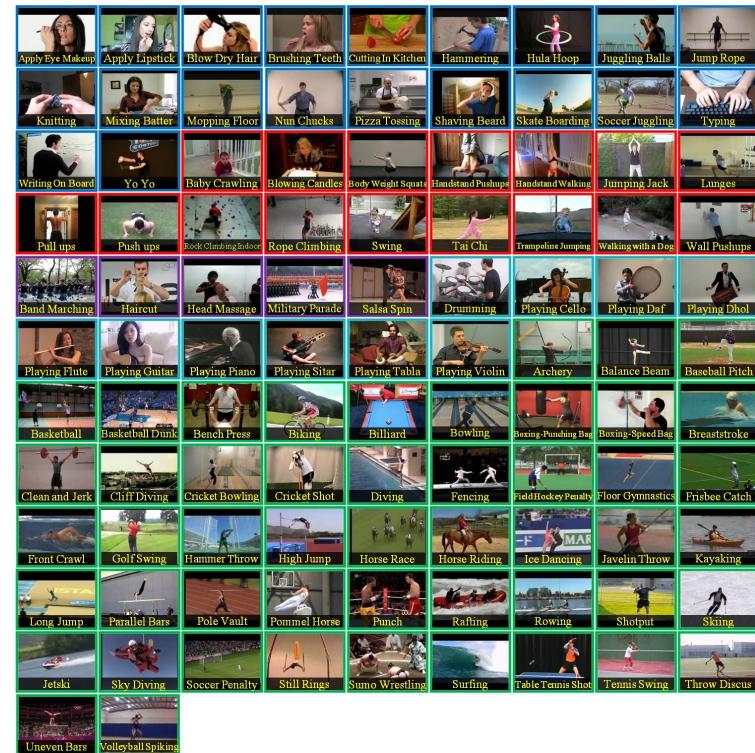


Actions datasets: sub-100 classes

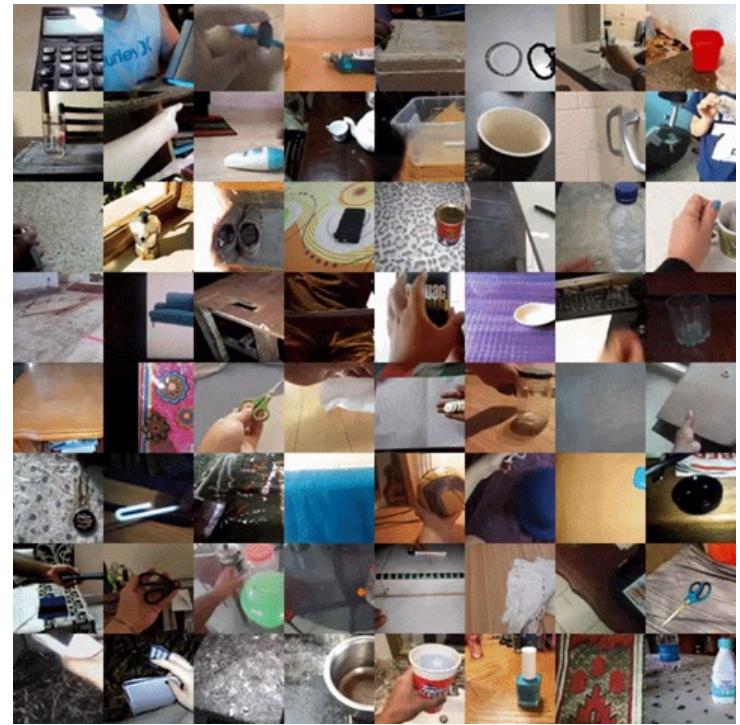
HMDB (51)



UCF (101)



Something-Sth. (174)



Kuehne, et al. **HMDB**: A Large Video Database for Human Motion Recognition. ICCV 2011

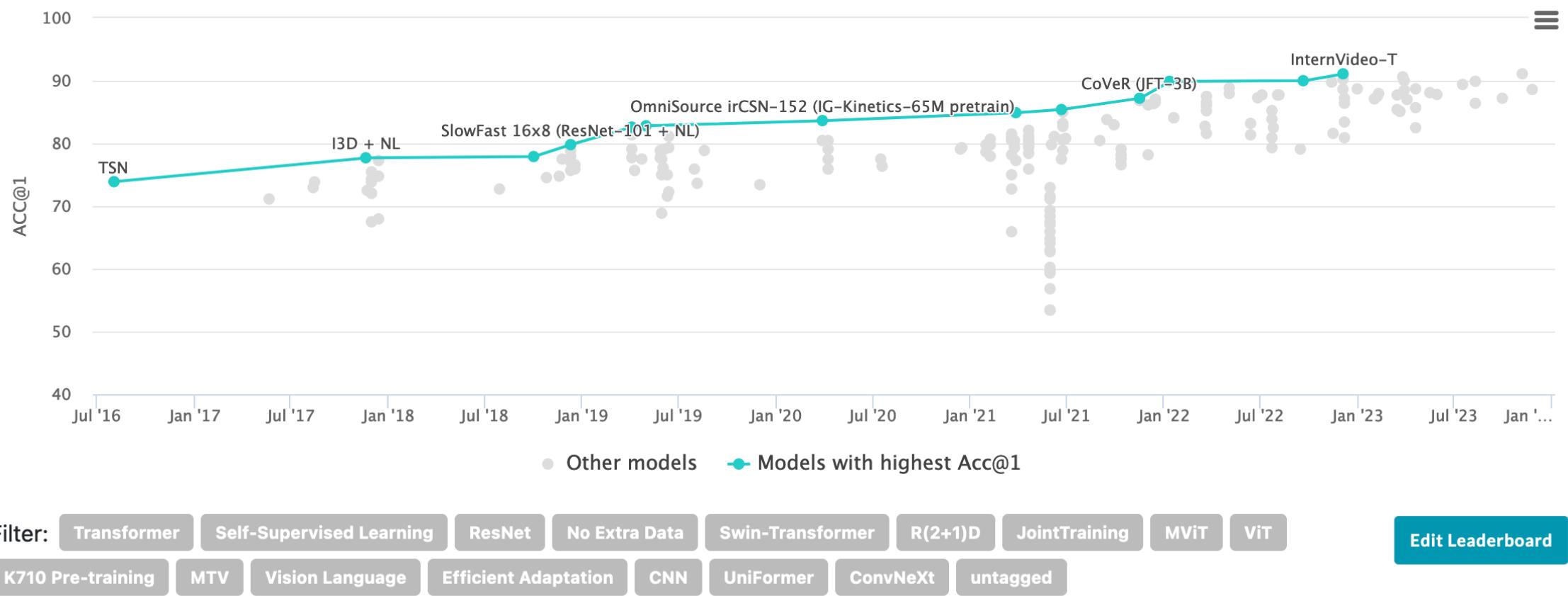
Soomro, et al. **UCF101**: A Dataset of 101 Human Action Classes from Videos in the Wild. CRCV-TR 2012

Goyal, et al. The **something something** video database for learning and evaluating visual common sense. ICCV 2017

Is there an ImageNet for videos?



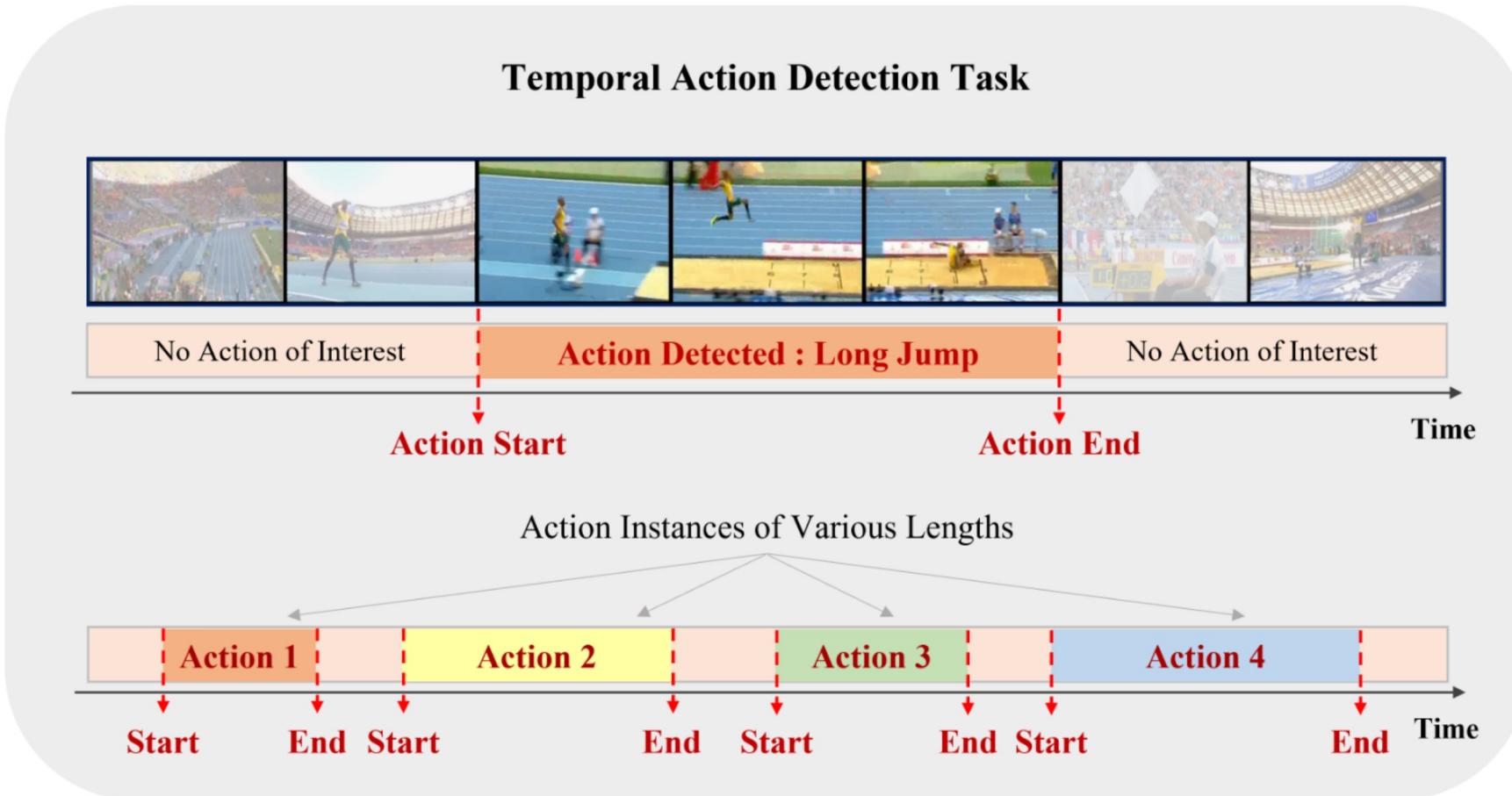
Kinetics dataset



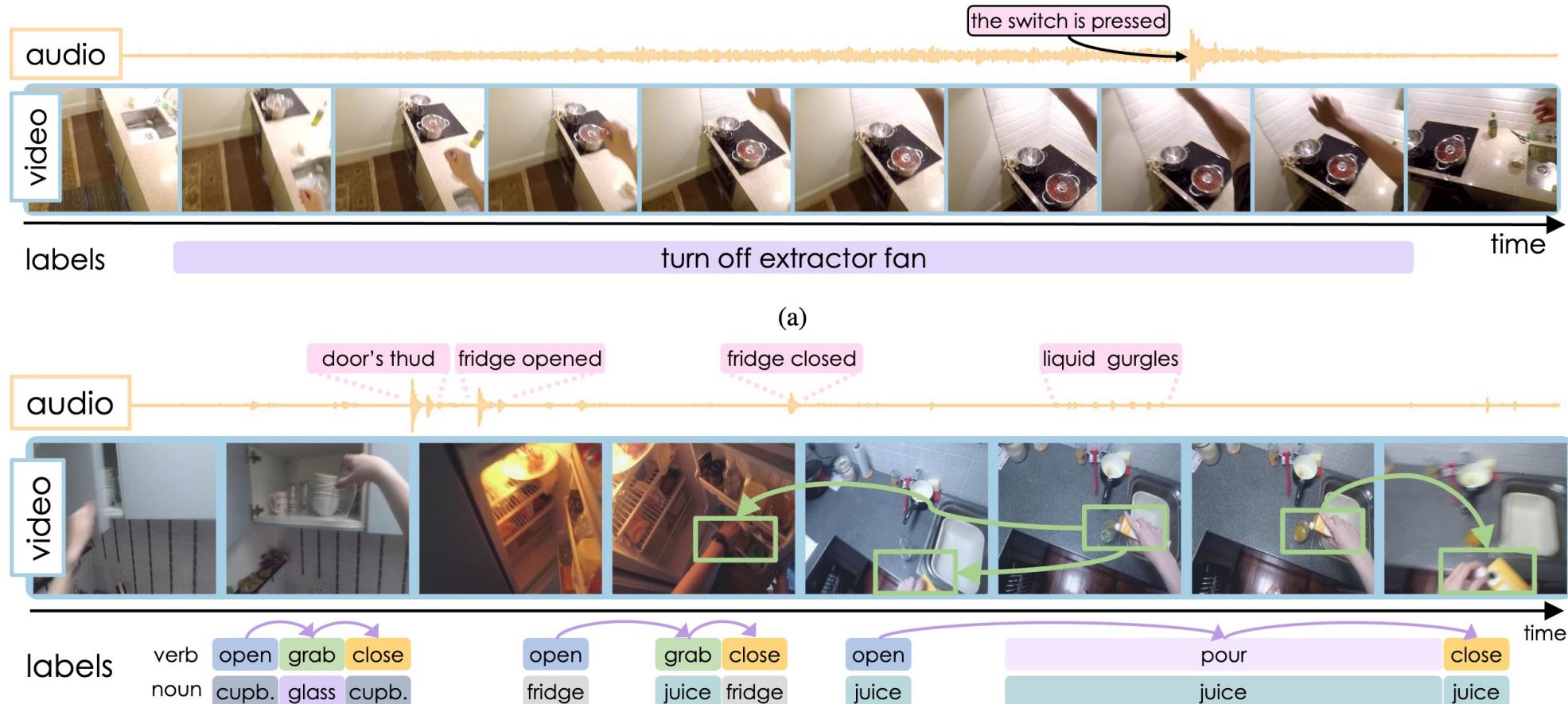
Challenges of action understanding

- who is doing the action?
- when does an action start?
- how long is the action?
- how to deal with actions vs. interactions?
- what are essential components of any action?
- what is the importance of the background scene?

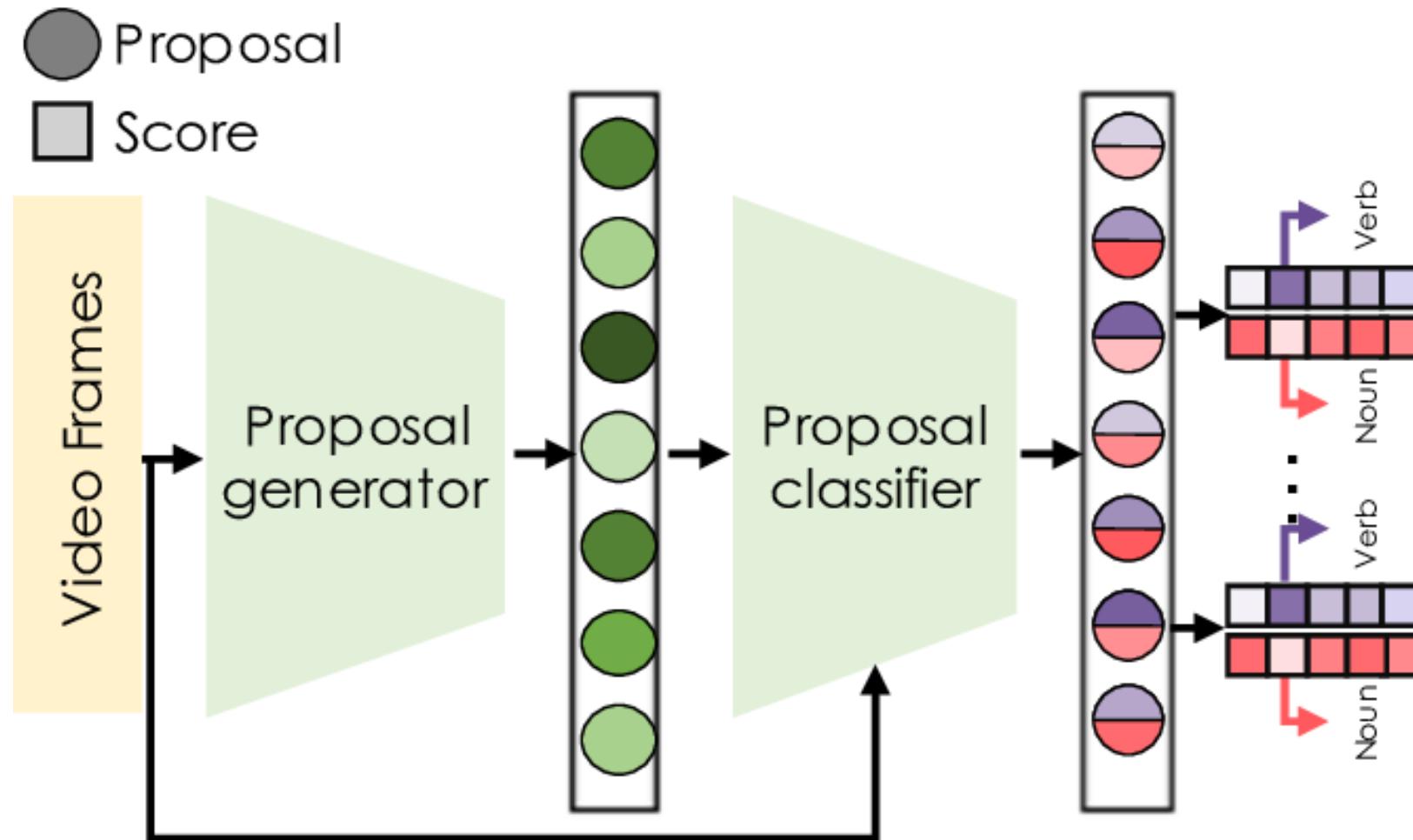
Temporal Action Localization



(Detour) Audio-visual Action Localization



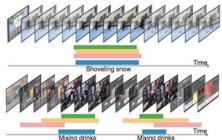
Temporal Action Localization





ACTIVITYNET

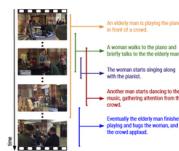
Large Scale Activity Recognition Challenge



ActivityNet Temporal Action Localization

This task is intended to evaluate the ability of algorithms to **temporally localize activities in untrimmed video sequences**.

Here, videos can contain more than one activity instance, and multiple activity categories can appear in the video.



ActivityNet Event Dense-Captioning

This task involves both **detecting and describing events** in a video. For this task, participants will use the ActivityNet Captions dataset, a new large-scale benchmark for dense-captioning events.

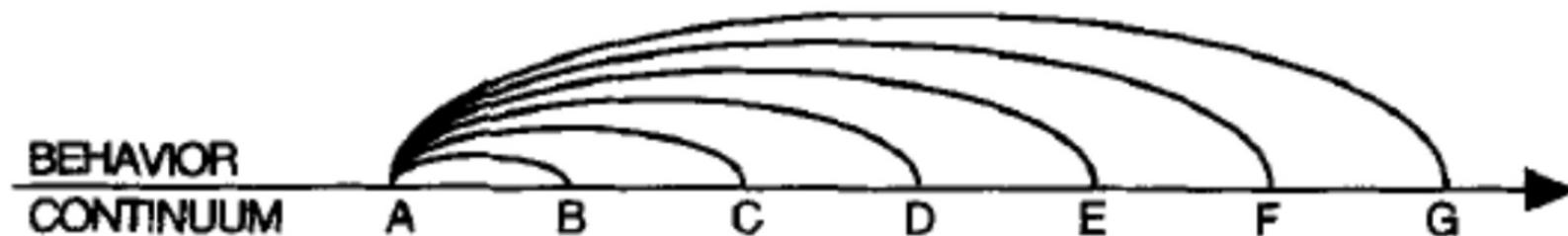


AVA-Kinetics & Active Speakers

This challenge addresses two fundamental problems for spatio-temporal video understanding: (i) localize actions extents in space and time, and (ii) densely detect active speakers in video sequences.

[DETAILS](#)

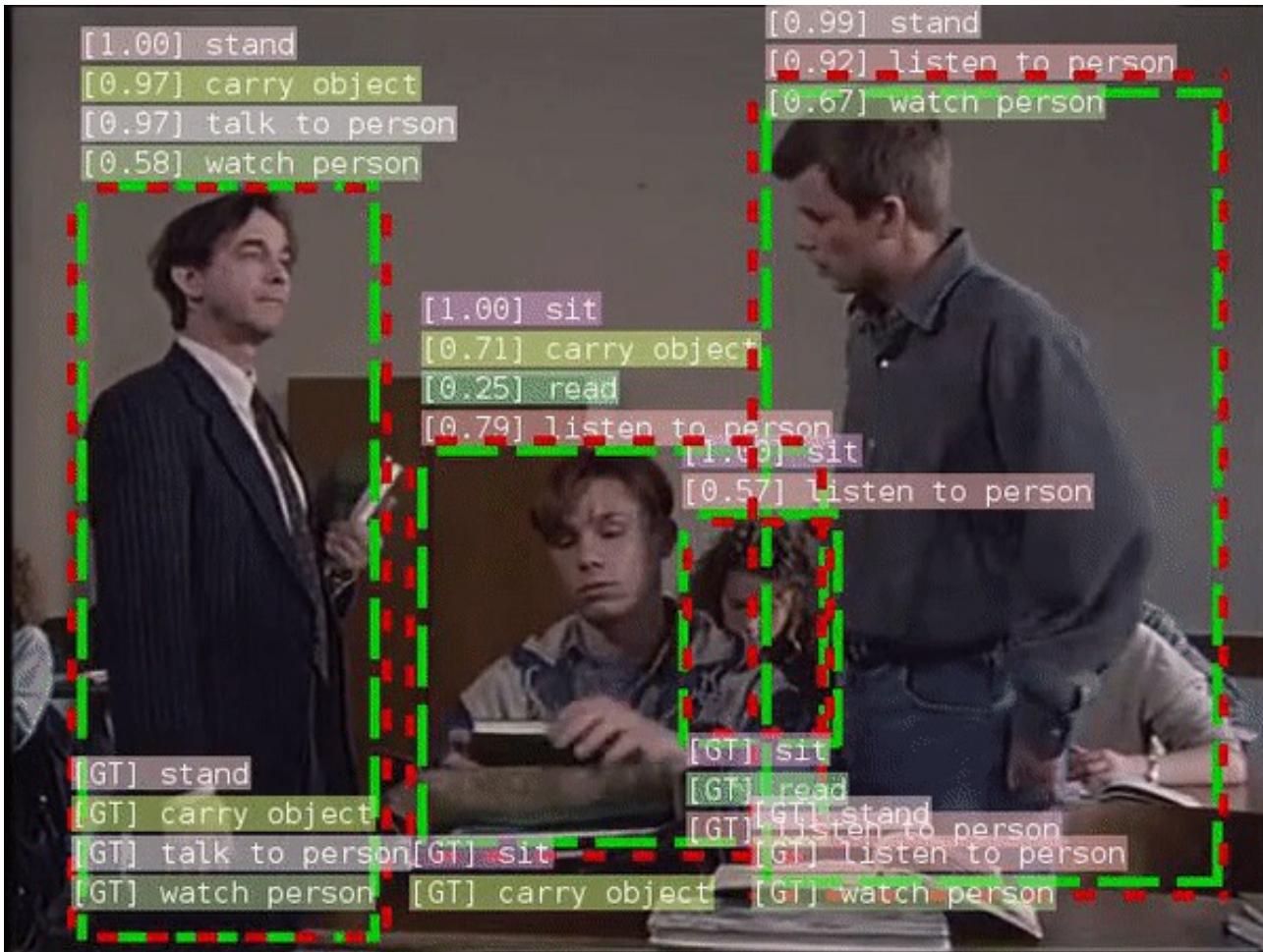
Atomic Visual Actions



- A TO B: STEPPING DOWN FROM THE CURB
- A TO C: CROSSING STREET
- A TO D: WALKING TO SCHOOL
- A TO E: WORKING TO "PASS" FROM THE THIRD GRADE
- A TO F: GETTING AN EDUCATION
- A TO G: CLIMBING TO THE TOP IN LIFE

Figure 2. This figure illustrates the hierarchical nature of an activity. From Barker and Wright [3], pg. 247.

Atomic Visual Actions



- Pose
- Person-Object
- Person-Person

Charades dataset

Hollywood in Homes: Crowdsourcing Data Collection



Sampled Words

Kitchen

vacuum
groceries
chair
refrigerator
pillow

laughing
drinking
putting
washing
closing

AMT

Scripts

"A person is washing their refrigerator. Then, opening it, the person begins putting away their groceries."

"A person opens a refrigerator, and begins drinking out of a jug of milk before closing it."

AMT

Recorded Videos



AMT

Annotations

"A person stands in the kitchen and cleans the fridge. Then start to put groceries away from a bag"

Opening a refrigerator

Putting groceries somewhere

Closing a refrigerator

"person drinks milk from a fridge, they then walk out of the room."

Opening a refrigerator

Drinking from cup/bottle

Why use a single label?

Let's be descriptive ...

Image captioning



The man at bat readies to swing at the pitch while the umpire looks on.



A large bus sitting next to a very tall building.



A horse carrying a large load of hay and two people sitting on it.



Bunk bed with a narrow shelf sitting underneath it.

Video captioning

MSR-VTT



1. A black and white horse runs around.
2. A horse galloping through an open field.
3. A horse is running around in green lush grass.
4. There is a horse running on the grassland.
5. A horse is riding in the grass.

Audio Descriptions (LSMDC)



**Mike leans over and
sees how high they are.**

**Abby clasps her hands
around his face and
kisses him
passionately.**

Text-to-video retrieval

Children and adults are performing various forms of martial arts



A reporter is talking about a movie scene from the wolverines



A man playing guitar and a group of people dancing with him



A person is melting chocolate in a oven



Video Situation Recognition

Event 1



Event 2



Event 3



Event 4



Event 5



Video Situations



Long-form videos

what is long?!

Identity-aware Captioning (LSMDC)



His brow furrowed,
SOMEONE looks down
at the ground.

SOMEONE eyes him
angrily, her jaw
clenched.

SOMEONE heads off.

SOMEONE folds her
arms.

SOMEONE approaches
SOMEONE, who leans
against the wall of the
house.

His brow furrowed, [...]
looks down at the
ground.

[...] eyes him angrily,
her jaw clenched.

[...] heads off.

[...] folds her arms.

[...] approaches [...],
who leans against the
wall of the house.

[PERSON1]

[PERSON2]

[PERSON1]

[PERSON2]

[PERSON1],
[PERSON3]



His brow furrowed,
[PERSON1] looks down
at the ground.

[PERSON2] eyes him
angrily, her jaw
clenched.

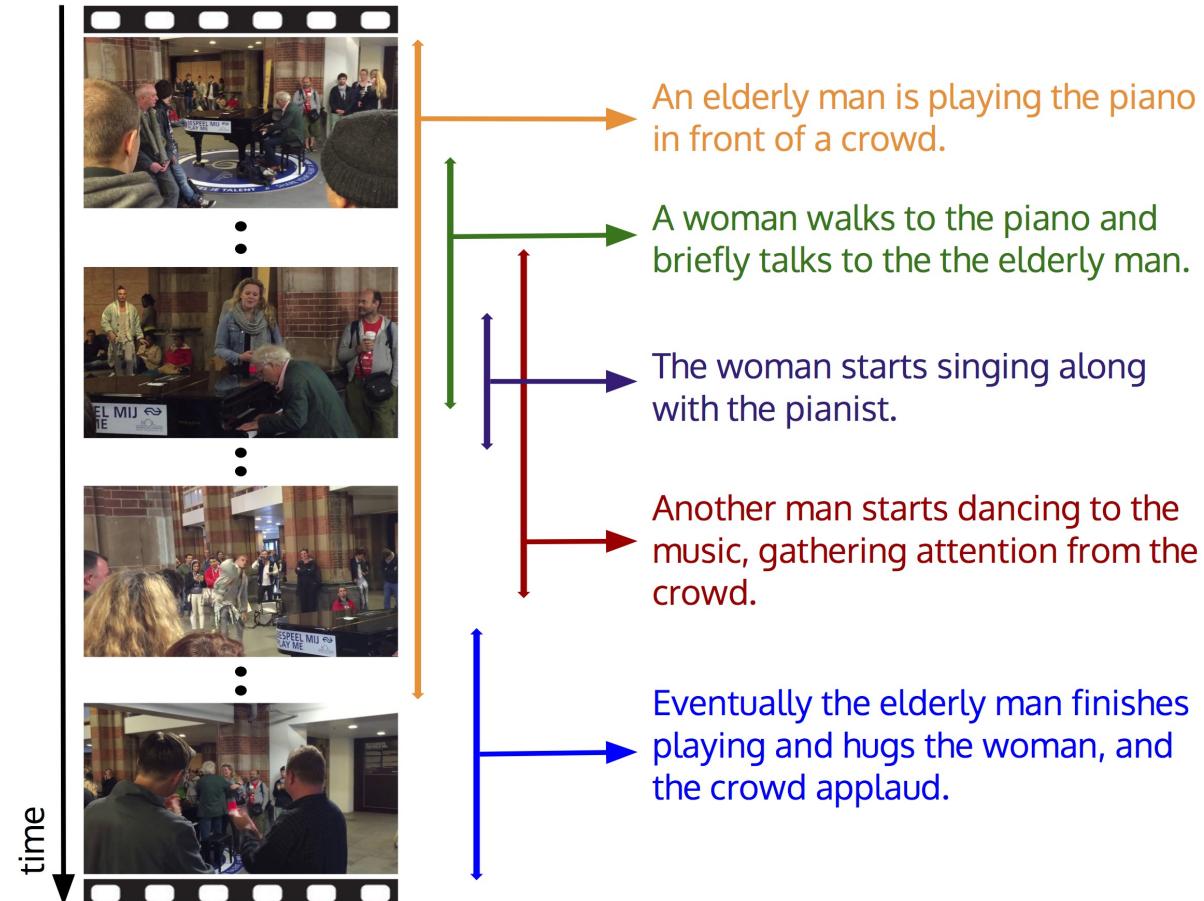
[PERSON1] heads off.

[PERSON2] folds her
arms.

[PERSON1] approaches
[PERSON3], who leans
against the wall of the
house.



Dense event captioning



Long-form Video Understanding

