

HD-EPIC



HD-EPIC: A Highly-Detailed Egocentric Video Dataset



Toby Perrett



Ahmad Darkhalil



Saptarshi Sinha



Omar Emara



Sam Pollard



Kranti Parida



Kaiting Liu



Prajwal Gatti



Siddhant Bansal



Kevin Flanagan



Jacob Chalk



Zhifan Zhu



Rhodri Guerrier



Fahd Abdelazim



Bin Zhu



Davide Moltisanti



Michael Wray

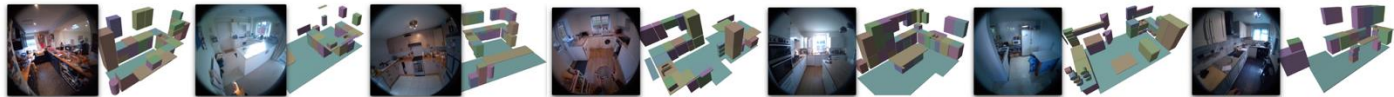


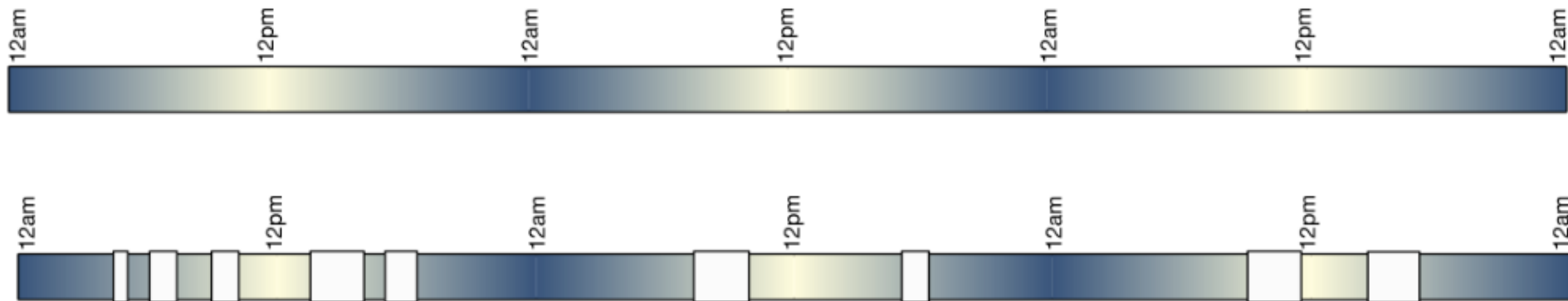
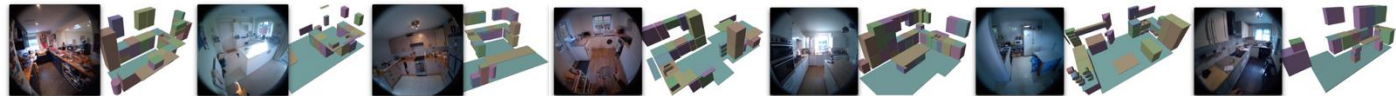
Hazel Doughty

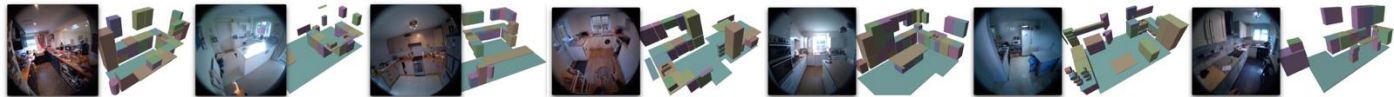


Dima Damen

Dima Damen
April 2025





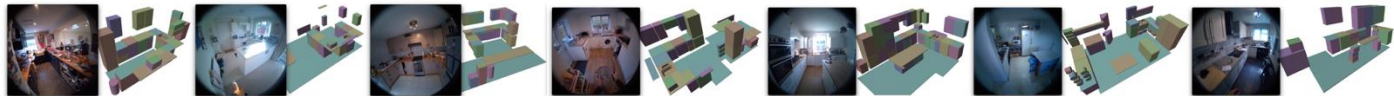


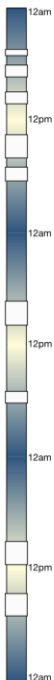
Recorded over 3 days

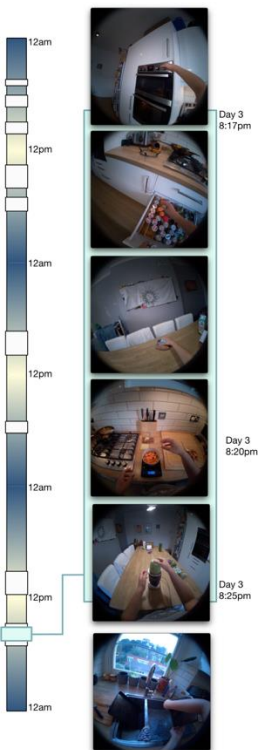


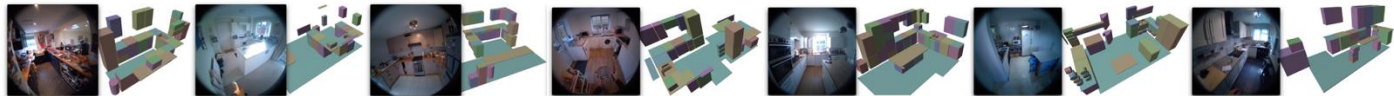


- 9 participants (6 male, 3 female)
- Participants recorded 3.5 to 7.2 hours (avg. 4.6 hours).
- We collected 41.3 hours in 156 videos.
- Long-time commitment (avg. 50 hours): training, data collection, reviewing footage, providing activities/recipe/nutrition and detailed narrations.
- Frames are manually anonymised in RGB as well as SLAM cameras (no easy way to directly map)









Recipe: Southwestern Salad

1: Preheat the oven to 400F

2: Wash and peel the sweet potatoes and chop into bite-sized pieces. Put the sweet potatoes in a bowl and add the olive oil, cumin, and chili powder. Pour onto tray and roast for 10 mins.

3: Pulse all the dressing ingredients in a food processor until mostly smooth.

**Recipe
and nutrition**

Day 3
17pm



Cacio e Pepe (modified)

Ingredients:

~~200 g~~

→ penne

~~400g~~ of pasta of your choice
(we recommend bucatini)

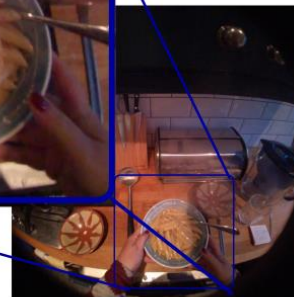
~~2~~ tablespoon of black peppercorn

~~30 g~~

→ parmigiano

200g of freshly grated pecorino cheese

+25g of slightly salted butter



Steps:

1. Toast the peppercorns until fragrant in a dry frying pan over medium heat, about 2 minutes. Keep them moving to prevent them from burning.

~~Once toasted, roughly crush.~~

→ step 2

2. Cook your choice of pasta in a large pot of generously salted boiling water for ~~around 4-6 minutes~~, or until al dente.

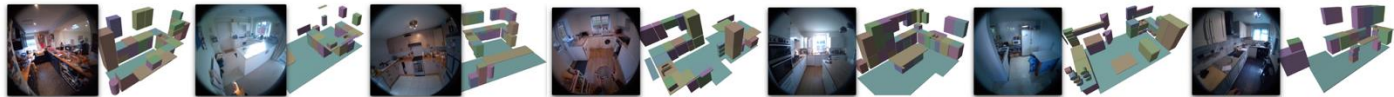
→ step 1

3. While the pasta cooks, add freshly grated cheese and crushed black
 on very low heat

peppercorns to a large serving bowl. Gradually add a cup of the boiling cooking water constantly mixing to obtain a silky, smooth sauce that's able to completely coat the pasta.

→ step 3





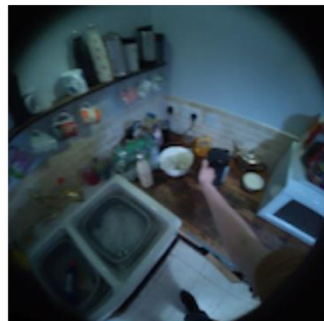
- The **prep** of a corresponding **step** is defined as all essential actions the participant takes to get ready to execute a given step.
- For example, the **step** 'chop tomato':
 - **Prep:** retrieve tomato from storage, wash tomato, retrieve a knife and chopping board.
- the **step** 'add chopped onions and stir':
 - **Prep:** retrieve tomato from storage, wash tomato, retrieve a knife and chopping board, **and chop the onions.**



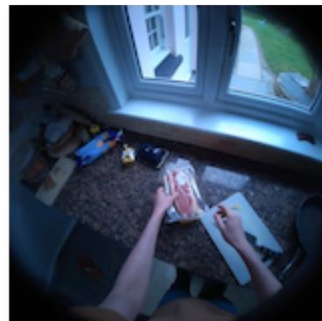
- Prep



- Step



pick up kettle from its base on the counter with my right hand



pick up packet of bacon



pour water from kettle into the pan with my right hand

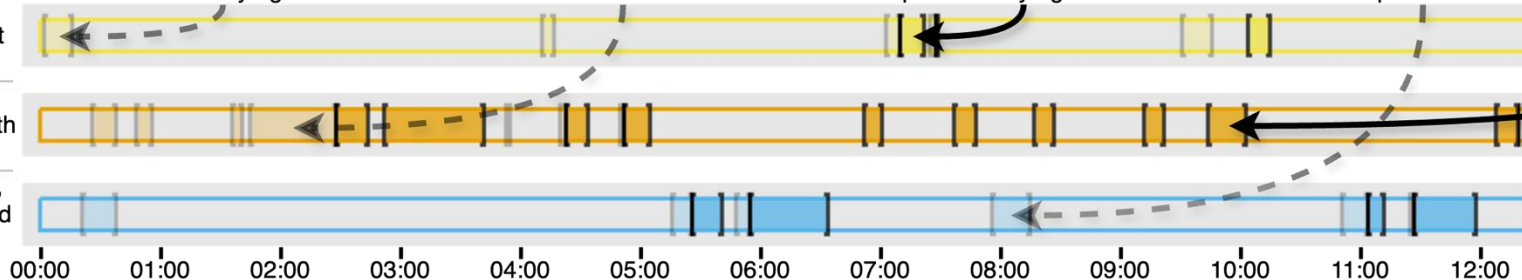


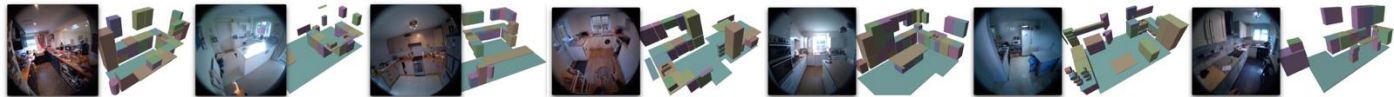
pick up block of cheese that from the top shelf of the

Cook the pasta in a pan of boiling salted water according to the packet instructions.

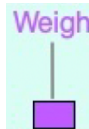
Slice the bacon and place in a non-stick frying pan on a medium heat with half a tablespoon of olive oil and ...

Meanwhile, beat the eggs in a bowl, then finely grate in the Parmesan and mix well.





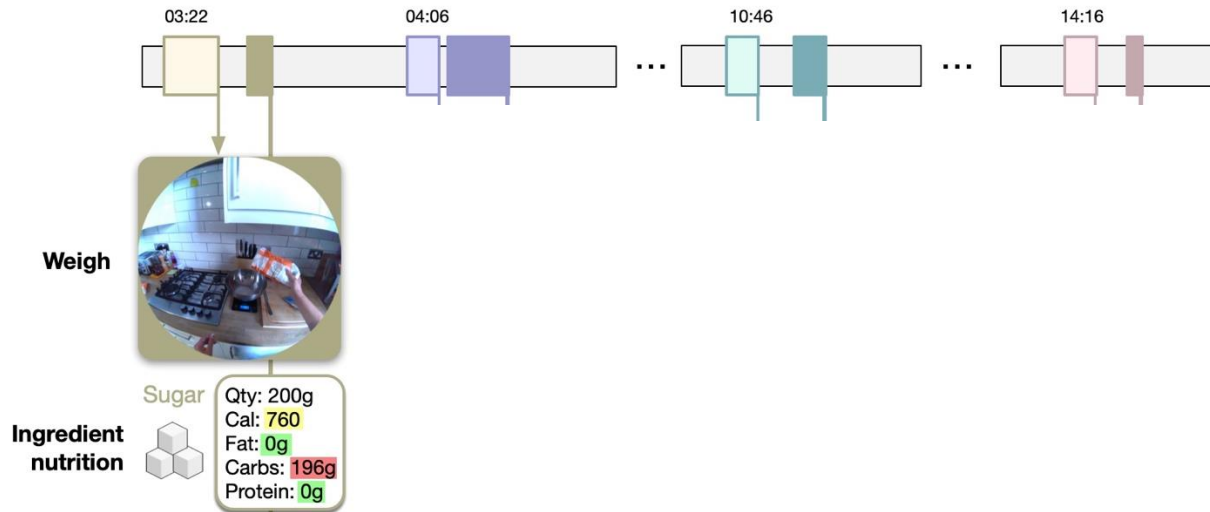
```
"P01_R03_I01": {
  "name": "penne pasta",
  "amount": 125,
  "amount_unit": "g",
  "calories": 445,
  "fat": 1.9,
  "carbs": 90,
  "protein": 15,
```





```
"P01_R03_I01": {
  "name": "penne pasta",
  "amount": 125,
  "amount_unit": "g",
  "calories": 445,
  "fat": 1.9,
  "carbs": 90,
  "protein": 15,
```







- 69 recipes. Avg: 6.6 steps, 8.1 ingredients, 4.8 hours, 2.1 videos.
- Our longest recipe took 2 days and 6 hours to complete.

HD-EPIC Videos	Recipe Videos
Multiple recipes / video, Multiple videos / recipe	One video == one recipe
All the action including preparatory and cleaning/clearing	Only the steps (prep edited out)
Faster (no explanations)	Slower (with descriptions)
Ingredients weighed on camera	Ingredients pre-weighed
Usually average skill level	Usually above average skill level



Recipe and Nutrition

Recipe: Lazy Cake

Step 1: Start by mixing the melted butter, sweetened condensed milk with the cocoa powder. Combine them really well together, but it's okay if there are a few small lumps.

Step 2: Taste the mixture and adjust it to your desired sweetness by either adding more condensed milk for more sweetness or more cacao powder. I prefer it with the amounts stated in the recipe.

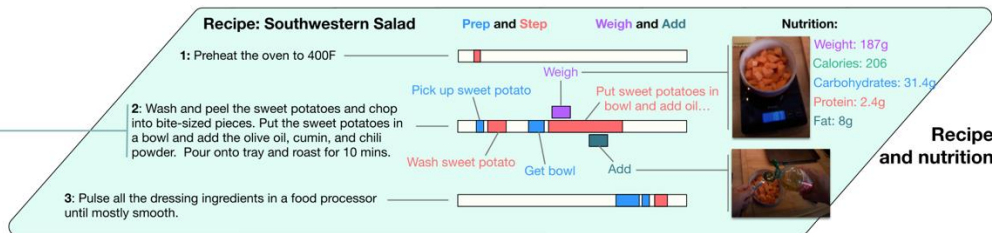
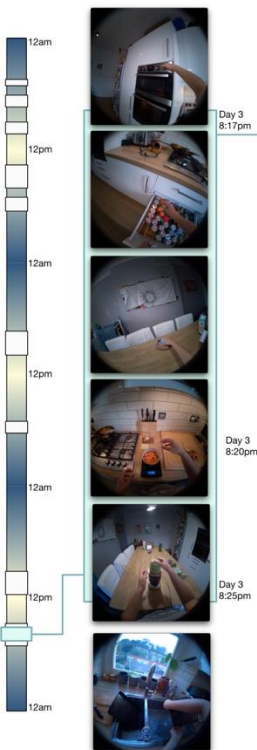
Step 3: Next you need to crumble up the biscuit into the mixture.

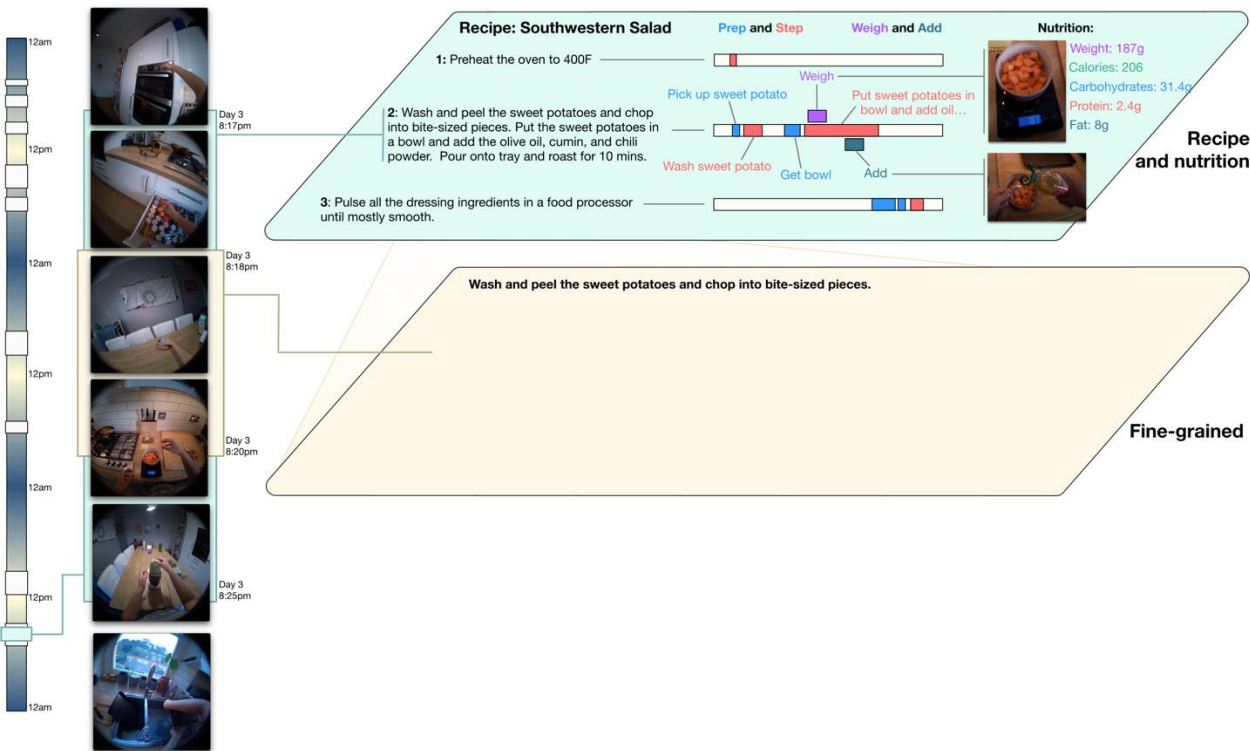
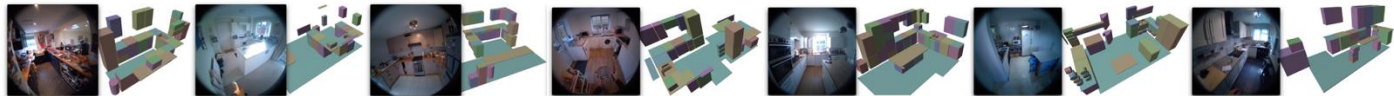


Step 4: It's important to use small sized biscuit crumbs like in the picture, but also try and incorporate really small crumbs in there as well to help the mixture hold.

Step 5: The last step is shaping the mixture. Do this by spreading out some plastic wrap (or cling film) and spooning the mixture in the centre in the shape of a log. Slowly wrap the plastic around the mixture and tighten as you do this to shape into a cylindrical log.

Step 6: Next wrap the log with foil to help it maintain its shape and place it in the fridge for at least 5 hours (depending on your freezer).







Wash and peel the sweet potatoes and chop into bite-sized pieces.



Narrations

Reach the back of the counter, pick up a packet of avocados and a sweet potato and slide them forward to the front of the kitchen island.

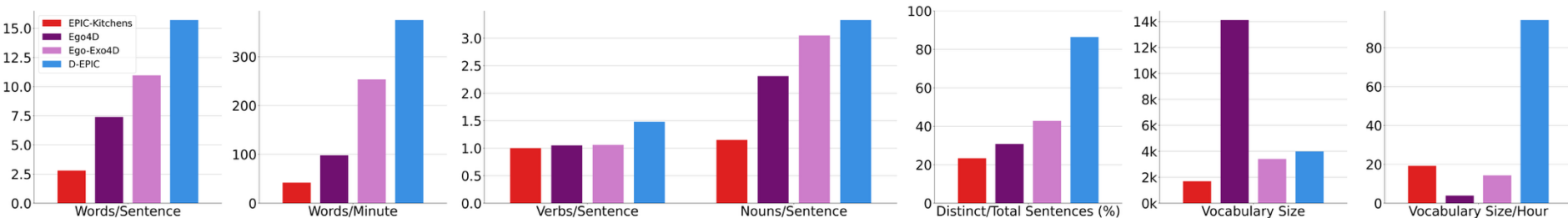
Fine-grained

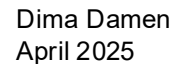
Highly-Detailed Narrations



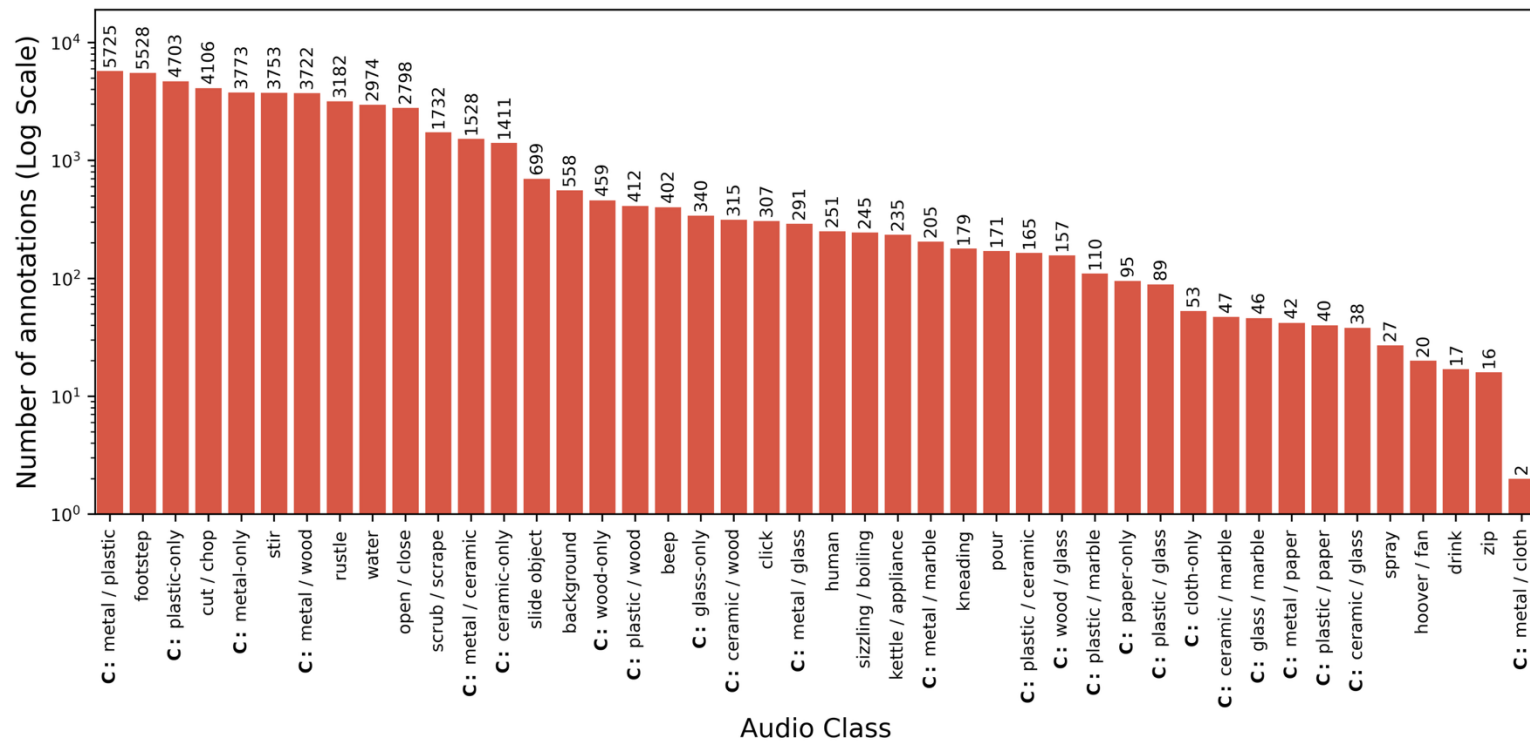


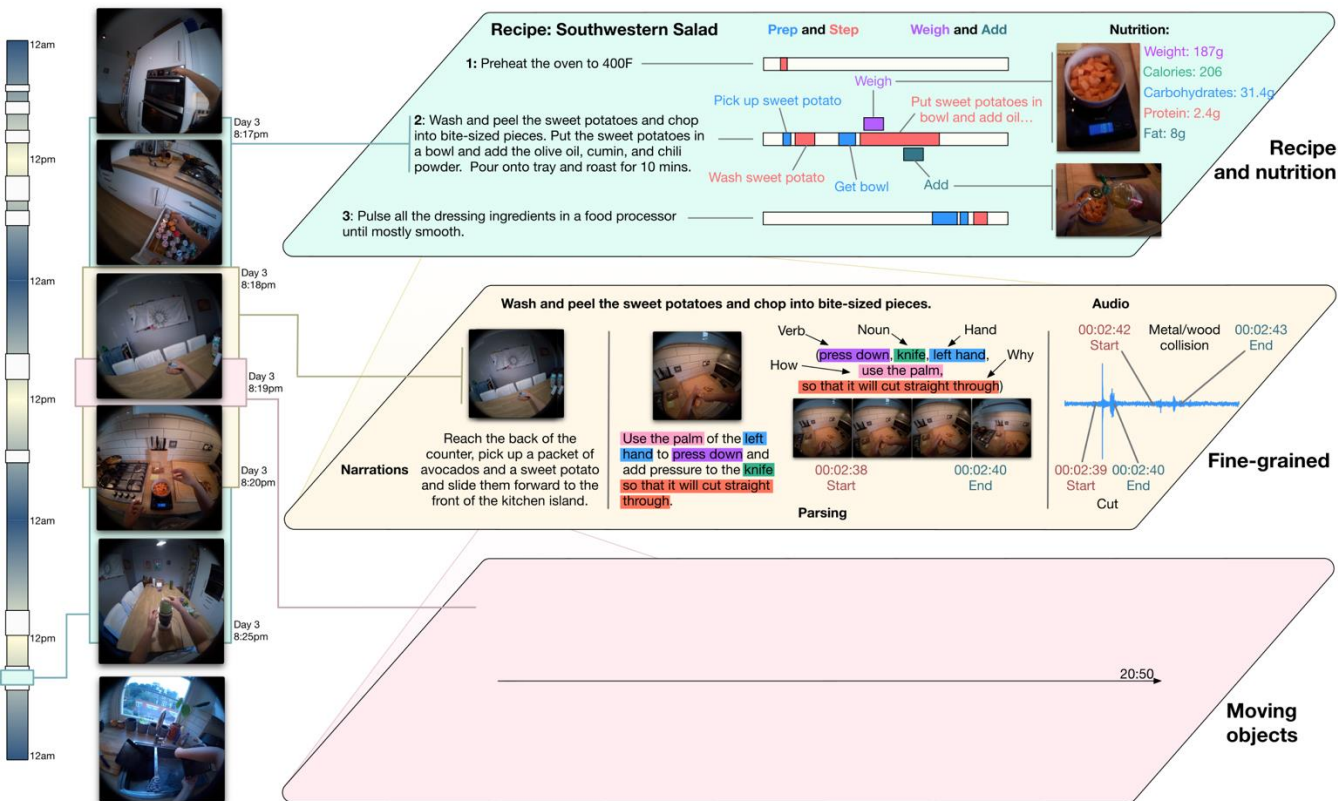
- 59,454 fine-grained actions, with a mean duration of 2.0s (± 3.4 s).

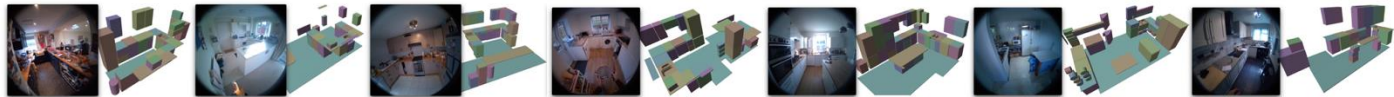




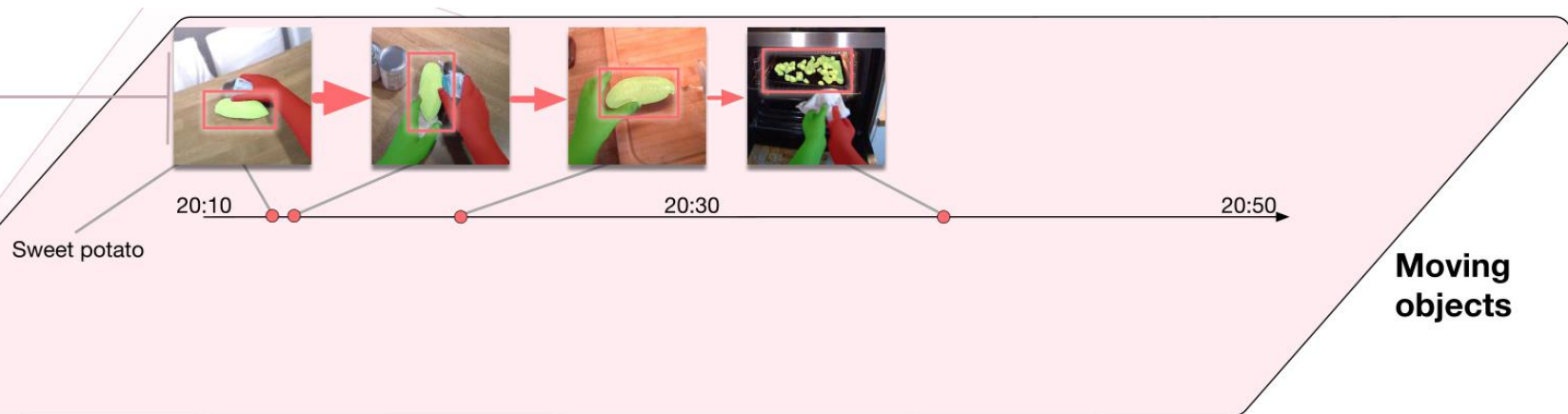


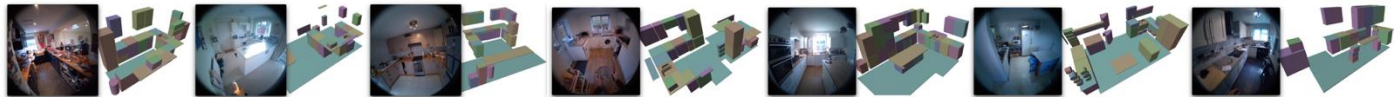






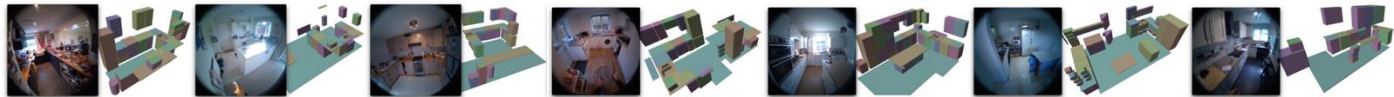
y 3
5pm



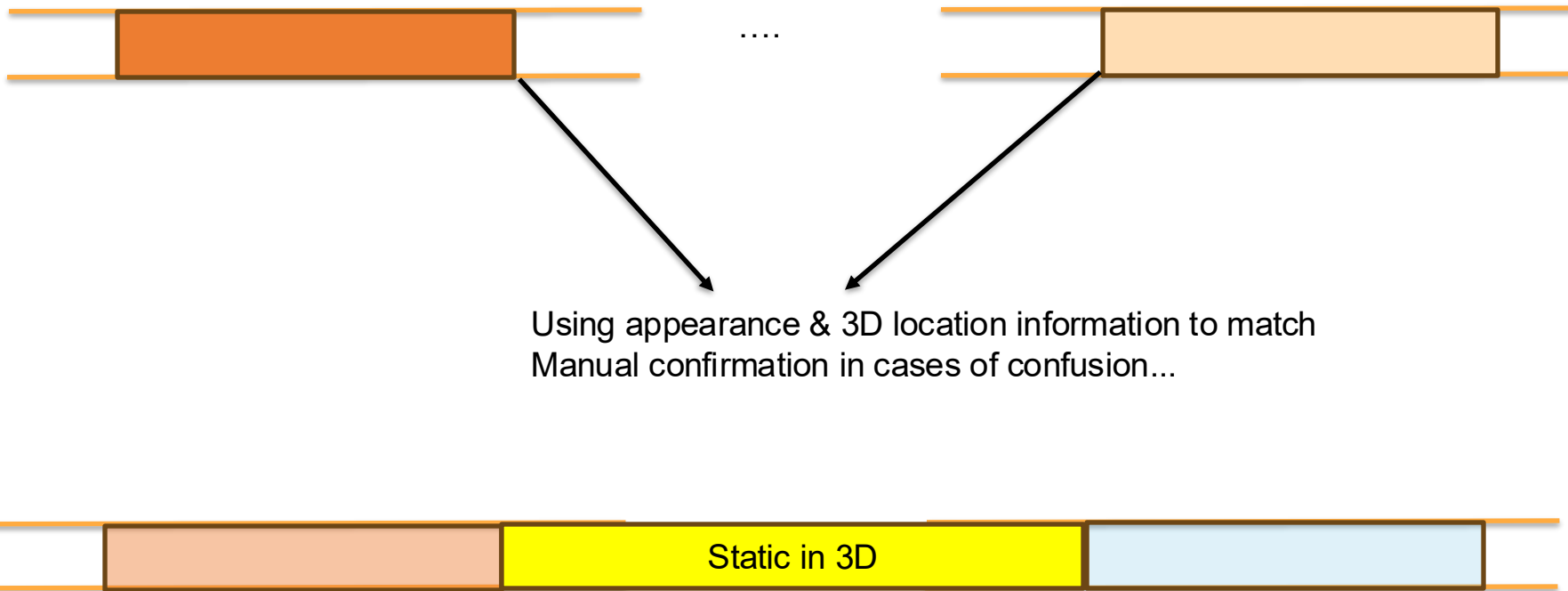


- How to minimize the annotations for tracking objects...





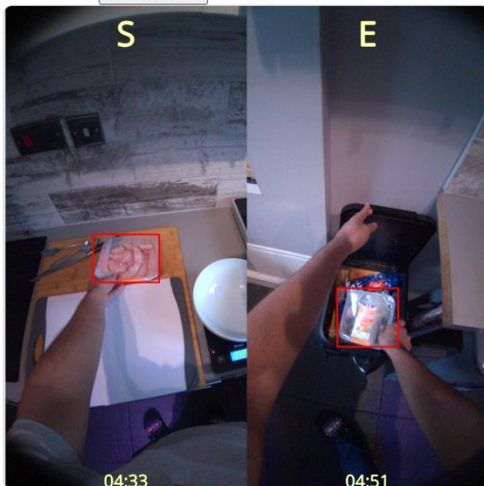
- How to minimize the annotations for tracking objects...





Query Image

Choose Files 201 files



42 / 199

← Previous

Next →

Undo

▼ rubbish bin box of chicken wooden chopping board

Enter Track Name (optional)

Create New Track

Inconsistent Query

Object Tracks

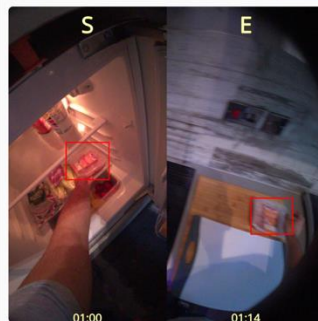
Sort by Distance



Save Tracks

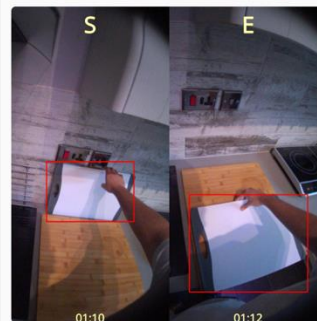
box of chicken (0.0m)

Add



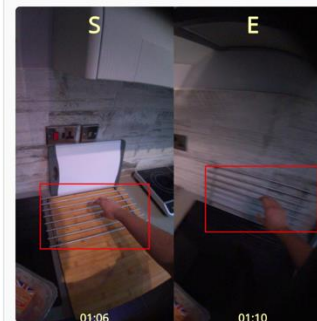
plastic chopping board (0.3m)

Add



metal cooling rack (0.6m)

Add



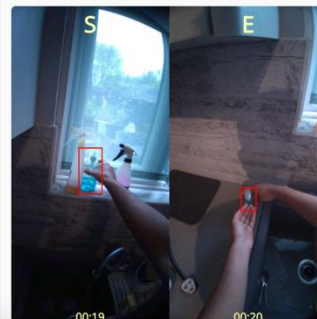
plastic measuring cup (1.0m)

Add



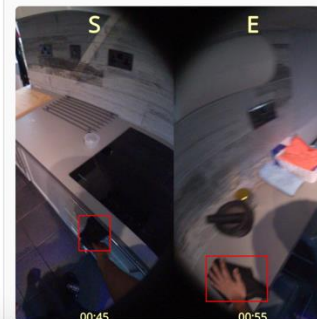
hand washing liquid (1.3m)

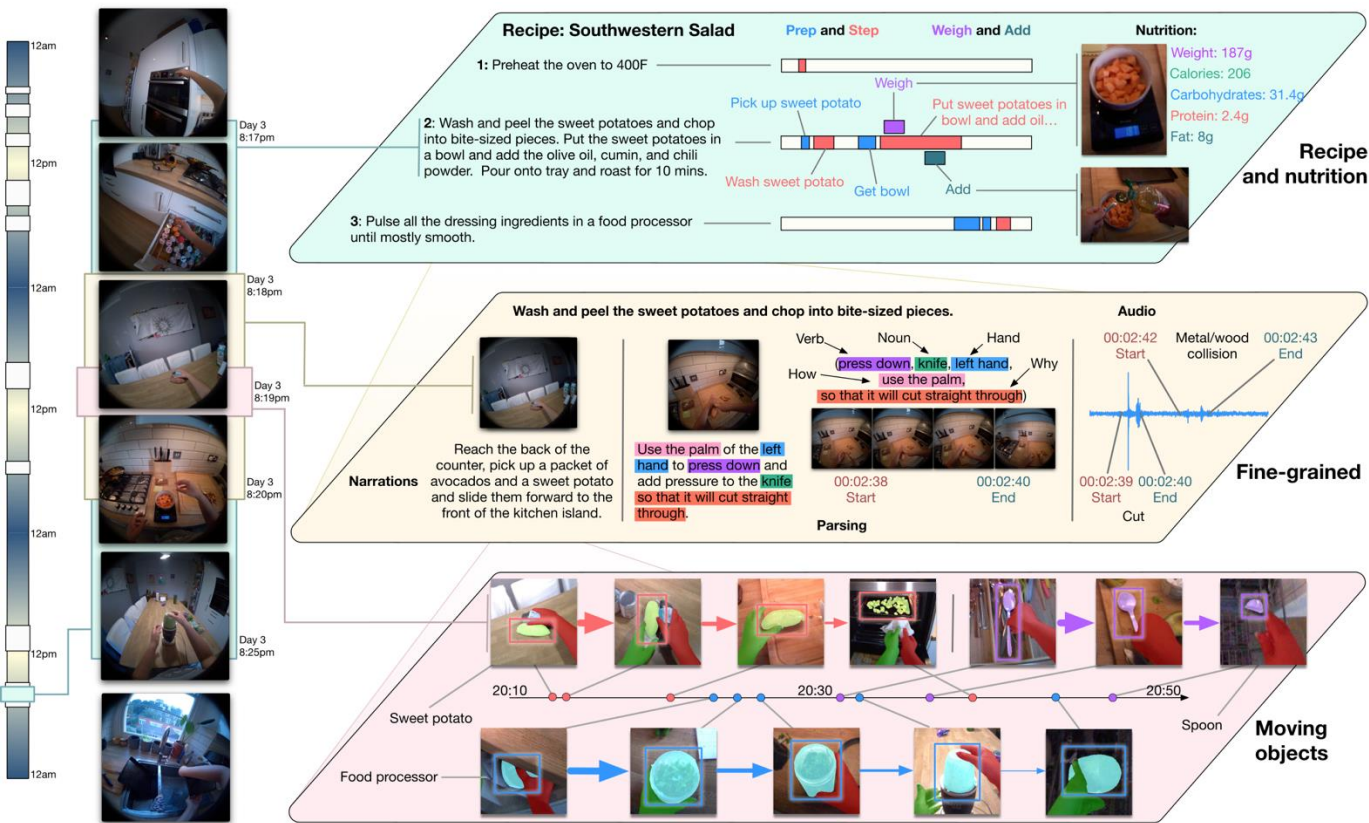
Add



kitchen towel (1.5m)

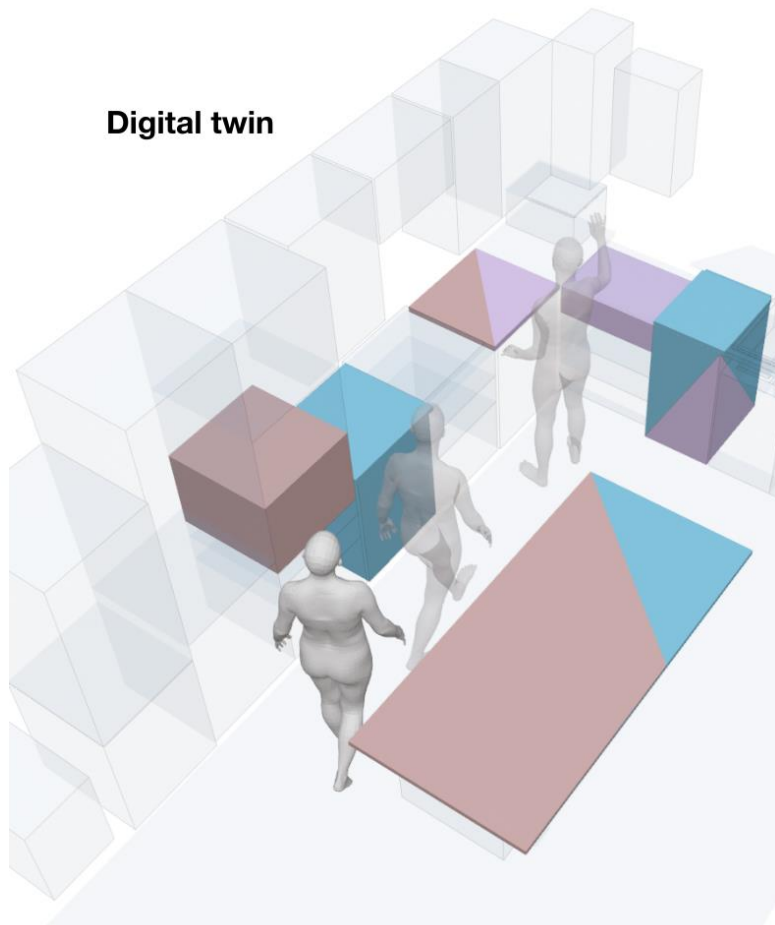
Add





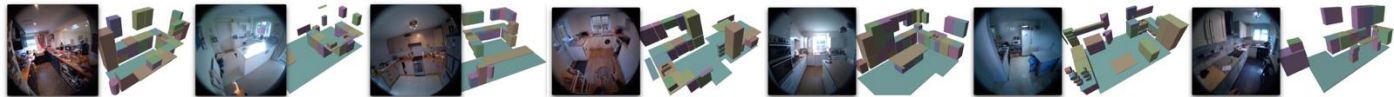


Digital twin

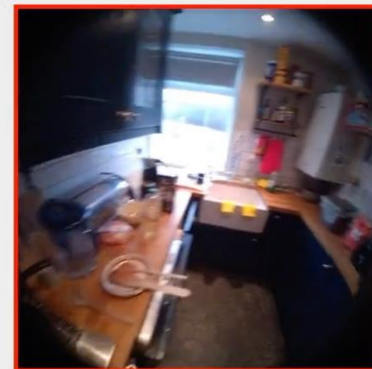
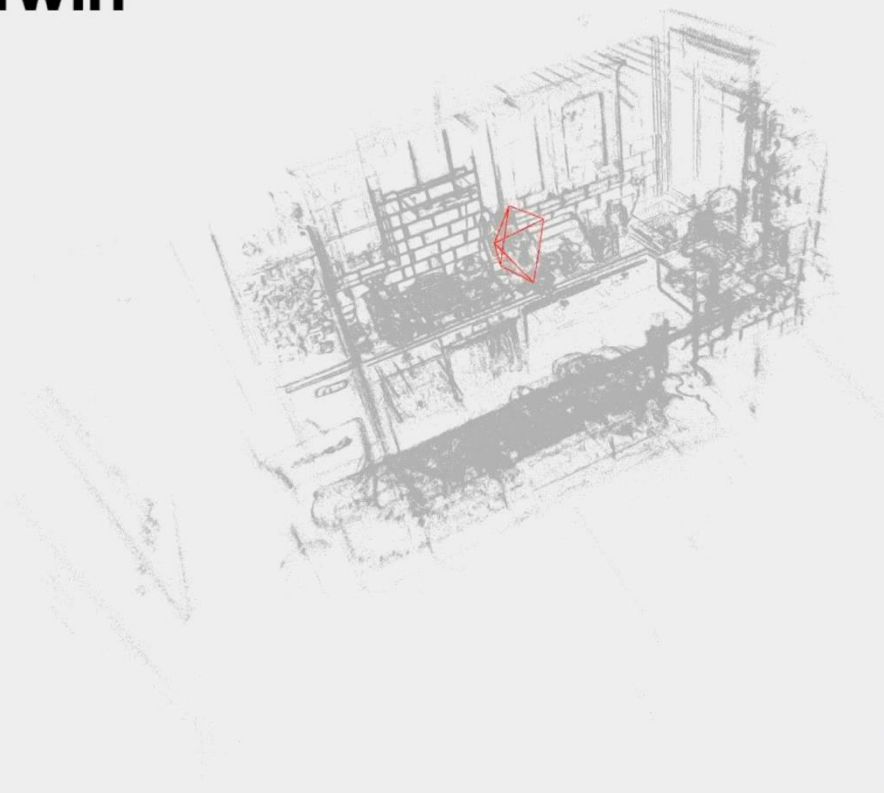


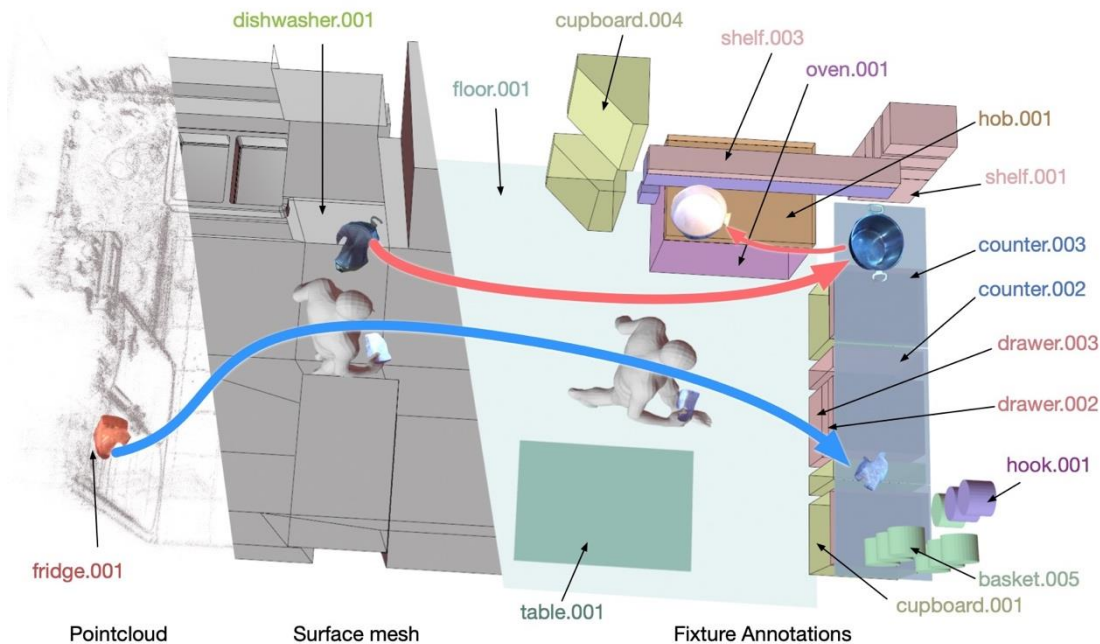
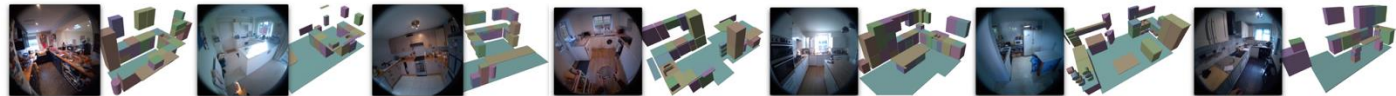
Perrett e

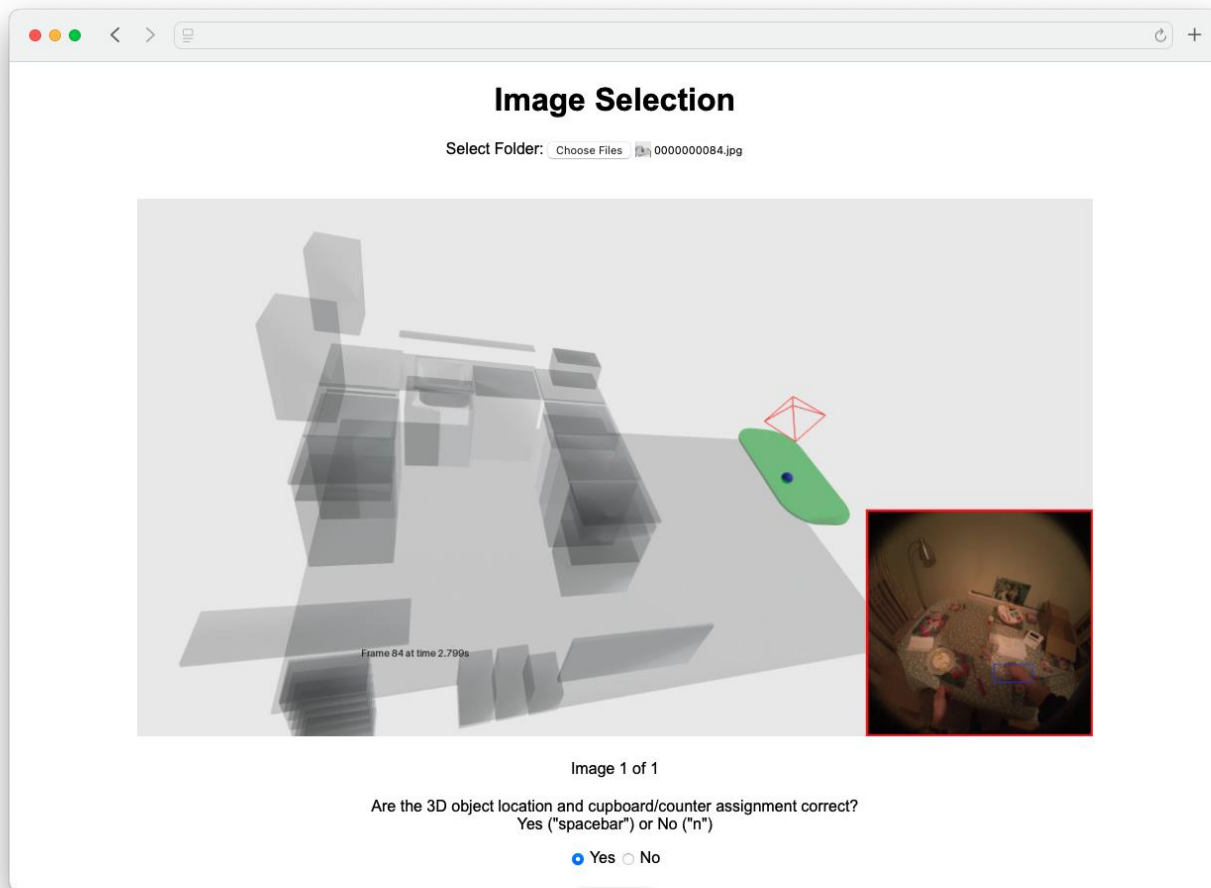
Dima Damen
April 2025

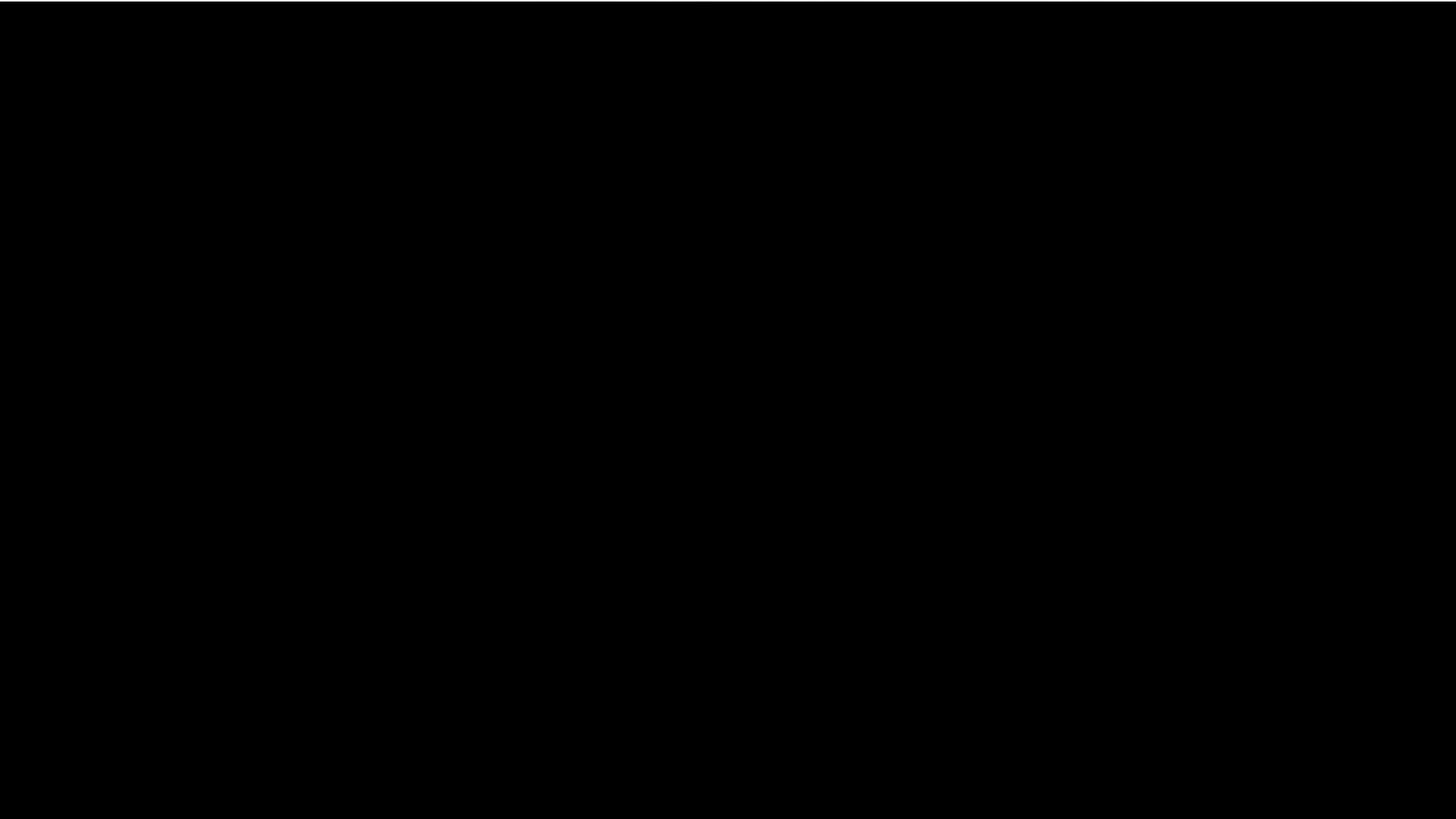


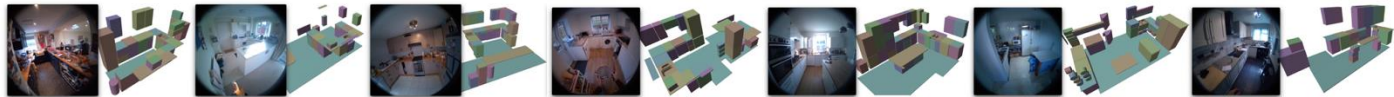
Digital Twin



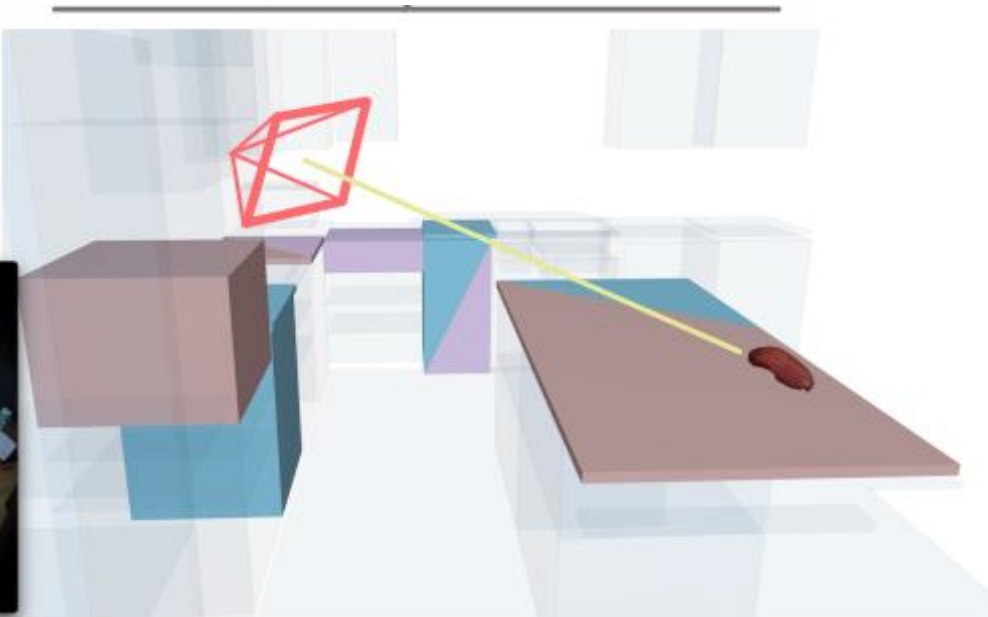


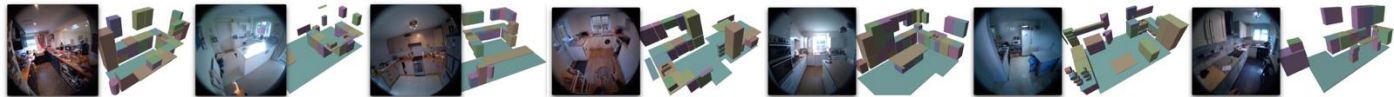




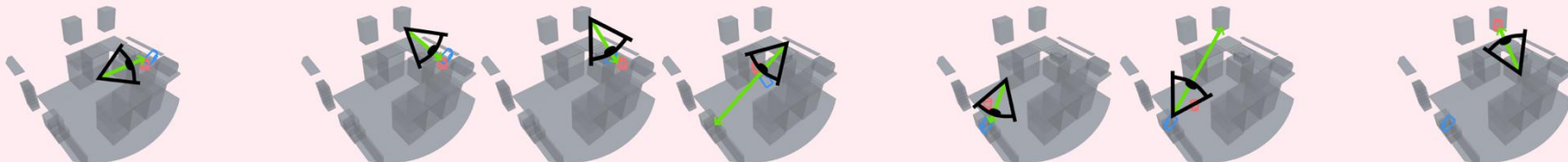


Gaze priming

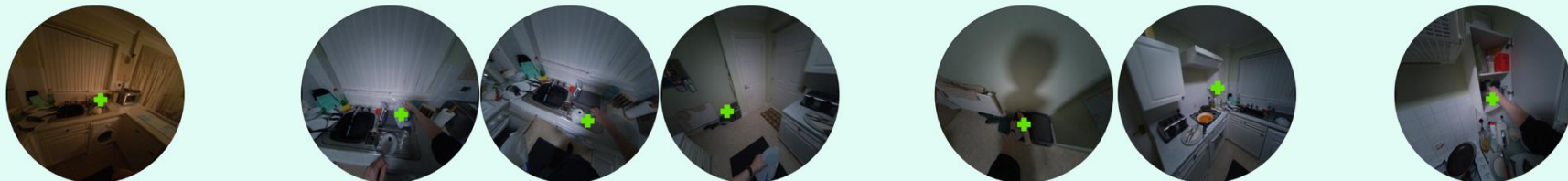




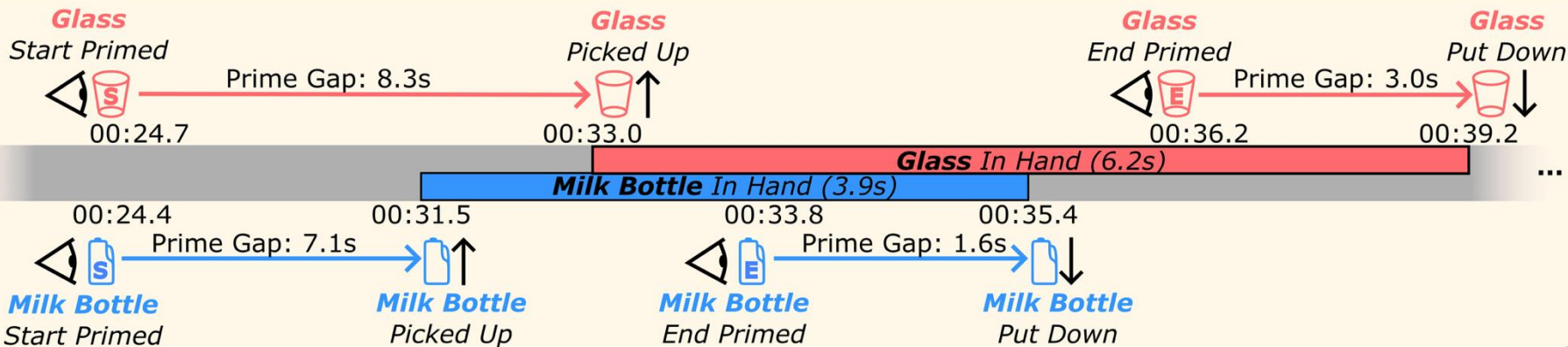
3D Scene

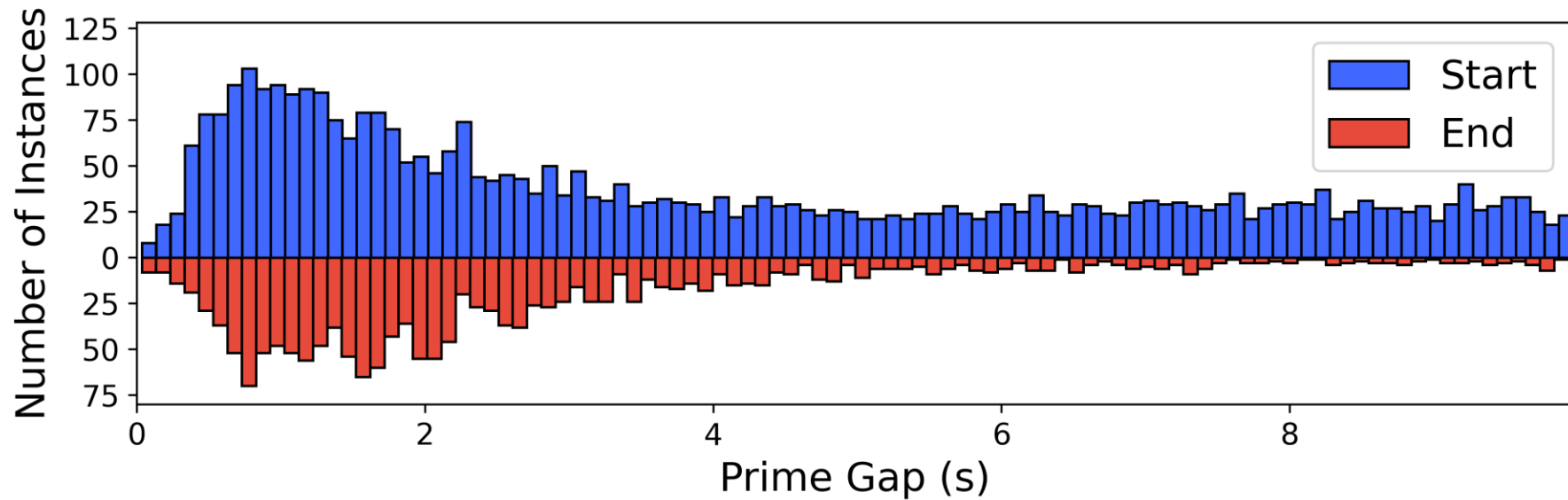


Frames
w/ 2D Gaze



Object Movement

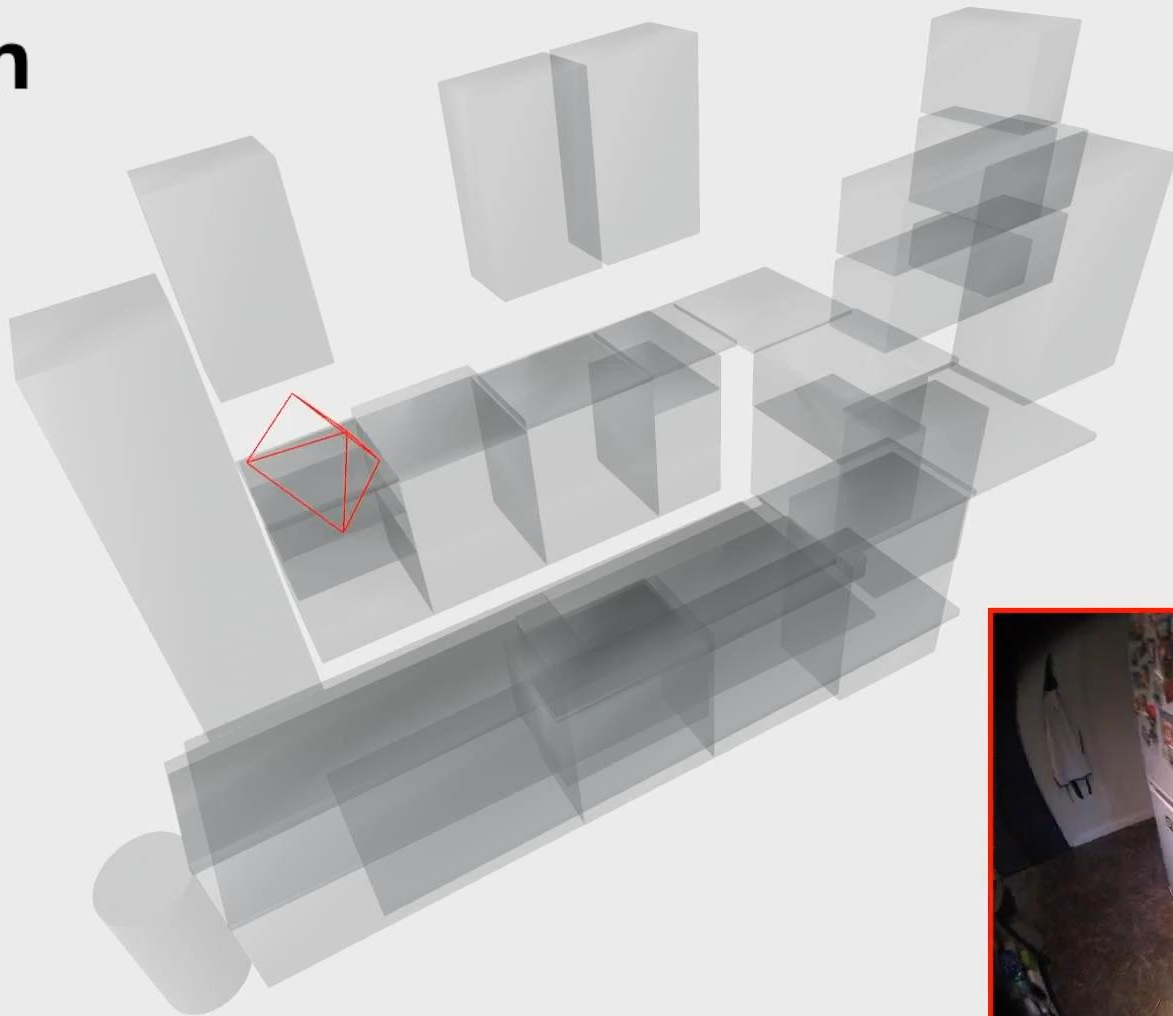


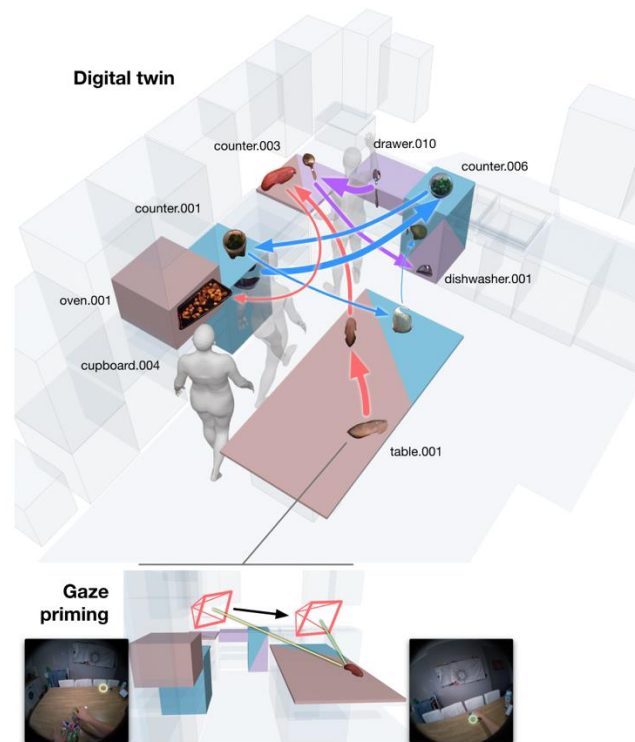
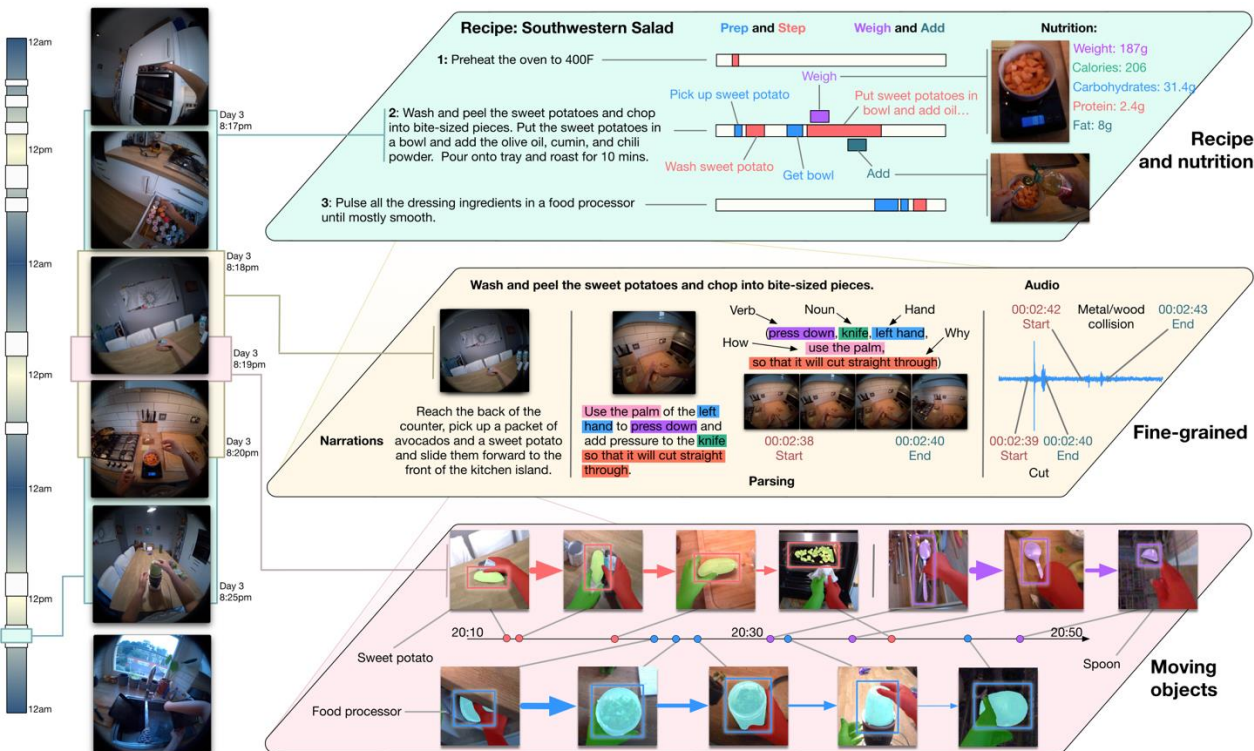
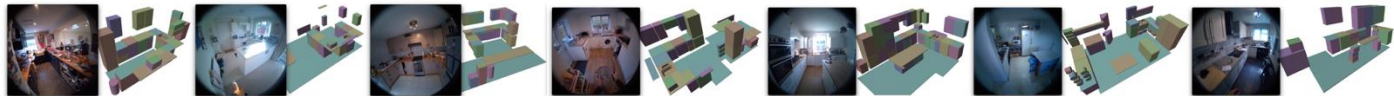


Digital Twin

Fixtures

Open drawer

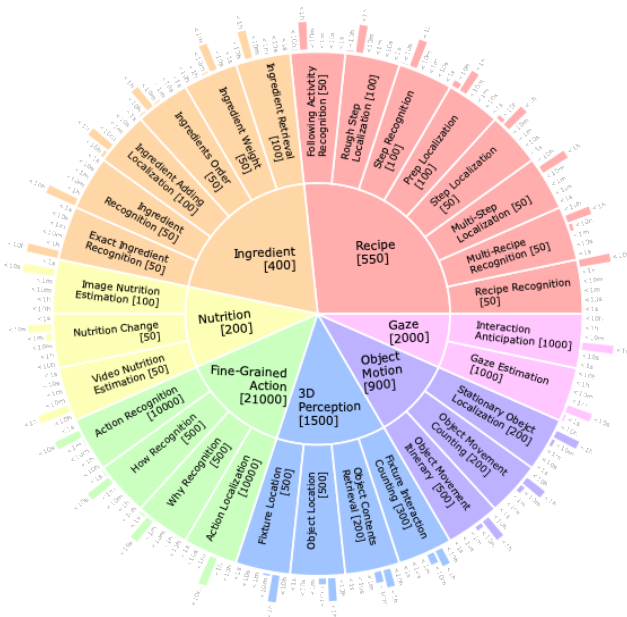
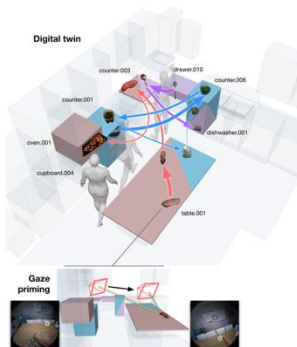
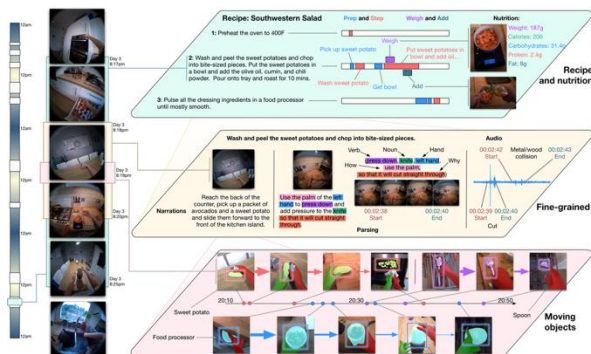






Annotation Type	Total annotations	Annotations/min
Narrations	59,454	24.0
Parsing (Verbs + Nouns + Hands + How + Why)	303,968	122.7
Recipes (Preps + Steps)	4,052	1.6
Sound	50,968	20.6
Action boundaries	59,454	24.0
Object Motion (Pick up + Put down + Fixtures + Bboxes + Masks)	153,480	62.0
Object Itinerary	4,881	2.0
Object Priming (Starts + Ends)	18,264	7.4
Total		263.2

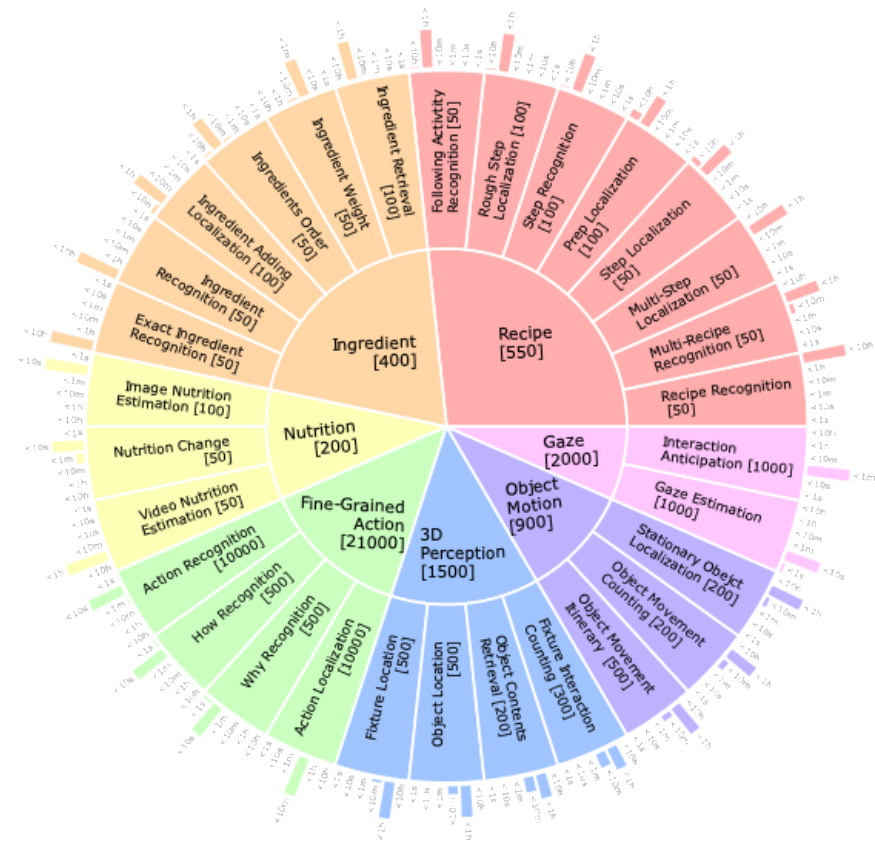
Table A3. HD-EPIC annotations per minute



Sec 2: HD-EPIC VQA Benchmark

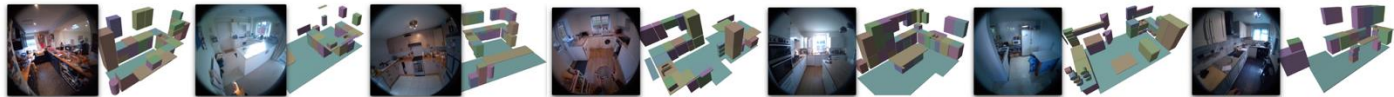


1. **Recipe**. Questions on temporally localising, retrieving, or recognising recipes and their steps.
2. **Ingredient**. Questions on the ingredients used, their weight, their adding time and order.
3. **Nutrition**. Questions on nutrition of ingredients and nutritional changes as ingredients are added to recipes.
4. **Fine-grained action**. What, how, and why of actions and their temporal localisation.
5. **3D perception**. Questions that require the understanding of relative positions of objects in the 3D scene.
6. **Object motion**. Questions on where, when and how many times objects are moved across long videos.
7. **Gaze**. Questions on estimating the fixation on large landmarks and anticipating future object interactions.





Model	Recipe	Ingredient	Nutrition	Action	3D	Motion	Gaze	Avg.
Blind - Language Only								
Llama 3.2	33.5	25.0	36.7	23.3	22.3	25.5	19.5	26.5
Gemini Pro	38.0	26.8	30.0	22.1	21.5	27.7	20.5	26.7
Video-Language								
VideoLlama 2	30.8	25.7	32.7	27.2	25.7	28.5	21.2	27.4
LongVA	29.6	30.8	33.7	30.7	32.9	22.7	24.5	29.3
LLaVA-Video	36.3	33.5	38.7	43.0	27.3	18.9	29.3	32.4
Gemini Pro	60.5	46.2	34.7	39.6	32.5	20.8	28.7	37.6
<i>Sample Human Baseline</i>	96.7	96.7	85.0	92.5	93.8	92.7	75.0	90.3



Input	Recipe	Ingredient	Nutrition	Action	3D	Motion	Gaze	Avg.
Llama 3.2								
A only	26.8	23.8	14.0	20.2	14.9	15.4	17.8	19.0
Q + A	33.5	25.0	36.7	23.3	22.3	25.3	19.5	26.5
GT Narrations + Q + A	70.8	46.3	34.0	62.5	42.9	28.7	29.4	45.0
Gemini Pro								
A only	29.6	21.0	17.7	19.2	18.9	16.3	18.0	20.1
Q + A	38.0	26.8	30.0	22.1	21.5	27.7	20.5	26.7
GT Actions + Q + A	79.0	54.8	36.3	31.3	42.5	32.8	25.5	43.2
GT Narrations + Q + A	82.6	57.5	36.7	63.6	47.6	38.5	29.0	50.8
Video + Q + A	60.5	46.2	34.7	39.6	32.5	20.8	28.7	37.6

Table A4. VQA Input Ablation Our benchmark cannot be solved by analysing Q+A pairs or external knowledge and is a challenge for state-of-the-art closed and open source video VLM models.

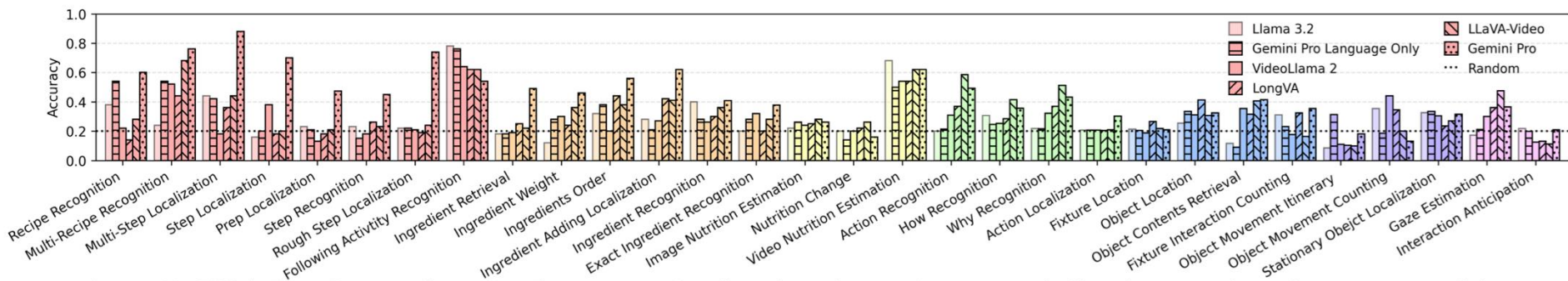
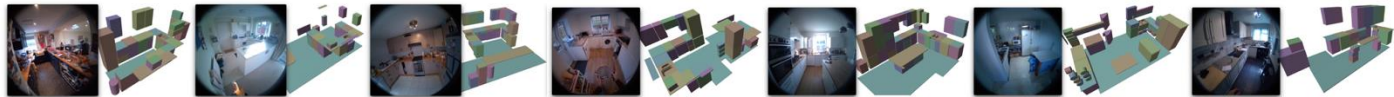


Figure 11. **VQA Results per Question Prototype.** Our benchmark contains many challenging questions for current models.

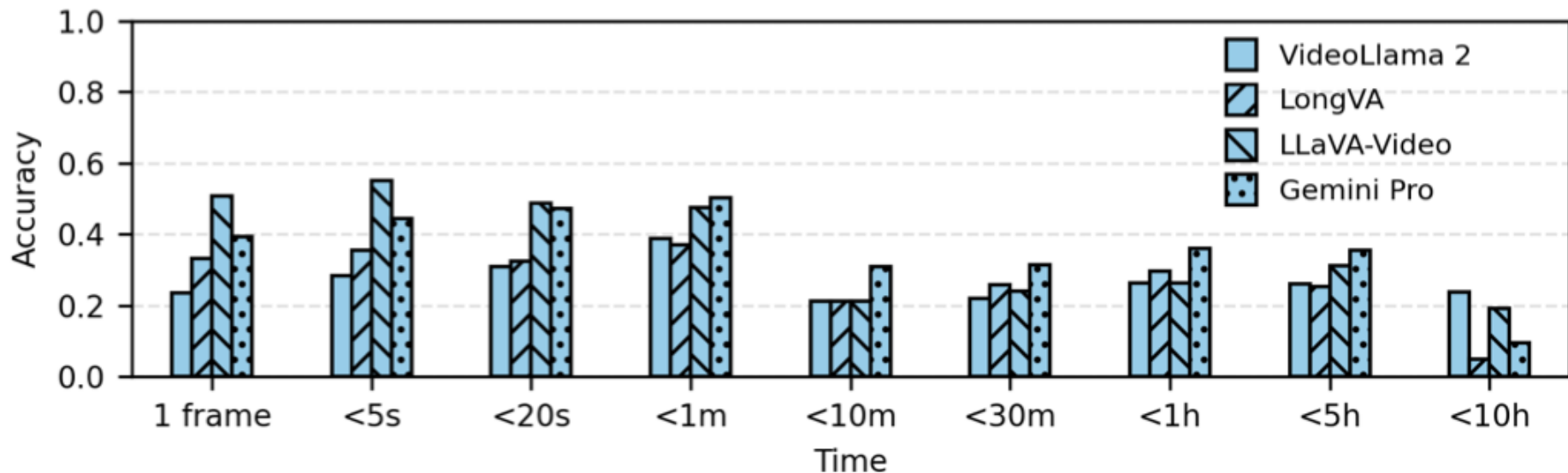


Figure 12. Effect of Input Length. Models struggle with questions of all video input lengths. s=second, m=minute, h=hour.

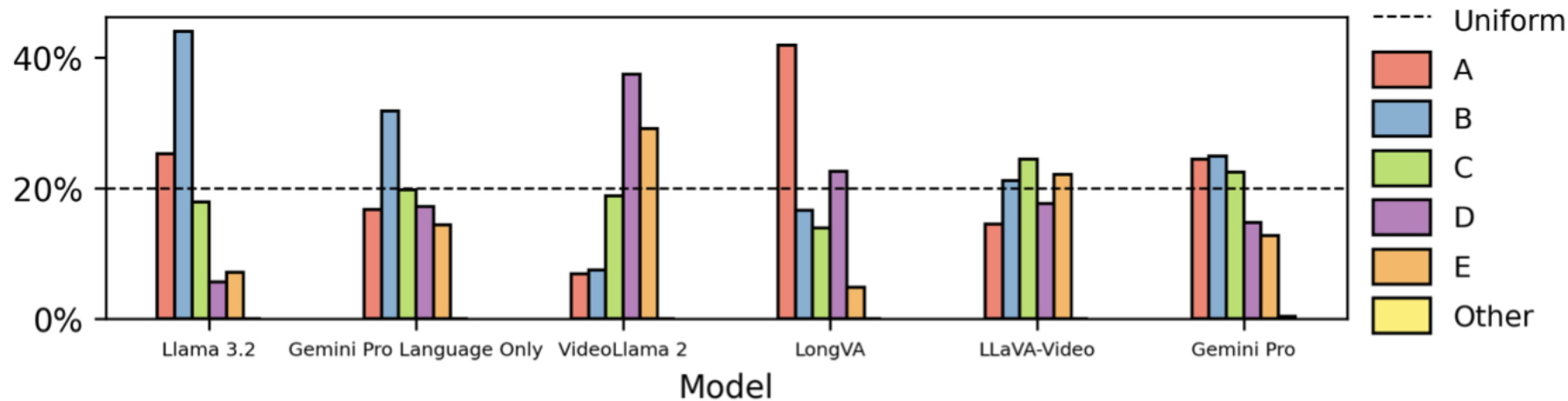
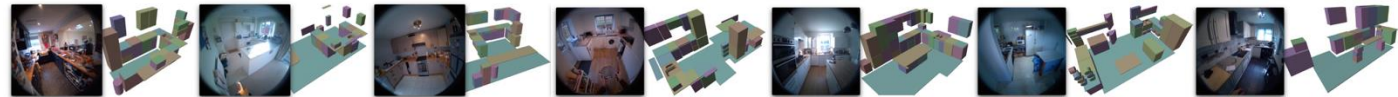
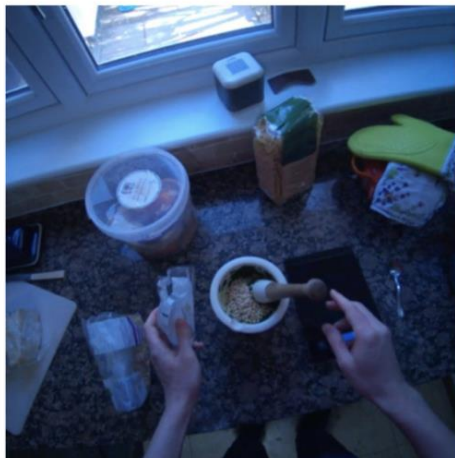
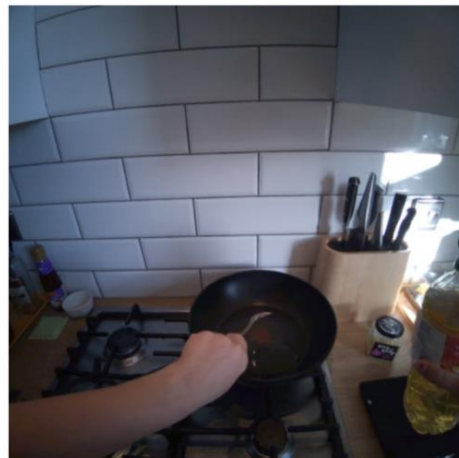


Figure A11. **Prediction Bias of Models.** Most models have a bias in answer, although it is different for each model.



Which of the ingredients in these images showcase higher **fat**?

A. Oil

B. Pine nuts

C. Salted butter

D. Butter

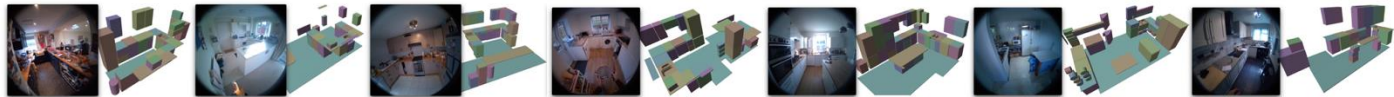
E. Vegetable oil





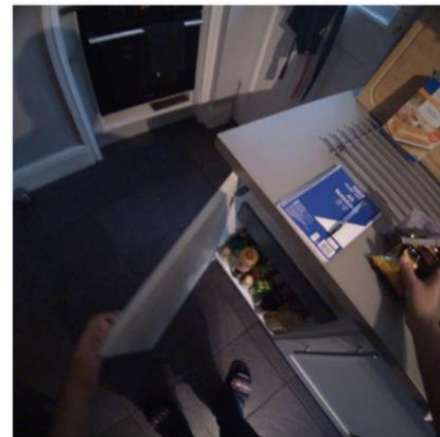
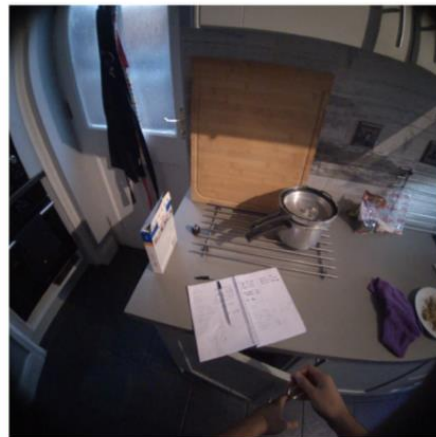
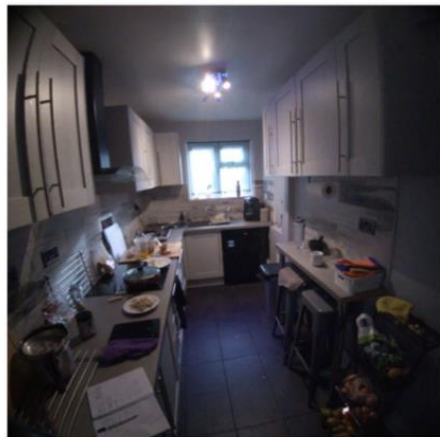
Which of these sentences best describe the action(s) in the video? [00:03:56 - 00:04:03]

- ☒ A. Wash the cutting board using the sponge in right hand, then, rotate the cutting board so that the back side can be washed
- ☐ B. With sponge in right hand, clean cutting board while holding board steady with left hand, then with left hand put cutting board under water to clean from soap
- ☐ C. With left hand, grab cutting board from dish rack, then, with both hands put cutting board down on kitchen counter
- ☐ D. With my left hand, pick up cutting board, then, with both hands, run the cutting board under water to clean
- ☐ E. Pick up cutting board from drying rack using right hand, then, dry the cutting board using tea towel in left hand whilst flipping and rotating the cutting board with right hand



What is the best description for how the person carried out the action **pick up bowl of coconut milk** in this video segment? [00:18:44 - 00:18:46]

- ☐ A. Using both hands holding the bowl from bowl rim.
- ☐ B. By holding both sides using the oven gloves.
- ☒ C. using the right hand and lift the large white bowl up.
- ☐ D. using left hand and removing the fork used to stir it using right hand.
- ☐ E. using both hands from the kitchen top above the dishwasher.



How many times did I **open** the item at bounding box (165, 452, 1408, 1408) in 00:00:57?

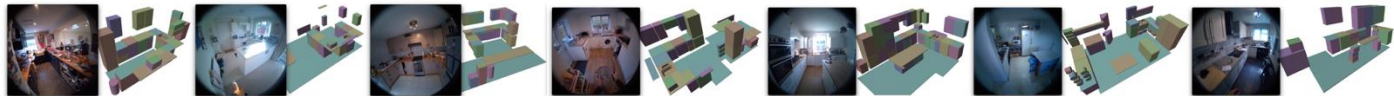
A. 3

 B. 1

C. 4

D. **5**

 E. 2

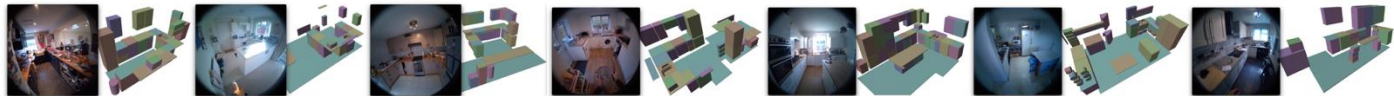




- ✓ video 1
- video 2
- video 3
- video 4
- video 5
- video 6
- video 7
- video 8
- video 9
- video 10
- video 11
- video 12
- video 13

Question: Which of these recipes were carried out by the participant?

- ☐ Chicken Korma
- ☐ Mutton Vindaloo
- ☐ Keema Matar
- ☐ Mangsho Bhuna
- ☐ Katsu Chicken



HD-EPIC

A Highly-Detailed Egocentric Video
Dataset

[Paper \(ArXiv\)](#)

[The Dataset](#)

[Explore Samples](#)

[Watch Video](#)

[VQA benchmark](#)

[Explore VQA](#)

[Download](#)

[Team](#)



HD-EPIC: A Highly-Detailed Egocentric Video Dataset



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