

HD-EPIC: A Highly-Detailed Egocentric Video Dataset

Supplementary Material

Contents

A Video Showcase	1
B Data Collection	1
B.1. Recruitment and Equipment	1
B.2. Instructions and Collected Data	1
B.3. Narrations	2
B.4. Post-Processing: Multi-Video SLAM and Gaze	2
C Annotation Pipeline	3
C.1. Annotating Recipe Steps and Ingredients . .	3
C.2. Fine-Grained Actions	3
C.3. Digital Twin	6
C.4. Density of Our Annotations	8
D Benchmark and Results	8
D.1. VQA Benchmark Details	9
D.2. Additional VQA Experiments	13
D.3. VQA Model Details	14
D.4. Long-Term VOS Benchmark	15
D.5. Recognition Benchmarks	15

A. Video Showcase

We share a video showcasing our annotations on our webpage: <http://hd-epic.github.io> containing three parts.

start → 01:16. Annotations overview. This is a walk-through of Fig. 1, showing our annotation hierarchy and how annotations are linked to the 3D Digital Twin.

01:16 → 04:24. Annotation showcase on one sequence from one kitchen. For one video, we show a walk through of the various annotations on this same video.

04:25 → end. Annotation examples from other kitchens. First are examples of our highly-detailed narrations, followed by examples of hand and object segmentations. We then present audio annotations with waveforms. Next, examples of object-fixture assignment are shown in the Digital Twin. Finally, examples of ingredient nutrition, running total nutrition, and total nutrition for whole recipes are shown.

B. Data Collection

In this section, we provide more details of: the recruitment and equipment used in Sec. B.1; instructions and collected data in Sec. B.2; narrations in Sec. B.3; and Post-Processing in Sec. B.4.



Figure A1. Participant P04 unpacking the data recording equipment.

B.1. Recruitment and Equipment

Each participant engaged in a long-time commitment (avg. 50 hours): data collection, reviewing footage, and providing detailed narrations including activities/recipe/nutrition information. Participants signed a consent form and a data release form after which they were rewarded monetarily. In Fig. A1, we show the devices provided to participants for the data collection, packaged tightly in one backpack that the participants took home. Data collection was carried out using Project Aria devices [72]—a multi-sensor platform, in the form factor of a pair of glasses¹, with 3 front facing cameras (1 RGB and 2 SLAM), 7 microphones and inward facing cameras for gaze estimation². Due to the storage limits, four devices were provided to each participant. We supplied a smartphone with the Aria Mobile Companion app, with all four devices already registered, as well as the MyFitnessPal app [3] installed. Due to battery life limits, devices were sometimes wired to a pocket-carried power bank whilst recording long sessions. Participants also received recording instructions and a set of digital weighing scales for nutritional tracking. The scales’ display is clear to read, to enable OCR readings.

B.2. Instructions and Collected Data

Participants were recruited to record their daily kitchen activities for at least three consecutive days. All recordings were unscripted: participants were asked to wear the glasses every time they walked into their kitchen, pressing the recording button upon entering, and stopping the recording when they left the kitchen. Videos were recorded between January and July 2024. The number of hours collected per participant ranged from 3.5 to 7.2, with an av-

¹Participants who require glasses and cannot use lenses were excluded

²We used Profile 28—1408 × 1408 resolution 30 FPS RGB footage, 60 FPS eye tracking, 30 FPS SLAM

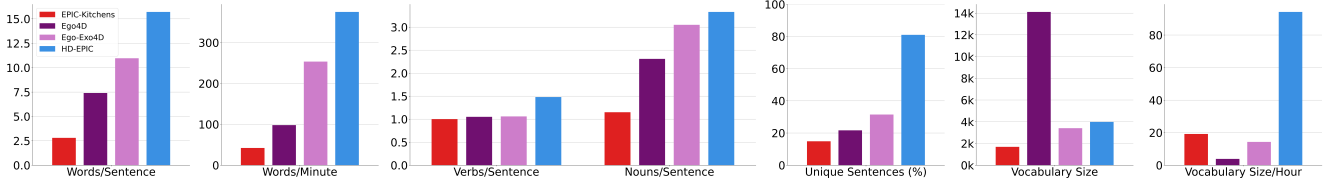


Figure A2. Dataset Statistics: Comparison with EPIC-Kitchens, Ego4D and Ego-Exo4D.

erage of 4.6. We collected a total of 156 videos, with an average length of 15.9 (± 14.5) mins. In total, we collected 41.3 hours (4.46M frames) of footage from 9 participants.

Following data collection, participants provided the recipe they freely prepared. They cited the source of the recipe (e.g. online website or cookbook) and any modifications to ingredients or steps. We collected a total of 69 recipes from the 9 participants. Our recipes cover various cuisines including: Masala Dosa (Indian), Cacio e Pepe (Italian), Sfesiha (Lebanese), and Banana Bread Chocolate Chip Cookies (American) with an average (max) of 6.6 (18) steps and 8.1 (24) ingredients per recipe. Additionally, participants provided time segments (as start-end times) covering each step of their recipes. Recipes typically expand over several videos, on average taking 4 hours and 48 minutes to complete across 2.1 videos from starting preparation to conclusion³. Our longest captured recipe (P03_R03) took 2 days and 6 hours to complete.

We found that participants would interleave the preparation of multiple dishes, often preparing several courses simultaneously (e.g. salad or dessert while preparing a main). Whilst natural in daily routines, this has not been captured in previous recordings which consist of edited online videos [93] or short recipes captured under controlled settings [29].

To track the nutrition values of prepared meals, we instructed participants to weigh ingredients using the digital scales and log nutrition values via the provided MyFitnessPal app. By combining the weight and ingredient information, detailed nutrition information associated with recipes were collected, adding an interesting dimension where we can estimate the weighed ingredients from visual data. While spices are often included as ingredients, we do not capture their nutrition values as they rarely alter these values and are mainly included to improve the dish’s taste. In total, participants used 558 ingredients including: high protein ingredients such as tuna and kidney beans, high carb such as dates and flour, and high fat such as sour cream and pine nuts. Participants prepared both high calorific meals such as Lazy Cake Recipe (P01_R05, 4.8K calories) and low-calorie meals including Crispy Cucumber Salad (P08_R09, 274 calories).

B.3. Narrations

We collect sentences detailing the actions carried out by participants in each video. Following [18], we instruct the participants to record spoken narrations of all actions throughout the videos and re-purpose the web-based narrator tool used in [29] for expert commentary. Using the tool, participants pause the video and record a narration whenever they observe an action and provide detailed narrations which include at least one verb and noun, and information about *what* they did *why* and *how* they did it when relevant.

This results in a rich set of narrations which are denser and more detailed than previous egocentric datasets. Fig. A2 shows that HD-EPIC narrations are much denser with a higher words per sentence and words per minute than Ego-Exo4D, the next densest dataset, demonstrating a significant increase in detail. We also see increases in the numbers of nouns and verbs per sentence as sentences describe actions in more detail, even common actions are described uniquely—e.g. ‘pick plate from a pile of plates using the right hand’ and ‘pick plate from a lower shelf using the left hand’. In total, 81% of all narrations in HD-EPIC are unique sentences, compared to only 31% in Ego-Exo4D, 22% in Ego4D, and 15% in EPIC-KITCHENS. Non-unique narrations have recently been identified as a concern in recent works [54].

B.4. Post-Processing: Multi-Video SLAM and Gaze

We make use of the Multi-Video SLAM API provided by the Aria Machine Perception Services (MPS) [1] to reconstruct a single point cloud from all recordings obtained in a kitchen, across multiple days. It is important to highlight that no additional scanning is used to reconstruct the scene—it is simply reconstructed from the unscripted recordings. We obtain the 3D point clouds and 2D tracked points along with camera trajectories consisting of 1kHz frequency and 6DoF, estimated using odometry, for all videos. Out of the 156 videos, 3 videos could not be successfully processed via SLAM API. For these videos, we run COLMAP [65, 66] following EPIC Fields [77] pre-processing to obtain camera poses and manually align them to the corresponding kitchen obtained from MPS.

We also use the MPS services to obtain yaw and pitch

³In contrast, Ego-Exo4D recipes are 8.5 mins on average.

Dataset	Val&Test Hours	Action Segments	Unscripted	Free Setting	Recipe	Nutrition	Gaze	Audio Labels	Object	Hand	3D object over time	Labelled 3D environment	Camera pose	Fully annotated
CMU-MMAC [20]	16.6	✓	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗
ADL [55]	5.8	✓	✓	✓	✗	✗	✗	✗	B-Box	✗	✗	✗	✗	✓
UTE [41]	16.9	✗	✓	✓	✗	✗	✗	✗	Polygon	✗	✗	✗	✗	✓
KrishnaCam [71]	14.0	✗	✓	✓	✗	✗	✗	✗	✗	✗	✗	✗	✗	✓
Charades-EGO [69]	6.7 ¹	✓	✗	✓	✗	✗	✗	✗	✗	✗	✗	✗	✗	✓
EGO-CH [61]	26.5	✓	✓	✗	✗	✗	✗	✗	B-Box	✗	✗	✗	✗	✗
MECCANO [62]	3.0	✓	✗	✗	✗	✗	✓	✗	B-Box	B-Box	✗	✗	✗	✓
EGTEA Gaze+ [45]	19.1	✓	✗	✗	✗	✗	✓	✗	✗	Mask	✗	✗	✗	✓
HOI4D [46]	11.4	✓	✗	✗	✗	✗	✗	✗	Mask	Mask	✓	✓	✗	✓
Assembly101 [68]	66.8 ¹	✓	✗	✗	✗	✗	✗	✗	✗	3D pose	✗	✗	✗	✓
EPIC-KITCHENS-100 [18]	25.3	✓	✓	✓	✗	✗	✗	✓	Mask	Mask	✓	✗	✓	✗
Ego4D [28]	288.7 ²	✓	✓	✓	✗	✗	✗	✗	B-Box	B-Box	✗	✗	✗	✗
HoloAssist [80]	49.8	✓	✗	✗	✗	✗	✓	✗	✗	3D pose	✗	✗	✓	✓
Aria Digital Twin [50]	8.1	✗	✗	✗	✗	✗	✓	✗	Mask	✗	✓	✓	✓	✓
Aria Everyday Activities [47]	7.3	✗	✗	✓	✗	✗	✓	✗	✗	✗	✗	✗	✓	✗
Aria Everyday Objects [74]	0.4	✗	✓	✓	✗	✗	✗	✗	B-Box	✗	✓	✗	✓	✓
IndustReal [67]	3.5	✓	✗	✗	✗	✗	✓	✗	✗	3D pose	✗	✗	✗	✓
IT3DEgo [91]	4.6	✗	✓	✓	✗	✗	✗	✗	B-Box	✗	✓	✗	✓	✓
CaptainCook4D [52]	42.5	✓ ³	✗	✗	✓ ⁴	✗	✗	✗	✗	3D pose	✗	✗	✓	✗
Ego-Exo4D [29]	85.1	✗	✓	✓	✓ ³	✗	✓	✗	Mask	3D pose	✗	✗	✓	✗
HD-EPIC	41.3	✓	✓	✓	✓	✓	✓	✓	Mask	Mask	✓	✓	✓	✓

Table A1. Comparing egocentric video datasets. ¹Subset of egocentric videos. ²Episodic Memory, Hands+Object and Forecasting benchmarks. ³Only key steps annotated.

angle for eye gaze and obtain a unit vector along the gaze direction by anchoring it to the central pupil frame. The depth of the gaze vector is estimated by projecting it into the world frame and obtaining its first intersection with any surface on the 3D space. We then use the device calibration and project the ray to obtain 2D location on the image frame. After gaze is calculated, we remove the gaze camera input from all VRS files for anonymity.

We also provide additional information on the anonymisation process for videos. To ensure the privacy of the participants, faces were anonymised in the rare case where they are visible due to reflections or accidental capture. Participants were asked to view their videos and notify us of cases where faces or identifiable information was present, we also anonymised videos which we found to need anonymisation. We then used black boxes to mask these faces/identifiable information to avoid potential regeneration. In total 8 videos were marked for anonymisation (6% of our 156 videos)—we redacted a total of 2774 frames. Correspondingly, we manually marked and anonymised these frames in the two gray-scaled SLAM cameras stored in the VRS files.

We separately post-process the VRS files to convert them to mp4. For the audio, we select channels 5 and 6 to get a binaural audio out of the 7 audio channels in the input recording. We select these channels as they are closer to the ears and remove other audio channels to mitigate breathing noises. We then extract the frames from VRS matching the audio timestamp. Finally, we obtain a video with 30 FPS and binaural audio with 48 kHz sampling rate. These videos are used in all follow-up annotations and benchmarking.

C. Annotation Pipeline

C.1. Annotating Recipe Steps and Ingredients

We match the annotated start-end times to recipes steps to automatically compile ‘recipe videos’, which we use for annotating recipe steps and ingredients. We further link recipes, steps, and ingredients by annotating time segments where ingredients are added to the recipe. For each ingredient, we temporally annotate the weighing and adding actions. The weighing segments allow tracking the amount of each ingredient from zero to the total used in the dish whereas the adding segments mark when an ingredient is incorporated in the recipe. When weighing is performed with the scales, we obtain automatic OCR readings of the scales which are manually verified and checked to ensure the full amount matches the quantity in the nutrition information provided by the participant. Fig. A3 visualises the annotations per ingredient.

C.2. Fine-Grained Actions

C.2.1 Transcriptions

The audio narrations, provided by the participants are first converted into timestamped transcriptions. We use an ensemble of audio transcription models, namely Whisper [60], Qwen-Audio [14] and Nemo [38], to improve the quality of transcriptions. We find that each of these models makes a different set of errors in understanding the diverse spoken accent of participants, so we opt for the consensus of models along with manual checks. For each narration, we select a transcription from one of these models based on a number of checks. Transcriptions with spelling errors or words

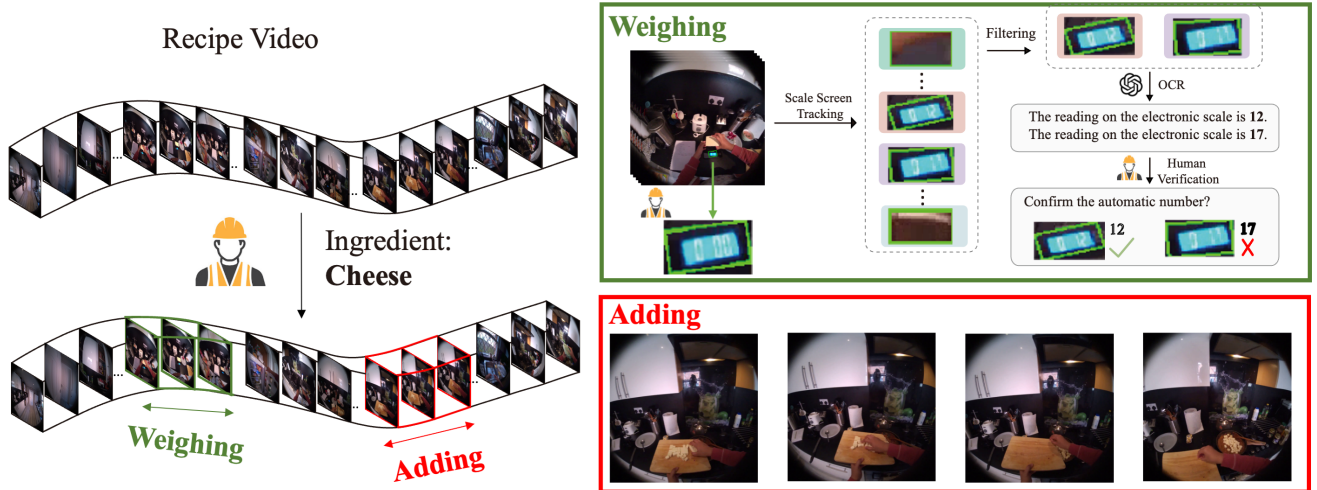


Figure A3. Per ingredient, we temporally label weighing (if present) and adding to recipe. Weighing is OCRed and manually verified.

not in a whitelisted kitchen vocabulary, which we populate during transcription cleaning, are excluded. If there is more than one unique transcription, we select the highest model confidence. If all transcriptions are excluded by the spelling or whitelist checks, but multiple models have produced the same transcription, then we select this transcription. Otherwise the Whisper transcription is selected as a default. We then manually check all candidate transcriptions using dedicated human checkers, who listen to the audio and verify or correct the transcription. We then fix typos and obtain a set of cleaned detailed transcriptions.

C.2.2 Action Annotations

For each narration, we label the precise time segment bounded by when the action in the narration starts and concludes, using AMT workers. We divide the video into HITs of 10 consecutive narrations. Multiple annotators ($K_a \geq 3$) are tasked to annotate both start and end times of the segment containing the narrated action, primed by the timestamp, denoted by $A_i = [t_{s_i}, t_{e_i}]$. The option to skip an action is also provided to the annotators if they are unable to find it in the video; accounting for rare cases when the action occurs off-camera. The interface used to annotate action boundaries is shown in Fig. A4. We ensure we have sufficient inter-annotator agreements between at least three annotators, adding annotators until satisfied.

We select the best annotation by measuring inter-annotator agreement. For narration i , let $A_i(j)$ represent the annotation from annotator j . The agreement score for annotator j on narration i is $\alpha_i(j) = \frac{1}{K_a} \sum_{k=1}^{K_a} IOU(A_i(j), A_i(k))$, where IOU is temporal intersection-over-union. The average agreement score across the HIT is $\alpha(j) = \frac{1}{N_H} \sum_{i=1}^{N_H} \alpha_i(j)$, with N_H ac-

tions per HIT. We select the best annotator \hat{j} as the one with the highest average agreement on their three lowest-scoring actions and similarly find \hat{k} as the second-best annotator.

The ground-truth action segment A_i is then found as:

$$A_i = \begin{cases} \text{Union}(A_i(\hat{j}), A_i(\hat{k})), & \text{if } IOU(A_i(\hat{j}), A_i(\hat{k})) > 0.5 \\ A_i(\hat{j}), & \text{otherwise} \end{cases} \quad (1)$$

where we choose to merge annotations from the two best annotators for an action segment if they have a strong agreement, which helps to avoid overly tight segments [17].

To ensure annotation quality, we require minimum annotator agreement $\min(\alpha(j)) \geq 0.3$ for each HIT. Failed HITs are re-attempted with up to 3 new annotators, then manually annotated if still below threshold. We also perform quality checks on randomly sampled action segments.

C.2.3 Parsing and Clustering

We next parse the open vocabulary narrations so they can be used for tasks that expect a closed vocabulary, such as action recognition. Additionally, we wish to identify the parts of the narration where the participant is detailing the manner in which the action was carried out and the goal achieved by carrying out the action. We achieve these by parsing these open vocabulary narrations as follows. As a result of the fine-grained nature of the collected narrations, it was expected that LLMs with their impressive performance in NLP tasks would provide better results compared to smaller language models. To that end different open source LLMs were tested for the task of extracting the pairs of nouns and verbs and identifying the main action. Initial testing showed mixed results with hallucinations being a major issue. Even though a combination of different


Please mark the start and the end of the action described below.

You might see the same action repeated a few times. Please make sure you annotate each repetition listed on the right.

Please pause the video (shortcut: P) and use the backwards/forwards buttons to be more precise.

Click the 'Start' and 'End' buttons below to mark the start and end of the action.

Note: You cannot directly click on the YouTube video. Please use the buttons provided below to pause, play or seek.



YouTube video player interface with controls: <<< << < play > >> >>>. Video ID: P03 20240218 144430.

- 1) put down plate containing lemon sprinkled with baking soda and seaweed on the kitchen top next to the sink
- 2) with my right hand, open the cupboard above the kitchen hob on the right hand side
- 3) close the cupboard above the sink by pushing the cupboard with my right hand
- 4) close the cupboard above the kitchen hob with my right hand by pushing the door
- 5) pick up dirty saucepan from the kitchen hob using my right hand
- 6) with my right hand pick up dishwasher brush from the sink
- 7) use the dishwasher brush to scrub the saucepan to remove food bits that are still attached to the bottom of the saucepan
- 8) put down the dishwasher brush on the sink by dropping it with my right hand
- 9) pick up sponge from the sink
- 10) because the food bits are hard to remove from the saucepan with the dishwasher brush, i use now the abrasive part of the sponge to scrub the saucepan and remove food bits

1) put down plate containing lemon sprinkled with baking soda and seaweed on the kitchen top next to the sink

Figure A4. Interface for annotating action boundaries

prompting techniques such as few-shot prompting improved the results, given the number of narrations, the results were too unpredictable.

Each narration transcription is parsed into nouns, verbs, and any noted hands (Left/Right/Both) using spaCy's Part-of-Speech tagging [30]. To better identify compound nouns, we supplement spaCy's noun chunks with Constituency Parsing [37], and use pattern matching to filter out irrelevant parts, such as articles. For pronouns (e.g. 'it's'), we use the co-reference resolution within the spaCy pipeline along with heuristics to replace these with corresponding nouns, selecting from the directly preceding narration. We also identify the primary verb-noun pairs, as the 'ROOT' verb in spaCy with its corresponding noun, to form the main action. To enhance accuracy, we also apply heuristics to de-prioritize some verbs (e.g. "use something") from being considered the main action. We show some parsing examples in Table A2.

Similar to the POS tagging, clustering nouns and verbs using LLMs was too noisy and inconsistent. The clusters from [18] were used in addition to a combination of heuristics, word embeddings, and manual assignment to assign any new verbs or nouns to these clusters. We assigned 1,172 new verbs (such as *brew*, *simmer*, *burn*) and 16,812 new nouns (e.g. *mortar*, *nutmeg*, *glue*) to the clusters. We additionally added 9 new verb clusters and 3 new noun clusters

which did not fit in the clusters from [18].

C.2.4 Sound Annotations

We follow a procedure similar to [31] to collect audio annotations. First, we only use a stereo 2-channel audio, extracted from the 5th and 7th microphones. These offer the closest audio input to the two ears and in our experience decreases the presence of breathing sounds in the audio signal, thus better allowing annotators to focus on the sound event. We divide the untrimmed audio, for one video, into chunks based on a silence threshold to provide manageable project sizes to the annotators, whilst avoiding splitting an ongoing sound.

We use a modified version of the VIA Interface [22]. Annotators were tasked with listening to the audio and labelling the start and end times of any distinct audio event they heard. They selected labels from the classes from [31], though annotators could still provide a free-form text description. We post-process the annotations, typo-correcting free-form text descriptions and grouping into the classes from [31]. We did not find a sufficient number of new free-form text descriptions to warrant additional classes, hence the classes match that of [31].

To reduce label noise, we used a customised version of the LISA [4] interface. Here, the annotators were provided

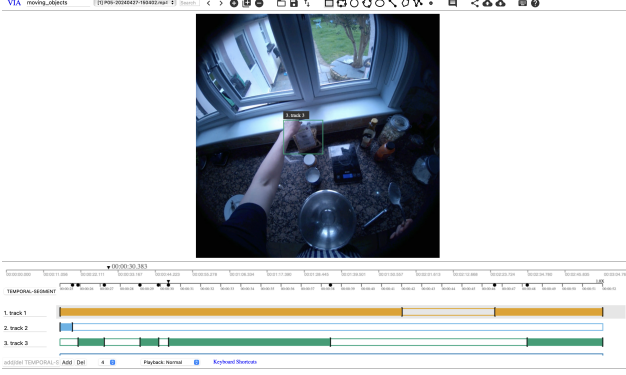


Figure A7. The VIA interface used for annotations object movements in 2D.

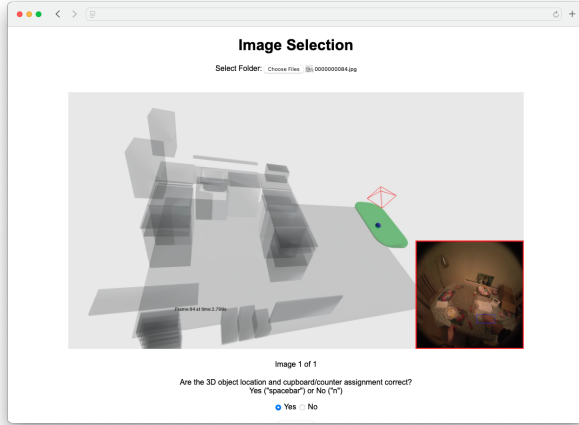


Figure A8. Interface for reviewing object/fixture assignments.

Fig. A7. In addition to the start and end timestamps of the object motion, we ask the annotators to provide bounding boxes around the object at the onset and end of the motion.

Object/fixture review. Fig. A8 shows the interface for manually reviewing object/fixture assignments. Reviewers were asked to verify whether the 3D location and assigned fixture look correct. This process was repeated - after each iteration corrections to the Fixture annotations (*e.g.* including missing bins, hooks *etc.*) were made in Blender, as well as adding heuristics to the assignment process.

Fixture transitions. Frequent transitions are cross-fixtures (*i.e.* excluding objects picked up and placed down on the same fixture), sink→counter and counter→cupboard. Fig. A9 shows common transitions and locations where the 20 most common objects are placed on/in, normalised per object. Again, counters are the most used fixtures.

Long Term Object Tracking. We link object movements over time to provide long term tracking for all objects, *i.e.* object itineraries. These capture sequences of an object’s

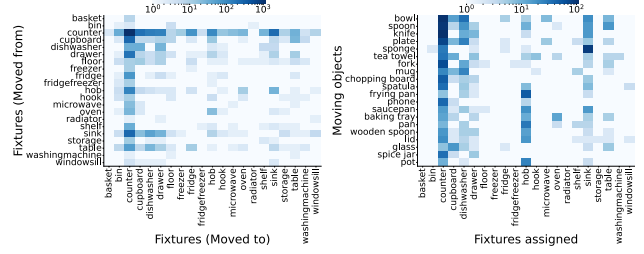


Figure A9. Left: Common source→destination object movements. Right: Common take/put locations of the top-20 moving objects.

movement throughout the video. Fig. A10 shows the object associations interface used to connect objects throughout the video. It utilises our lifted 3D locations to speed up the annotation process. For each track, the annotation interface re-orders objects based on the comparison between the query track’s initial 3D location (pick up) and previous tracks’ final 3D location (put down). This approach was used in [56] for egocentric tracking, and we use it as an initial step in our interface. Around 12% of the objects moved at least 5 times in the video making it suitable for long-term tracking. On average, an object moves $2.4 (\pm 2.6)$ times in HD-EPIC.

Hand and Object Segmentations To improve performance, we first manually annotate a small set of images from each video (5-10 frames) for both left and right hands. To cover the diverse locations and lighting conditions, we use the camera pose estimates to identify clusters of hotspots and select a set of images that represent hands in action at various locations. This is sent to dedicated annotators who manually segment these. We follow the same approach in [19], where hands and arms are segmented jointly to avoid the ad hoc decision of the hand boundary. We use the TORAS [33] annotation tool for all our manual segmentation. We manually annotate 768 frames across 156 videos, which accounts for roughly 0.02% of the total frames.

For each video, we utilise SAM2 [64] to predict 2D hand masks for every frame. The manually annotated frames are placed in the conditional memory (roughly 5 frames per video) whilst the previous 6 predictions are used in the non-conditional memory in order to improve temporal consistency. We notice that SAM2, like most Video Object Segmentation (VOS) models, suffers from identifying the hand late when it enters the scene. In order to address this, we run the aforementioned process both forwards and backwards through the video.

Correcting Noisy Hand Segmentations. We devise a simple, non-learned procedure to automatically detect noisy masks based on the rate of change of the mask’s area over consecutive frames. A major change in size often indicates noise. If the previous frame has no mask and the current mask is large, then it is likely that a frame has missed a

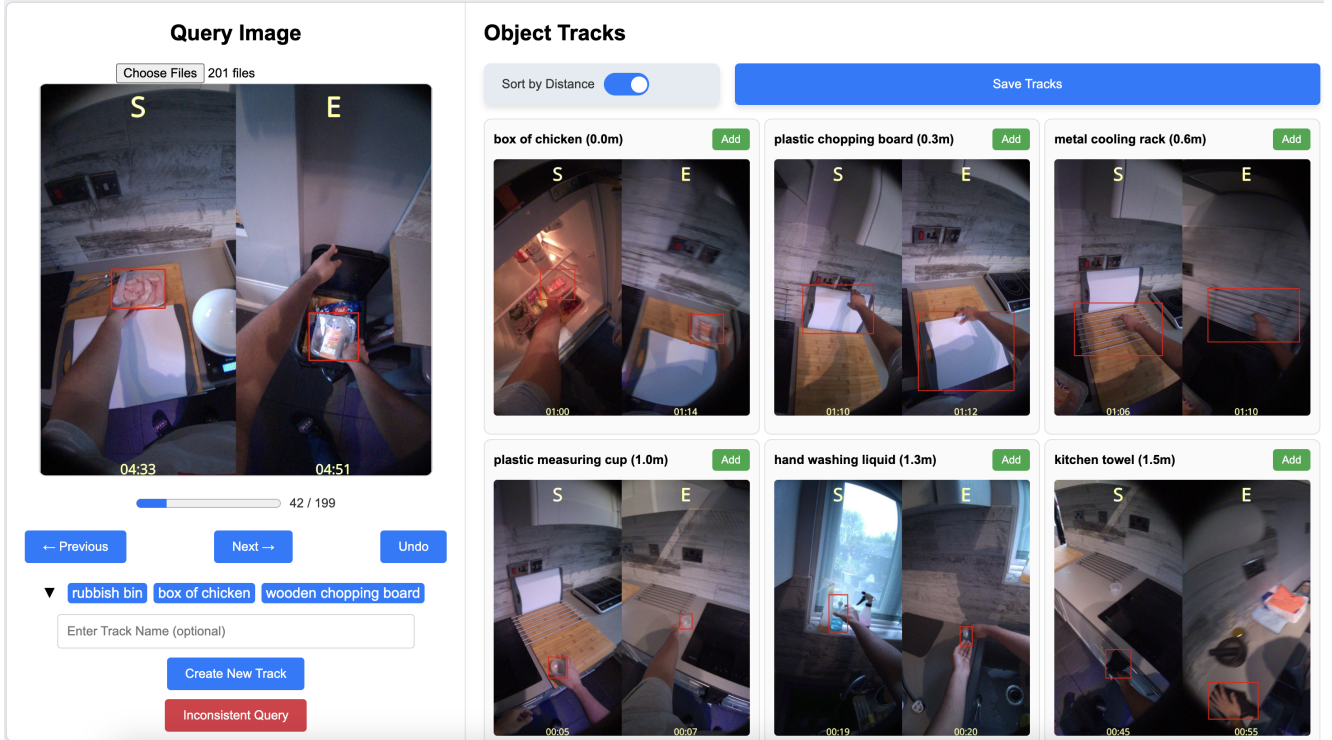


Figure A10. **Object associations interface.** For each query track, a set of possible object names extracted from the narrations is suggested, we provide these names as an additional annotations. Importantly, all previous tracks are sorted by their distance. It is trivial for the annotator to confirm that the track on the left, matches the box of chicken on the right, which has been positioned in exactly the same location (0.0m). By clicking ‘Add’, this track is correctly associated connecting the object movement over time.

mask or under-segmented. We flag frames for noise correction. Additionally, if the previous mask is reasonable but the current mask is much larger or smaller, then this is another indication of errors. We use a dynamic threshold which catches when masks expand or contract too quickly based on the area of the previous mask.

Once we have flagged noisy frames, we sample from these frames for manual annotations. In order to propagate the new knowledge from the manual correction phase, we re-run the SAM2 procedure within a small window of frames around the original error, with the newly corrected frame in the memory (alongside the original 5 ground-truth images). The idea here is to correct frames around the original noise with the new manually annotated frame as a conditioning frame. In total, this correction procedure produces roughly 7.5k right hand masks and roughly 7k left hand masks for manual correction (0.19% of total hand masks). As this is still a rather large amount of annotation to conduct, we further reduce this by uniform temporal sampling, achieving an average of roughly 20 hand side masks per video to be corrected. In total, we manually clean and annotate a further 5163 right hand masks and 4332 left hand masks (0.12% of total hand masks).

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

C.4. Density of Our Annotations

On average, we have 263 annotations per minute of the dataset. Here in Tab. A3, we show the breakdown of this number. Note that this excludes the hand masks that we generate at 30 fps.

D. Benchmark and Results

We first detail how each of the question prototypes in the VQA benchmark are formed in Sec. D.1. We then provide

additional experiments for the VQA demonstrating an ablation of the inputs and the prediction bias of models in Sec. D.2. Sec. D.3 gives additional details of the models used in the VQA benchmark. Finally, we provide additional details of the long-term VOS benchmark and recognition benchmark in Sec. D.4 and Sec. D.5.

D.1. VQA Benchmark Details

In this section, we detail how all question prototypes are formed and sampled.

D.1.1 Recipes

These questions test the ability to understand, retrieve, localize, and recognize recipes and their steps.

Recipe Recognition. To test a model’s ability to recognise entire recipes from long-term video, we formulated the question prototype “Which of these recipes were carried out by the participant?” As input, the model receives every video recorded by a given participant, concatenated chronologically. The positive answer is the name of one of the recipes prepared in these videos. To sample difficult negative answers, we generated recipe names similar to the positive recipe via ChatGPT [5], with the prompt “You are a professional chef. Your goal is to identify 4 similar recipes for each line of the text file. These recipes must technically be different dishes, they should not be different names for the same meal. The recipes can make use of the same ingredients or new ones.” The results were manually edited to ensure validity, which allowed us to create more difficult negatives than were otherwise present in the dataset.

Multi-Recipe Recognition. Extending this, the question prototype “Which of these recipes were carried out in this video?” tests a model’s ability to recognise every recipe occurring within a video. As input, the model now receives an entire video containing one or more recipes. The positive answer is the names of all recipes occurring in the video, with the negatives being generated in the same way as in **Recipe Recognition** above.

Multi-Step Localisation. With the question prototype “In this video, when did the participant perform each of the following <recipe name> recipe steps: “<recipe step a>”, “<recipe step b>”, “<recipe step c>”?”, we then test a model’s ability to localise a subset of time segments of 3 steps of a recipe in a video. One video is supplied as input and in the question text we supply the text for 3 recipe steps selected in a random order. The positive answer is a list of the largest segment start and end times for the 3 steps from a given recipe that occur within the video. To generate negative answers, we first collect all the time segments for recipe steps that don’t correspond to the recipe in the question text (if any exist). Then we concatenate these with all the time segments for recipe prep

occurring in the video. Finally, for each negative answer we randomly sample 3 segments.

Step Localisation. Similarly, the question prototype “When did the participant perform step <recipe step> from recipe <recipe name>?” assesses a model’s ability to localise all major segments corresponding to a recipe step. Again, the input is one video which must contain ≥ 5 unique recipe steps. The positive answer is the list of all annotated start and end times for a randomly chosen recipe step. In the question text, the name of the recipe and the recipe step text is provided. To form negative answers, we randomly sample other recipe steps from the same video and use the start and end times of the corresponding time segments.

Prep Localisation. Likewise, the question prototype “When did the participant perform prep for <recipe step> from recipe <recipe name>?” assesses a model’s ability to localise all major segments corresponding to prep for a given recipe step. Again, the input is one video which must contain ≥ 5 unique recipe steps. The positive answer is the list of all annotated start and end times for the prep relating to a randomly chosen recipe step. In the question text, the name of the recipe and the recipe step text is provided. To form negative answers, we randomly sample other recipe steps from the same video and use the start and end times of the corresponding prep time segments.

Step Recognition. As recipes are comprised of many high-level steps, the question prototype “What step did the participant do between <TIME HH:SS:MM video 1> and <TIME HH:SS:MM video 1> in this video?” tests a model’s ability to recognise these long activities. One video is supplied as input, which must contain at least 5 unique recipe steps. The positive answer is randomly sampled from the steps occurring in the video. In the question text, the provided start time is the start of longest segment relating to the positive step and the end time is the end of the longest segment. To form negatives, we randomly sample recipe steps from the same video as the positive answer.

Rough Step Localisation. Where the prior question prototypes consider the localisation of all segments of a recipe step/prep, “Which of these time segments belongs to the <recipe name> recipe step <recipe step> in this video?” tests whether a model can localise recipe steps from a subset of the possible time segments. Again, the input is one video with ≥ 5 unique recipe steps. Here, the positive answer is the start and end time of the longest time segment relating to the recipe step. The negative answers are formed by randomly sampling recipe steps from the same video and using their corresponding start and end times of the longest segment.

Following Activity Recognition. Finally, we test a model’s ability to recognise high-level activities that follow recipe steps with the question prototype “Which high-level

activity did the participant do while completing recipe step `<recipe step>` in this video?” As input, the model receives one video and the recipe step text in the question. The positive answer is determined by finding the high-level activity that occurs while the given recipe step is concluding. Negative answers are sampled from other activities occurring in the video that do not overlap with the end time of the given step.

All of these question prototypes were sampled uniformly at random from the space of possible positive answers. However, as some simpler recipes occur more frequently than others, these were excluded from the recipe recognition and step localisation tasks.

D.1.2 Ingredients

The ingredient questions assess a model’s ability to understand and identify the ingredients used, their exact amounts and in what order they are added.

Ingredient Retrieval. We test a model’s ability to recognise added ingredients, with the question prototype “Between `<TIME HH:SS:MM video 1>` and `<TIME HH:SS:MM video 1>`, which of these ingredients were added to the dish being prepared?” One video is used as input, with the start and end time of an ingredient-adding segment supplied in the question text. Instead of using the exact annotated times, we pad the start and end times by adding 5 seconds to either side of the segment. If the new times overlap with the adding of a new ingredient, we instead add the smallest possible time that avoids the introduction of a new ingredient to the segment. The positive answer is randomly selected from the ingredients added in the video. We generate negatives by randomly selecting from ingredients that are added in the video and do not overlap with the positive answer.

Ingredient Weight. This question tests the ability to identify ingredient weight with the prototype “How much did the participant weigh of `<ingredient>` in this video?”. The input is a video segment annotated as the weighing sequence for a given ingredient and the positive is the correct weight of this ingredient. It should be noted that humans can trivially perform this task by reading the scale while current models struggle. The negatives are generated by sampling a random multiplier between (0.5, 5) for the positive answer.

Ingredient Order. This question tests the ability to identify the order in which ingredients are added with the prototype “What is the order of ingredients added to the dish in this video?”. The input is a video segment where at least three ingredients were added to the dish and the positive answer is the list of the added ingredients in the correct order. For recipes with more than five ingredients, a subset of five consecutive ingredients was sampled randomly. The nega-

tives are random permutations of the positive order.

Ingredient Adding Localisation. We assess the ability to temporally localise when ingredients are added with “When was ingredient `<ingredient>` added to recipe `<recipe name>`?”. The input is an untrimmed video and positive is the start and end time of the annotated adding segment. Negative answers are the start and end time of *action* segments from the input video where the narration mentions the ingredient. Action segments overlapping with the adding time segment are discarded. If less than 4 action segments include the ingredient, we use action segments containing similar ingredients. Similarity is calculated using cosine distance > 0.5 of spaCy [30] (model `en_core_web_lg`) word embeddings.

Ingredient Recognition This tests the ability to identify the ingredients used with “Which of these ingredients is used in `<recipe name>`?”. The input is all videos of one participant, with one of the completed recipes selected as `<recipe name>`. Only recipes with at least four ingredients were considered. The positive answer is obtained by randomly sampling a single ingredient from the ingredients used. The negatives were generated using an LLM which is asked to provide potential ingredients for a given recipe that do not include the ingredients used in the dish. This was done to ensure more challenging negatives than those sampled from other recipes in the dataset. We also created a negative version of this question: “Which of these ingredients is not used in `<recipe name>`?” was also created. Here, the input video selection is the same and positive and negative selection strategies are swapped.

Exact Ingredient Recognition This question tests the ability of a model to identify the exact quantity of an ingredient with “What was the exact quantity of `<ingredient>` used in `<recipe name>`?”. For this question, the input is all videos relating to a randomly selected recipe. The positive answer is randomly selected from the used ingredients while excluding ingredients used in the **Ingredient Weight** question. The negatives are generated by sampling a random multiplier between (0.5, 5) applied to the positive answer.

To select final questions for our benchmark we uniformly sample over participants.

D.1.3 Nutrition

The nutrition questions assess a model’s ability to understand the nutritional values of ingredients used in the videos.

Image Nutrition Estimation. Our first question is “Which of the ingredients in these images showcase higher `<nutrition>`?”, where `<nutrition>` is one of “calories”, “fat”, “carbs” or “protein”. The input is 5 video frames showing different ingredients. To obtain these we manually label a frame clearly showing the ingredient.

We search this frame within the ingredient-adding sequence as this allows a clear view of the correct quantity of the ingredient. In total we annotate 200 frames/ingredients, from which we generate 834 questions. The positive is a randomly selected ingredient. To sample the 4 negatives and their corresponding frame, we first calculate the value-to-amount ratio of each ingredient, *e.g.* if we have 30g of protein in a 100g ingredient, the ratio is 0.3. For a nutritional value we form negatives by finding ingredients that have a *higher* value-to-amount ratio but a *lower* nutritional value compared to the positive answer. This strategy samples hard negatives that are not answerable by text alone. To ensure a reasonable separation between positive and negative answers, the difference in nutritional value is at least 20% of the positive answer’s value.

Nutrition Change. We use “From <TIME HH:SS:MM video 1> to <TIME HH:SS:MM video 1>, what changed in the nutrition values of the dish with recipe <recipe name>?”. The input is a video containing a recipe (<recipe name>) and one or more adding segments. <TIME HH:SS:MM video 1> placeholders are replaced with the annotated start-end times of an ingredient-adding segment. Answers are in the format “calories changed by <a>, fat changed by , carbs changed by <c>, protein changed by <d>”. For the positive <a/b/c/d> are replaced with the correct ingredient nutritional values. For the negatives, the placeholders are multiples of the correct nutritional value. The random multiplier is sampled within (0.10, 0.75) or (1.25, 1.75).

Video Nutrition Estimation. Finally, we also use the prototype “What is the ingredient with the highest <nutrition> in this recipe?”, where <nutrition> is the same as above. For this question, the input is an untrimmed video containing one recipe. The positive answer is the ingredient with the highest calories/fat/carbs/protein in the recipe. Negative answers are sampled from ingredients that have a value-to-amount ratio similar to the positive. At least 2 negative answers ingredients come from the same recipe, while the rest can be sampled from other recipes when it is not possible to obtain enough negative answers from a single recipe.

Final benchmark questions are randomly sampled from the set of possible questions uniformly across all participants to ensure diversity.

D.1.4 Fine-Grained Action

With these questions, we assess a model’s ability to understand detailed fine-grained actions as described by the participants in our ground-truth narrations.

Action Recognition. We first formulate the question prototype “Which of these sentences best describe the ongoing action(s) in the video 1?”. The input is a short video clip

containing 1-3 actions. Positives contain the correct 1-3 action descriptions obtained from the narration transcriptions. We sample negatives by selecting other sequences of the same number of actions from other video clips in the dataset. We ensure difficult narrations by requiring negatives to have at least one verb and noun in common with the positive and prioritise those with more common verbs and nouns.

How Recognition. To assess the model’s ability to understand how fine-grained actions are performed, we use the prototype “What is the best description for how the person carried out the action <verb noun> in this video segment?”. We randomly sample a narration annotated with start and end times. The input is a video clip of this action as given by the start-end time annotations and <verb noun> in the question is the main action extracted from the narration by our parsing. The positive answer is the manner of the target action as provided by the ground-truth narration also extracted by the parsing. To sample negative answers, we first identify other narrations that share both the same verb and noun cluster (see Sec. B.3) with the main action of the target narration. We then extract ‘how’ from these narrations using our parsing. To avoid false positives, we filter out any ‘how’ with verb and noun clusters which are subsets or supersets of the positive answer and manually filter the remaining questions.

Why Recognition. To assess the model’s ability to understand why fine-grained actions are performed, we use “What is the best description for why the person performed the action <verb noun> in video 1?”. The input is a video clip for a randomly sampled narration and <verb noun> is the main action parsed from the narration. The positive answer contains the reason for performing the target action as extracted from the narration by our parsing. To sample negatives, we use ‘whys’ extracted from other narrations with the same main action for at least 2 negatives. When there are not enough of these negatives we add whys from actions that share the same verb. To avoid false positives, we remove any negative ‘why’ with verb and noun clusters which are subsets or supersets of the positive answer and manually filter the remaining questions.

Action Localisation. To assess ability to localize fine-grained actions, we created the question prototype “When did the action <verb noun> happen in the video 1?”. The input is a 30 second to 15 minute video clip. Answers are in time period format (*e.g.*, <TIME HH:MM:SS.MS> to <TIME HH:MM:SS.MS>). The positive answer corresponds to the time period when a randomly sampled target action occurs, while the negative answers are time periods when different actions occur. To make the questions challenging, at least two of the negative timestamps are from actions that share the same noun as the target action but differ in the verb. To ensure 4 negatives we add timestamps for

actions that share the same verb but have a different noun, or those that share either noun or verb cluster.

D.1.5 3D Perception

The 3D perception questions assess the model’s capability to understand 3D environment or 3D movements from only video.

Fixture Location. These questions test the ability to understand a fixture’s relative location in 3D. We use the fixture mesh annotations for each kitchen to generate questions of type “Given the direction I am looking at <TIME HH:MM:SS video 1>, where is the <query fixture> located?” We select the <query fixture> as one of the unique fixtures in the kitchen *e.g.* fridge, hob, sink, microwave. For the positive answer, we calculate the angle between the direction where the camera is pointed and the direction of the <query fixture> from the camera location. We represent this angle as position of the hour hand on the clock *e.g.* 2 o’clock, 7 o’clock. For the negatives, we randomly sample from other possible positions of the hour hand on the clock. We don’t choose negatives with an absolute ≤ 1 o’clock difference to the positive, in case it might lead to confusion, *e.g.* if the positive is 2 o’clock, then the negatives can be sampled from 4 o’clock to 12 o’clock. For the timestamp in the question, we ignore frames that have a tilt greater than 45° in the camera because it would lead to incorrect calculation of the positive.

Object Location. Here we test the model’s ability to understand object location as objects move with the prototype “Where did I take/put the object identified by <BBOX Y1 X1 Y2 X2> at <TIME HH:MM:SS video 1> from before/after putting/taking it at <TIME HH:MM:SS>?” The bounding box and timestamp in the question is directly taken from the annotations of moving objects in 2D (Sec. C.3). If the question asks where the object is put down, the bounding box and timestamp are when the object is picked up before that. If the question asks where the object is taken from, the box and timestamp are from when the object is put down after that. For the positives, we use the name of the correct fixture associated with the object bounding boxes at the time of pick up/put down. For the negatives, we sample different fixtures that other moving objects have been associated with in the video. We use relative location to another unique fixture *e.g.* sink to identify fixtures without a unique name *e.g.* the cupboard top left of the sink.

Object Contents Retrieval. Here we test the model’s ability to understand 3D locations of fixtures as well as recognise objects that were put/taken from the fixture. We use the fixture meshes to generate questions of type “Which of these objects did the person take/put from/in/on the item indicated by bounding box <BBOX Y1 X1 Y2 X2> in

<TIME HH:MM:SS video 1>?” We randomly sample a fixture in the kitchen and project the 3D vertices of its mesh onto the image plane. We use the maximum and minimum of the projected vertices to find the 2D bounding box. We repeat this for all frames in the video and search for frames where the sampled fixture is occluded less than 30% and the projected bounding box lies fully within the image. We then choose one of the identified frames and the corresponding bounding box for the question. We use the long-term object track and fixture transition annotations Appendix C.3 to identify objects that were either taken from or put on/in the fixture in question. We treat these objects as positive objects and sample few of them as our positive answer. For negatives, we sample objects that were taken from or put on the same fixture in different videos but do not overlap with the positives. We also sample objects that were assigned to other fixtures as negatives.

Fixture Interaction Counting. Here we assess the model’s ability to remember the interactions between participants and fixtures with the question “How many times did I open/close the item indicated by bounding box <BBOX Y1 X1 Y2 X2> at <TIME HH:MM:SS video 1>?” The timestamp and bounding box for the query fixture is obtained following the same procedure as for **Object Retrieval** questions. For the positive, we use the correct number of times the fixture was opened or closed. We obtain this by first using the parsed narrations to identify phrases relevant to fixture interactions, *i.e.* where the verb includes ‘open’ or ‘close’ and the noun includes one of the fixtures. For relevant narrations, we use a 1s temporal window before the corresponding narration timestamp to find the fixture with the highest cumulative gaze intersection. This identifies the exact fixture that is being interacted with at that time. The negatives are sampled randomly from ± 4 of the positive.

D.1.6 Object Motion

Here, we test the model’s ability to correctly reason over various moving objects given a full video as input. We use the long-term object tracking annotations mentioned in Sec. C.3 for questions in this category.

Object Movement Itinerary. We design the question as “Where was the object <BBOX X1 Y1 X2 Y2> seen at time <TIME HH:MM:SS video 1> moved from/to throughout the video?”. We provide bounding boxes around the object instead of giving the object name as text. As positives, we selected the correct order of the fixture names involved as either source or destination in all trajectories of the highlighted object in the query video *e.g.* “from fridge to sink, then from sink to microwave”. As negatives, we select the fixtures associated with the itineraries of other objects that are moving in the video. For fixtures that do not have

a unique name *e.g.* cupboards, we follow the same procedure as **Object Location** and generate descriptions of the relative position.

Object Movement Counting. We begin by testing the ability to count how many times an object moves with the prototype “How many times did the object <BBOX Y1 X1 Y2 X2> seen at <TIME HH:MM:SS video 1> change locations in the video?”. The input is the full video, and the object is specified by its bounding box at a particular timestamp from the moving objects in 2D annotations (Sec. C.3). Positive answers are derived by counting distinct movements using the long-term object tracking annotations in Sec. C.3. To mitigate rapid movements and erroneous annotations, we retain only objects satisfying two criteria: (1) consecutive tracks must be separated by at least 1.1 seconds, and (2) the spatial displacement between the end of one track and the start of the next must not exceed 20 cm. For negative answers, we randomly sample values around the positive answer. To avoid text-based shortcuts, the sampling window is varied randomly from [1,8], and the maximum negative value is restricted to 8. We additionally enforce that the negative values for cases having more than 2 movements should be at least ± 2 value apart from the ground truth answer.

Stationary Object Localization. Here, we evaluate the capability to perform the reverse task: determining when an object remains static in the video. The question is: “After the object <BBOX Y1 X1 Y2 X2> seen at <TIME HH:MM:SS video 1> is first moved, from which of the following starting times does the object remain static for more than <X> seconds?”. Similar to the previous task, the input is the entire video, and the object is specified by its bounding box at a given timestamp. The static duration is determined using long-term object tracking annotations (Sec. C.3). We find t_1 and t_2 , the two longest static durations of an object. A value is then randomly sampled from the range $(t_2 + 5, t_1)$ to fill the placeholder <X>, which specifies the required static duration for the question. We also use the criteria mentioned in the previous question to filter objects with rapid movements and spatial inconsistency across tasks. For the positive answer, a timestamp is randomly selected from the interval starting at the beginning of the maximum static duration and ending at a point determined by subtracting the previously chosen static duration from its endpoint. Negative answers are randomly sampled from the video from the start of the video to the last movement of the object, ensuring no overlap with the positive. All answers are provided as timestamps formatted as <TIME HH:MM:SS video 1>.

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	38.5

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.

D.1.7 Gaze

The Gaze questions evaluate the model’s ability to understand the subject of the camera wearer’s visual attention, which can also be a signifier for future interactions.

Gaze Estimation. We assess the model’s capabilities to understand where the camera wearer is fixating their gaze within a given video clip. We formulate this question as “What is the person looking at in this video segment?”. The input is a trimmed video clip where the camera wearer fixates on some large object, *i.e.* a landmark, for at least 0.5 seconds. Positives are the name of the landmark being fixated on, negatives were randomly sampled from all other landmarks visible in the video segment. To determine if a landmark is visible, we consider all 8 vertices of the 3D bounding box, as well as the centre of the object. An object is considered visible if, in any frame, 5 of the 9 points pass three checks: 1) the point projects onto the camera plane, 2) the point is in front of the camera and 3) the ray from the point to the camera centre does not intersect with another landmark’s bounding box *i.e.* it is not occluded. Common objects are given unique names according to relative positions similar to **Object Location**. Overall, there are 1220 possible questions, of which we randomly sample 1000.

Interaction Anticipation. We also evaluate how effectively a model can anticipate the next object to be interacted with. We format the question as “What object will the person interact with next, ignoring ongoing interactions?” and provide a 10 second video segment, concluding 0.3 seconds after an object is primed for pick-up using eye gaze. For these questions, the positive answer is the object which has just been primed, whilst the negatives are randomly sampled from all other objects that have been moved within the 2 minute video segment. We generate 1110 total questions and randomly select 1000.

D.2. Additional VQA Experiments

VQA Input Ablation. To explore how much information is contained in different inputs (video, ground-truth narrations, questions and answers) we ablate them in Tab. A4.

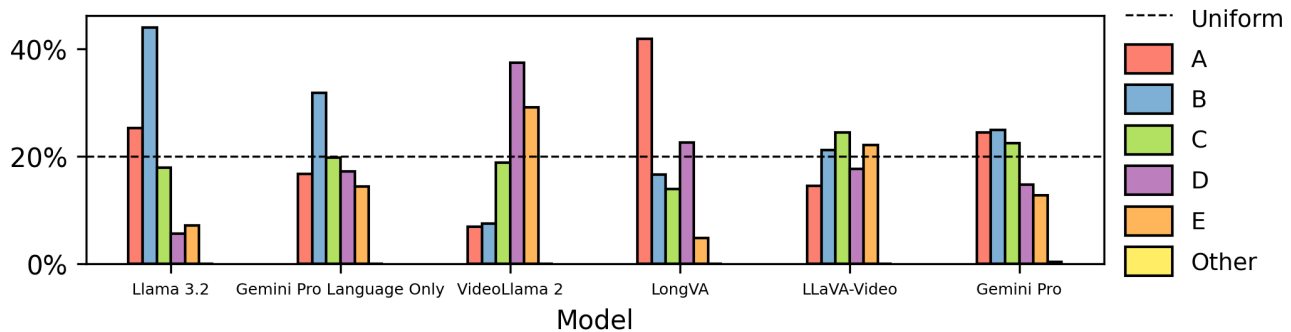


Figure A11. **Prediction Bias of Models.** Most models have a bias in answer, although it is different for each model.

For both Llama and Gemini, providing only the possible answers gives random performance. Recipe and ingredient prototypes are the exception, likely due to textual recipes within the training data. Giving the question improves performance slightly as question contain contextual information, e.g. **Why Recognition** provides the action corresponding to the possible reasons why.

Making predictions from the full ground-truth (GT) narrations is better than using GT action clusters for all categories demonstrating the benefit of more detailed short-term understanding. Using GT narrations instead of visual input gives better performance for all question categories, highlighting the challenge of visual understanding in HD-EPIC. Having GT narrations helps most with recipe and fine-grained action and object motion questions but far from solves these categories e.g. fine-grained action is 63.6%. Regardless, nutrition, 3D perception, object motion, and gaze still have low performance with GT narrations demonstrating their difficulty.

Prediction Bias of Models. Fig. A11 shows the distribution of output predictions for the 5 models tested. All models were able constrain their answers to the 5 options in almost all cases. Some models have clear significant bias (e.g. Llama 3.2 for option B and LongVA for option A).

Per Prototype Results In Table A6 we show the numerical values for the per-prototype results shown in Fig. 11.

D.3. VQA Model Details

Gemini Pro [75]. Version “gemini-1.5-pro-001” was used as it is the latest version which supports context caching - necessary for asking repeated questions about long videos to keep costs tractable. Gemini processes all videos at 1fps regardless of length, so we passed all videos at 768 x 768 at their default speed up to 6000s long, after which we re-encoded the videos at a higher speed to fit in the model’s 2M context window. In these cases we scaled all timestamps in questions to match this new speed to maintain the 1fps/time dependency noted in the documentation. The fol-

lowing prompt was used:

You are an expert video analyzer, and your job is to answer the multiple choice question by giving only the letter identifying the answer. Do not give any other information. For example, acceptable answers are ‘A’ or ‘B’ or ‘C’ etc. You must give an answer, even if you are not sure. Bounding boxes are in the format (ymin, xmin, ymax, xmax) relative to an image size of 1000x1000.

0.2% of questions were blocked/refused by the API due to safety reasons. These were counted as failure cases.

Llama. We used version “Llama-3.2-90B-Vision-Instruct” due to it being the largest and most up-to-date model we could run inference on. The model could run a majority of the questions, though occasionally when including the narrations, we had to evenly sample the input to be a maximum of 120K input tokens. The system prompt matched Gemini Pro.

LongVA [89]. Version “LongVA-7B-DPO” was used. We used the default resolution of 384x284. Videos were processed at 8fps up to a maximum context window of 400 frames (as large as would fit in VRAM with the model). After this, frames were uniformly sampled. The same prompt was used as for Gemini Pro.

VideoLLaMA 2 [13]. Version “VideoLLaMA2-7B-16F” was used. Videos were processed at 336x336 resolution with timestamps scaled similarly to Gemini Pro, and with the maximum supported context window of 16 frames. The same prompt was used as for Gemini Pro. In the case of multi-video input, if input exceeded the frame limit and one or more videos were not used, the prompt included: *N input videos have been concatenated together and the remaining videos have been truncated due to input length limit.*

LLaVA-Video [90]. Version “LLaVA-Video-7B-Qwen2” was used. Videos were processed the same as LongVA, but with the maximum number frames set to 128.

Sample Human Baseline. We also provide a sample human baseline to further assess the gap in understanding between humans and SOTA video-language models. For this

Dataset	Sequences	Avg. Duration	Total Duration	Objs/Seq	Annotated Frames
DAVIS [57]	30	2.8s	83.3s	2.0	1,999
YouTube-VOS [84]	507	21.9s	11,097.0s	2.1	13,710
Ours	1000	561.0s	561,034.1s	6.3	20,548

Table A5. Our HD-EPIC long term VOS benchmark compared to the validation sets of DAVIS and YouTubeVOS, where duration is measured in seconds and calculated using the first and last frame indices of the annotated frames.

we sampled 20 questions per question prototype, 600 questions in total, split between 3 participants.

D.4. Long-Term VOS Benchmark

The details of our long-term VOS benchmark are compared against the validation sets of two popular VOS benchmarks, DAVIS 2017 [57] and YouTube-VOS 2019 [84]. Our benchmark includes 1000 sequences, with an average duration of 561.0 seconds. This is calculated using the first and last frame indices of the annotated frames. When calculating durations for the DAVIS and YouTube-VOS we assume FPS values of 24 and 6, respectively, as stated in the original papers [53, 83]. It is clear that our benchmark provides a much longer temporal duration for video object segmentation evaluation. Furthermore, the table shows that the average number of objects in each sequence is nearly tripled compared to the prior benchmarks.

We evaluate three models on our long-term VOS benchmark (Static, SAM2 [64] and Cutie [12]) and describe the results in the table of Fig. 14. For SAM2 and Cutie, we insert each object’s first appearance frame into the working memory of the model and use the remaining frames as evaluation frames. For the method labelled Static, we copy the mask from the first appearance to all evaluation frames, to act as a naive baseline. Following [53], we use jaccard index \mathcal{J} , contour accuracy \mathcal{F} and their average $\mathcal{J}\&\mathcal{F}$.

Due to the large number of frames in our dataset, we only pass the memory and evaluation frames into the model. We acknowledge that this is limited and the results could be improved further by sampling more frames.

D.5. Recognition Benchmarks

Both action and sound recognition remain fundamental downstream tasks for video models. We thus evaluate strong video models for both tasks in this section. their performance in the following. Note that for all works we follow the data (frames/audio) preprocessing as described in the original papers.

Action Recognition. We assess 5 action recognition methods [9, 24, 26, 51, 76], using publicly available weights finetuned on EPIC-KITCHENS-100. We detail the models and For VideoMAE-L, we used the finetuned model from [9]. As customary, we perform test augmentations over 10 clips

during inference. The exception of this is TIM which also includes the audio modality, and averages the predictions are averaged across all context windows in which the action appears. We also include a chance baseline, where we randomly shuffle the ground-truth labels and compute the accuracy, giving a lower-bound for the action recognition challenge.

	Recipe Recognition	Multi-Recipe Recognition	Multi-Step Localisation	Step Localisation	Prep Localisation	Step Recognition	Rough Step Localisation	Following Activity Recognition	Ingredient Retrieval	Ingredient Weight	Ingredients Order	Ingredient Adding Localisation	Ingredient Recognition	Exact Ingredient Recognition	Image Nutrition Estimation	Nutrition Change	Video Nutrition Estimation	Action Recognition	How Recognition	Why Recognition	Action Localisation	Fixture Location	Object Location	Object Contents Retrieval	Fixture Interaction Counting	Object Movement Itinerary	Object Movement Counting	Stationary Object Localisation	Gaze Estimation	Interaction Anticipation
Blind - Language Only																														
Llama 3.2	38.0	24.0	44.0	16.0	23.0	23.0	22.0	78.0	18.0	12.0	32.0	28.0	40.0	20.0	22.0	20.0	68.0	20.1	30.8	21.8	20.6	21.4	25.4	11.5	31.0	8.6	35.5	32.5	17.3	21.7
Gemini Pro	54.0	54.0	42.0	20.0	21.0	15.0	22.0	76.0	18.0	28.0	38.0	21.0	28.0	28.0	26.0	14.0	50.0	21.3	24.6	21.6	20.7	20.4	33.4	9.0	23.0	31.2	18.5	33.5	21.2	19.8
Video-Language																														
VideoLlama 2	22.0	52.0	18.0	38.0	13.0	18.0	21.0	64.0	19.0	30.0	20.0	27.0	26.0	32.0	24.0	20.0	54.0	30.9	25.2	32.2	20.7	18.8	31.0	35.5	17.7	11.0	44.0	30.5	30.0	12.4
LongVA	14.0	44.0	36.0	18.0	18.0	26.0	19.0	62.0	25.0	24.0	44.0	42.0	30.0	20.0	25.0	22.0	54.0	36.9	28.4	37.0	20.5	26.6	41.2	31.5	32.3	10.2	34.5	23.5	36.0	13.0
LLaVA-Video	28.0	68.0	44.0	20.0	21.0	23.0	24.0	62.0	22.0	36.0	38.0	41.0	36.0	28.0	28.0	26.0	62.0	58.6	41.4	51.2	20.9	21.8	30.6	40.5	16.3	9.8	20.0	27.0	47.5	11.1
Gemini Pro	42.0	76.0	88.0	70.0	35.0	45.0	74.0	54.0	49.0	46.0	56.0	62.0	36.0	28.0	26.0	16.0	62.0	49.3	35.6	43.2	30.3	20.8	32.4	41.5	35.3	18.0	13.0	31.5	36.5	20.8

Table A6. Model results per question prototype