

# Ego-Exo4D: Understanding Skilled Human Activity from First- and Third-Person Perspectives

Kristen Grauman<sup>1,2</sup>, Andrew Westbury<sup>1</sup>, Lorenzo Torresani<sup>1</sup>, Kris Kitani<sup>1,3</sup>, Jitendra Malik<sup>1,4</sup>, Triantafyllos Afouras<sup>\*1</sup>, Kumar Ashutosh<sup>\*1,2</sup>, Vijay Baiyya<sup>\*5</sup>, Siddhant Bansal<sup>\*6,7</sup>, Bikram Boote<sup>\*8</sup>, Eugene Byrne<sup>\*1,9</sup>, Zach Chavis<sup>\*10</sup>, Joya Chen<sup>\*11</sup>, Feng Cheng<sup>\*1</sup>, Fu-Jen Chu<sup>\*1</sup>, Sean Crane<sup>\*9</sup>, Avijit Dasgupta<sup>\*7</sup>, Jing Dong<sup>\*5</sup>, Maria Escobar<sup>\*12</sup>, Cristhian Forigua<sup>\*12</sup>, Abrham Gebreselasie<sup>\*9</sup>, Sanjay Haresh<sup>\*13</sup>, Jing Huang<sup>\*1</sup>, Md Mohaiminul Islam<sup>\*14</sup>, Suyog Jain<sup>\*1</sup>, Rawal Khirodkar<sup>\*9</sup>, Devansh Kukreja<sup>\*1</sup>, Kevin J Liang<sup>\*1</sup>, Jia-Wei Liu<sup>\*11</sup>, Sagnik Majumder<sup>\*1,2</sup>, Yongsen Mao<sup>\*13</sup>, Miguel Martin<sup>\*1</sup>, Effrosyni Mavroudi<sup>\*1</sup>, Tushar Nagarajan<sup>\*1</sup>, Francesco Ragusa<sup>\*15</sup>, Santhosh Kumar Ramakrishnan<sup>\*2</sup>, Luigi Seminara<sup>\*15</sup>, Arjun Somayazulu<sup>\*2</sup>, Yale Song<sup>\*1</sup>, Shan Su<sup>\*16</sup>, Zihui Xue<sup>\*1,2</sup>, Edward Zhang<sup>\*16</sup>, Jinxu Zhang<sup>\*16</sup>, Angela Castillo<sup>12</sup>, Changan Chen<sup>2</sup>, Xinzhu Fu<sup>11</sup>, Ryosuke Furuta<sup>17</sup>, Cristina González<sup>12</sup>, Prince Gupta<sup>5</sup>, Jiabo Hu<sup>18</sup>, Yifei Huang<sup>17</sup>, Yiming Huang<sup>16</sup>, Weslie Khoo<sup>19</sup>, Anush Kumar<sup>10</sup>, Robert Kuo<sup>18</sup>, Sach Lakhavani<sup>5</sup>, Miao Liu<sup>18</sup>, Mi Luo<sup>2</sup>, Zhengyi Luo<sup>3</sup>, Brighid Meredith<sup>18</sup>, Austin Miller<sup>18</sup>, Oluwatumininu Oguntola<sup>14</sup>, Xiaqing Pan<sup>5</sup>, Penny Peng<sup>18</sup>, Shraman Pramanick<sup>20</sup>, Merey Ramazanova<sup>21</sup>, Fiona Ryan<sup>22</sup>, Wei Shan<sup>14</sup>, Kiran Somasundaram<sup>5</sup>, Chenan Song<sup>11</sup>, Audrey Southerland<sup>22</sup>, Masatoshi Tateno<sup>17</sup>, Huiyu Wang<sup>1</sup>, Yuchen Wang<sup>19</sup>, Takuma Yagi<sup>17</sup>, Mingfei Yan<sup>5</sup>, Xitong Yang<sup>1</sup>, Zecheng Yu<sup>17</sup>, Shengxin Cindy Zha<sup>18</sup>, Chen Zhao<sup>21</sup>, Ziwei Zhao<sup>19</sup>, Zhifan Zhu<sup>6</sup>, Jeff Zhuo<sup>14</sup>, Pablo Arbeláez<sup>†12</sup>, Gedas Bertasius<sup>†14</sup>, David Crandall<sup>†19</sup>, Dima Damen<sup>†6</sup>, Jakob Engel<sup>†5</sup>, Giovanni Maria Farinella<sup>†15</sup>, Antonino Furnari<sup>†15</sup>, Bernard Ghanem<sup>†21</sup>, Judy Hoffman<sup>†22</sup>, C. V. Jawahar<sup>†7</sup>, Richard Newcombe<sup>†5</sup>, Hyun Soo Park<sup>†10</sup>, James M. Rehg<sup>†8</sup>, Yoichi Sato<sup>†17</sup>, Manolis Savva<sup>†13</sup>, Jianbo Shi<sup>†16</sup>, Mike Zheng Shou<sup>†11</sup>, Michael Wray<sup>†6</sup>

<sup>1</sup>FAIR, Meta.

<sup>2</sup>University of Texas at Austin.

<sup>3</sup>Carnegie Mellon University.

<sup>4</sup>University of California, Berkeley.

<sup>5</sup>Project Aria, Meta.

<sup>6</sup>University of Bristol.

<sup>7</sup>IIIT, Hyderabad.

<sup>8</sup>University of Illinois, Urbana Champaign.

<sup>9</sup>Carnegie Mellon University.

<sup>10</sup>University of Minnesota.  
<sup>11</sup>National University of Singapore.

<sup>12</sup>Universidad de los Andes.

<sup>13</sup>Simon Fraser University.

<sup>14</sup>University of North Carolina, Chapel Hill.

<sup>15</sup>University of Catania.

<sup>16</sup>University of Pennsylvania.

<sup>17</sup>University of Tokyo.

<sup>18</sup>Meta.

<sup>19</sup>Indiana University.

<sup>20</sup>Johns Hopkins University.

<sup>21</sup>King Abdullah University of Science and Technology.

<sup>22</sup>Georgia Tech.

## Abstract

We present Ego-Exo4D, a diverse, large-scale multimodal multiview video dataset and benchmark challenge. Ego-Exo4D centers around simultaneously-captured egocentric and exocentric video of skilled human activities (e.g., sports, music, dance, bike repair). 740 participants from 13 cities worldwide performed these activities in 123 different natural scene contexts, yielding long-form captures from 1 to 42 minutes each and 1,286 hours of video combined. The multimodal nature of the dataset is unprecedented: the video is accompanied by multichannel audio, eye gaze, 3D point clouds, camera poses, IMU, and multiple paired language descriptions—including a novel “expert commentary” done by coaches and teachers and tailored to the skilled-activity domain. To push the frontier of first-person video understanding of skilled human activity, we also present a suite of benchmark tasks and their annotations, including fine-grained activity understanding, proficiency estimation, cross-view translation, and 3D hand/body pose. All resources are open sourced to fuel new research in the community.

<https://ego-exo4d-data.org/>

**Keywords:** video understanding, first-person video, egocentric, video-language, 3D, body pose

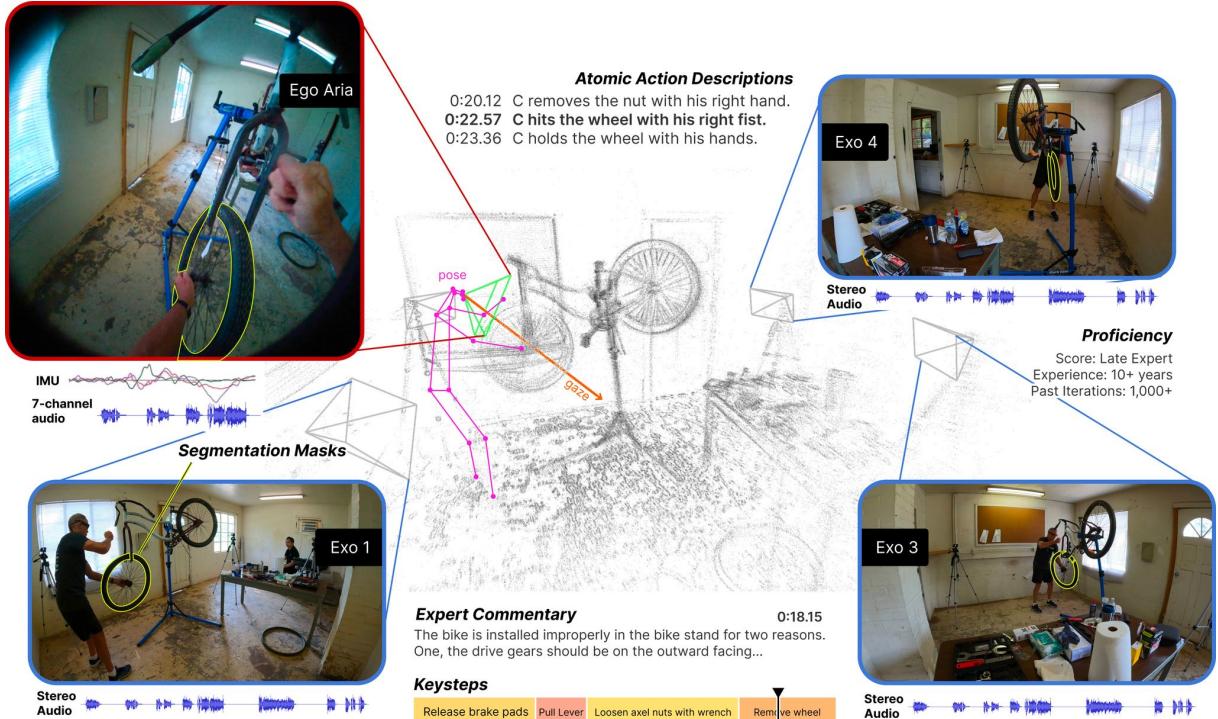
## 1 Introduction

A dancer leaps across a stage; Lionel Messi delivers a precise pass; your grandmother prepares her famous dumplings. We observe and seek human skills in a myriad of settings, from the practical (fixing a bike) to the aspirational (dancing beautifully). What would it mean for AI to understand human skills? And what would it take to get there?

Advances in AI understanding of human skill could facilitate many applications. In augmented reality (AR), a person wearing smart glasses could quickly pick up new skills with a virtual AI coach that provides real-time guidance. In robot learning, a robot watching people in its environment

could acquire new dexterous manipulation skills with less physical experience. In social networks, new communities could form based on how people share their expertise and complementary skills in video.

We contend that both the *egocentric* and *exocentric* viewpoints are critical for capturing human skill. Firstly, the two viewpoints are synergistic. The first-person (ego) perspective captures the details of close-by hand-object interactions and the camera wearer’s attention, whereas the third-person (exo) perspective captures the full body pose and surrounding environment context. See Figure 1. Not coincidentally, instructional or



**Fig. 1:** Ego-Exo4D offers egocentric video alongside multiple time-synchronized exocentric video streams for an array of skilled human activities—1,286 hours of ego and exo video in total. The data is both multiview and multimodal, and it is extensively annotated with language, 3D body and hand pose, keysteps, procedural dependencies, and proficiency ratings in support of our proposed benchmark tasks.

“how-to” videos often alternate between a third-person view of the demonstrator and a close-up view of their near-field demonstration. For example, a chef may describe their approach and the equipment from an exo view, then cut to clips showing their hands manipulating the ingredients and tools from an ego-like view.

Secondly, not only are the ego and exo viewpoints synergistic, but there is a need to *translate* fluently from one to the other when acquiring skill. For example, imagine watching an expert repair a bike tire, juggle a soccer ball, or fold an origami swan—then mapping their steps to your own body. Cognitive science tells us that even from a very young age we can observe others’ behavior (exo) and map it onto our own (ego) (Flavell et al., 1981, Newcombe, 1989), and this actor-observer translation remains the foundation of visual learning.

Realizing this potential, however, is not possible using today’s datasets and learning paradigms. Existing datasets comprised of both ego and exo

views (i.e., ego-exo) are few (Sigurdsson et al., 2018, Sener et al., 2022, Kwon et al., 2021, la Torre et al., 2009, Rai et al., 2021), small in scale, lack synchronization across cameras, and/or are too staged or curated to be resilient to the diversity of the real world. Thus the current literature for activity understanding primarily attends to *either* the ego (Damen et al., 2021, Grauman et al., 2022) or exo (Kay et al., 2017, Gu et al., 2018, Monfort et al., 2019, Soomro et al., 2012) view, leaving the ability to move fluidly between the first- and third-person perspectives out of reach. Instructional video datasets (Miech et al., 2019, Tang et al., 2020, Zhukov et al., 2019, Zhou et al., 2018) offer a compelling window into skilled human activity, but (like the above) are limited to single-viewpoint video, whether purely exocentric or mixed with “ego-like” views at certain time points.

We introduce Ego-Exo4D, a foundational dataset to support research on ego-exo video learning and multimodal perception. The result of a two-year effort by a consortium of 15

research institutions, Ego-Exo4D is a first-of-its-kind large-scale multimodal multiview dataset and benchmark suite. It constitutes the largest public dataset of time-synchronized first- and third-person video, captured by 740 diverse camera wearers in 123 distinct scenes and 13 cities worldwide. For every sequence, Ego-Exo4D provides both the camera wearer’s egocentric video, as well as *multiple* (4-5) exocentric videos from tripods placed around the camera wearer. All views are time-synchronized and precisely localized in a metric, gravity-aligned frame of reference. The total collection has 1,286 hours of video and 5,035 instances, each spanning 1 to 42 minutes of continuous capture.

Ego-Exo4D focuses on skilled single-person activities. The 740 participants perform skilled physical and/or procedural activities—dance, soccer, basketball, bouldering, music, cooking, bike repair, health care—in an unscripted manner and in natural settings (e.g., gym, soccer field, kitchens, bike shops, etc.), exhibiting a variety of skill levels from novice to expert. All video is recorded with rigorous privacy and ethics policies and formal consent of participants.

Ego-Exo4D is not only multiview, it is also multimodal. Captured with the unique open-source Aria glasses (Engel et al., 2023), all ego video is accompanied by 7-channel audio, IMU, eye gaze, both RGB and two grayscale SLAM cameras, and 3D environment point clouds. Additionally, Ego-Exo4D provides multiple new video-language resources, all time indexed: first-person narrations by the camera wearers describing their own actions; third-person play-by-play descriptions of every camera wearer action; and third-person spoken expert commentary critiquing their performance. The latter is particularly novel: performed by domain-specific experienced coaches and teachers, it focuses on *how* an activity is executed rather than merely *what* is being done, surfacing subtleties in skilled execution not perceivable by the untrained eye. All three language corpora are time-stamped against the video. To our knowledge, there is no prior video resource with such extensive and high quality multimodal data.

Alongside this data, we introduce benchmarks for foundational tasks for ego-exo video and we formalize them with annotations and evaluation

protocols to spur the community’s efforts. We propose four families of tasks:

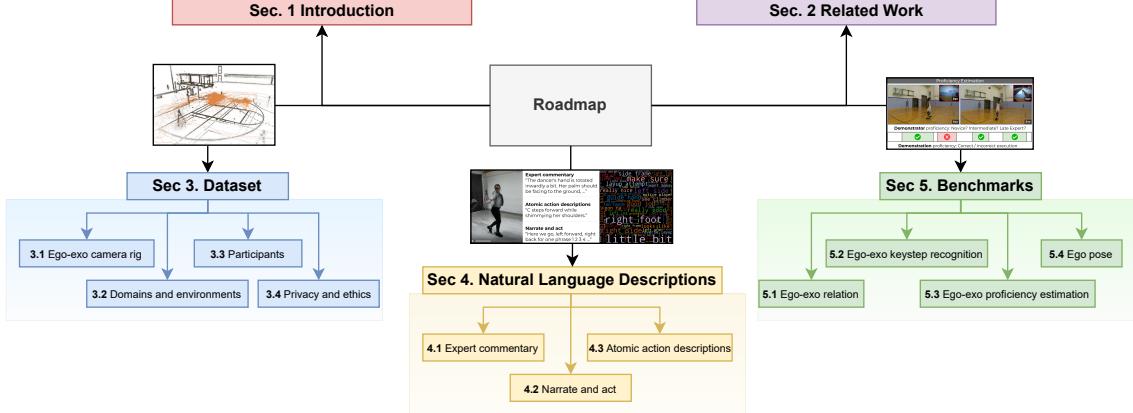
1. *ego-exo relation*, for relating the actions of a teacher (exo) to a learner (ego) by estimating semantic correspondences and translating viewpoints;
2. *ego(-exo) recognition*, for recognizing fine-grained keysteps and task structure;
3. *ego(-exo) proficiency estimation*, for inferring how well a person is executing a skill;
4. *ego pose*, for recovering skilled 3D body and hand movements from ego-video.

We provide annotations for each task—the result of more than 200,000 hours of annotator effort. To kickstart work in these new challenges, we also develop baseline models and report their results. We are hosting the first public benchmark challenges in 2024.

Though we are motivated by skill learning, Ego-Exo4D is poised for even broader influence, beyond the proposed benchmarks. Whereas existing datasets lack activity modeling in real-world 3D contexts (e.g., restricted to mocap suits and/or lab settings) and existing 3D datasets typically focus on static scenes and objects. Ego-Exo4D is a resource for **general 3D vision**—such as environment reconstruction, camera relocalization, audio-visual mapping, and many others. Similarly, our novel **video-language** resources will offer many opportunities for grounding of actions and objects, multimodal representation learning, and language generation. Finally, though our tasks prioritize perception from the “ego-only” perspective, the exo component of our data ensures its utility for the more **traditional exo viewpoint** too, e.g., for activity recognition and body pose estimation.

In summary, Ego-Exo4D is the community’s first diverse, large-scale multimodal multiview video resource. We have open sourced all the data, annotations, camera rig protocol, and benchmarks. With this release, we aim to fuel new research in ego-exo, multimodal activity, and beyond.

Figure 2 provides a roadmap for this paper. After reviewing related work (Sec. 2), we describe the dataset—its contents, camera setup, participants, and our approach to collection (Sec. 3)—followed by an overview of its three forms of natural language annotations (Sec. 4). Finally, we introduce the benchmark tasks organized into



**Fig. 2:** Overview of the paper and its sections, including the Ego-Exo4D dataset (Sec. 3.1, 3.2, 3.3), the natural language descriptions collected alongside the dataset (Sec. 4), and the benchmark tasks (Sec. 5).

the four families described above, outlining the motivation, task definitions, metrics, and baseline results for each (Sec. 5).

## 2 Related work

Next we review prior work in datasets, human skill, and cross-view analysis. Section 5 will discuss additional related work for each benchmark task.

### Egocentric datasets

There has been a surge of interest in egocentric video understanding, facilitated by recent ego-video datasets showing unscripted daily-life activity as in Ego4D (Grauman et al., 2022), EPIC-Kitchens (Damen et al., 2021, 2018, Tschernezki et al., 2023), UT Ego (Lee et al., 2012), ADL (Pirsiavash and Ramanan, 2012), and KrishnaCam (Singh et al., 2016), or procedural activities as in EGTea (Li et al., 2018), AssistQ (Wong et al., 2022), Meccano (Ragusa et al., 2021), CMU-MMAC (la Torre et al., 2009), EgoProcel (Bansal et al., 2022), and HoloAssist (Wang et al., 2023). Unlike any of the above, Ego-Exo4D focuses on multimodal ego *and* exo capture, and it is focused on the domain of skilled activities.

Many members of our Ego-Exo4D team worked together to create Ego4D (Grauman et al., 2022). The two datasets share some properties: both emphasize unscripted data, with long continuous captures in authentic environments with

diverse participants, and both offer novel benchmark tasks and video-language annotations. However, whereas Ego4D focuses on daily-life activity from the egocentric view alone, Ego-Exo4D focuses on skilled activity in specific domains, captures both egocentric and exocentric viewpoints, and is significantly more multimodal. Ego-Exo4D’s language annotations are also broader in scope compared to Ego4D, going beyond play-by-play action narrations to also include first-person how-to descriptions and third-person expert commentary about the skilled activities.

### Multiview and ego-exo datasets

Most existing multiview datasets focus on static scenes (Chang et al., 2017, Xia et al., 2018, Straub et al., 2019, Ramakrishnan et al., 2021, Xiao et al., 2013) and objects (Reizenstein et al., 2021, Wu et al., 2015), with limited (exo only) multiview human activity (Weinland et al., 2006, Corona et al., 2021). CMU-MMAC (la Torre et al., 2009) and CharadesEgo (Sigurdsson et al., 2018) are early efforts to capture both ego and exo video. CMU-MMAC (la Torre et al., 2009) features 43 participants in mocap suits who cook 5 recipes in a lab kitchen. In CharadesEgo (Sigurdsson et al., 2018), 71 Mechanical Turkers record 34 hours of scripted scenarios (e.g., “type on laptop, then pick up a pillow”) from the ego and exo perspectives sequentially, yielding unsynchronized videos with non-exact activity matches. More recent ego-exo efforts focus on specific activities in one or two environments. Assembly101 (Sener

Dataset	Year	Modalities	#Subj.	#Scenes	#Tasks	#Actions	#Masks	#BP	#HP	Nar.	EC
<i>Multimodal Egocentric Datasets</i>											
EGTEA-Gaze (Li et al., 2018)	2018	V,A,G	32	1	7	106	15k	-	-	x	x
MECCANO (Ragusa et al., 2021)	2021	V,D,G	20	2	1	61	-	-	-	x	x
EK100 (Damen et al., 2022)	2022	V,A	37	45	N/A	(97:300)*	-	-	-	✓	x
Ego4D (Grauman et al., 2022)	2022	V,A,3D,S,G,I	931	74	N/A	110†	-	-	-	✓	x
HoloAssist (Wang et al., 2023)	2024	V,A,D,G,3D,I	222	?	20	(49:165)*	-	?	?	✓	x
<i>Multiview Datasets</i>											
IXMAS (Weinland et al., 2006)	2006	V	10	1	N/A	11	-	-	-	x	x
MEVA (Corona et al., 2021)	2021	V,T,GPS	100	28	N/A	37	-	-	-	x	x
<i>Ego-Exo Datasets</i>											
CMU-MMAC (la Torre et al., 2009)	2009	V,A,M,I	43	1	5	-	-	-	-	x	x
Charades-Ego (Sigurdsson et al., 2018)	2018	TODO	71	N/A	N/A	157	-	-	-	x	x
LEMMa (Jia et al., 2020)	2020	V,D	8	14	15	(24:64)*	-	-	-	x	x
HOMAGE (Rai et al., 2021)	2020	V,A,T,B,Ma	27	10	70	453	-	-	-	x	x
H2O (Kwon et al., 2021)	2021	V,D	4	3	N/A	36	-	-	0.5M	x	x
Assembly101 (Sener et al., 2022, Ohkawa et al., 2023)	2022	TODO	53	1	101	(24:90)*	-	-	0.2M	x	x
EgoExoLearn (Huang et al., 2024)	2024	V,A,G,I	136	7	8	(95:254)*	-	-	-	✓	x
EgoExo4D	2024	V,A,I,G,3D,6D,B,Ma	740	123	43‡	689	2.2M	9.6M	4.4M	✓	✓

**Table 1:** Comparison between Ego-Exo4D and relevant datasets. Compared to existing datasets capturing both egocentric and exocentric views, Ego-Exo4D features more modalities, more subjects, and significantly larger scene diversity, as well as rich annotations including key-step segments, object masks, and three meticulously synchronized natural language descriptions paired with the videos (narrations, narrate-and-act, and expert commentary). To our knowledge, Ego-Exo4D also offers the largest available manual ground truth egocentric body pose annotations to date (in the above datasets or any others), and it has  $\sim 14$ M total frames of 3D pose annotations and pseudo-annotations. #*Tasks* denotes the number of tasks that subjects were asked to execute in each dataset, *Subj.* denotes recorded subjects, #*BP* refers to number of 3D body poses, #*HP* refers to number of 3D hand poses, *Nar.* denotes narrations, and *EC* refers to expert commentary annotations. Modality abbreviations: **V**ideo, **A**udio, **D**epth, **G**aze, **S**tereo, **IMU**, **3D** Environments, **T**hermal **I**R, **G**PS, **M**otion Capture, **6DOF**, **B**arometer, **M**agnetometer. \* denotes action taxonomies defined in terms of verbs and nouns, statistics reported as (number of verbs; number of nouns). † The number has been taken from the Moment Query benchmark. ‡ Number of tasks for Ego-Exo4D includes 21 procedural activities and 22 physical activities (listed in Table 2).

et al., 2022) and H2O (Kwon et al., 2021) provide time-synced ego and exo video at a lab tabletop where people assemble toy cars or manipulate handheld objects, with 53 and 4 participants, and 513 and 5 hours of footage, respectively. LEMMA (Jia et al., 2020) contains multi-agent, multi-task activities with 15 common daily tasks, performed by 8 individuals in 14 unique kitchens/living rooms. Homage (Rai et al., 2021) provides 30 hours of ego-exo video from 27 participants in 2 homes doing household activities like laundry. EgoExoLearn (Huang et al., 2024) provides 120 hours of egocentric videos emulating the human demonstration following process with exocentric demonstration videos.

Compared to any of the prior efforts, Ego-Exo4D offers an order of magnitude more participants, diverse locations, and hours of footage (740 participants, 123 unique scenes, 13 cities, 1,286 hours). Importantly, our focus on skilled tasks takes the participants out of the lab or home

and into settings like soccer fields, dance studios, rock climbing walls, and bike repair shops. Such activities also yield a wide variety of full body poses and movements within the scene, beyond using objects at a tabletop. This variety means Ego-Exo4D augments existing 3D human body pose datasets (Zhang et al., 2022, Li et al., 2023, Joo et al., 2017, Khirodkar et al., 2023, Guzov et al., 2021). Finally, compared to any prior ego-exo resource, Ego-Exo4D’s suite of modalities and benchmark tasks are novel and will expand the research directions the community can take for egocentric and/or exocentric video understanding. Table 1 summarizes Ego-Exo4D’s properties compared to those of existing datasets.

#### Human skill and video learning

Analyzing skill and action quality has received limited attention (Pirsavash et al., 2014, Bertasius et al., 2017, Parmar and Morris, 2019, Doughty et al., 2018, 2019, Zhang et al., 2023). Research in instructional or “how-to” videos

is facilitated by (largely exo) datasets like HowTo100M (Miech et al., 2019) and others (Tang et al., 2020, Zhukov et al., 2019, Zhou et al., 2018, Ben-Shabat et al., 2020). Challenges include grounding keysteps (Miech et al., 2019, Zhukov et al., 2019, Bansal et al., 2022, Elhamifar and Huynh, 2020, Xu et al., 2021, Miech et al., 2020, Dvornik et al., 2022, Lin et al., 2022), procedural planning (Chang et al., 2020, Bi et al., 2021, Zhao et al., 2022, Wang et al., 2023, Zhong et al., 2023, Shvetsova et al., 2022, Ko et al., 2022, Cao et al., 2022), learning task structure (Elhamifar and Huynh, 2020, Narasimhan et al., 2023, Zhou et al., 2018, Alayrac et al., 2016, Ashutosh et al., 2023, Zhou et al., 2023), and leveraging noisy narrations (Miech et al., 2019, 2020, Lin et al., 2022). A portion of Ego-Exo4D is procedural activities, but unlike the above, it offers simultaneous ego-exo capture. The scale and diversity of our data—including its three forms of language descriptions—widen the avenues for skilled activity understanding research.

#### **Ego-exo cross-view modeling**

There is limited prior work on ego-exo cross-view modeling, arguably due to a lack of high-quality synchronized real-world data. Prior work explores matching people between videos (Ardesir and Borji, 2016, 2018, Fan et al., 2017, Xu et al., 2018, Wen et al., 2021) and learning view-invariant (Sigurdsson et al., 2018, Ardesir and Borji, 2018, Sermanet et al., 2018, Yu et al., 2019, 2020, Xue and Grauman, 2023) or ego features (Li et al., 2021). Beyond the specific case of ego-exo, cross-view methods are explored for translation (Regmi and Borji, 2018, 2019, Tang et al., 2019, Ren et al., 2021, Luo et al., 2024, Cheng et al., 2024), novel view synthesis (Liu et al., 2021, Ren and Wang, 2022, Rombach et al., 2021, Wiles et al., 2020, Watson et al., 2022, Tseng et al., 2023, Chan et al., 2023), and aerial to ground matching (Regmi and Shah, 2019, Lin et al., 2015). Ego-Exo4D provides a testbed of unprecedented size and variety for cross-view modeling. In addition, our ego-exo relation tasks (cf. Section 5) surface new challenges in novel-view synthesis with widely varying viewpoints.



**Fig. 3:** The Aria device used for egocentric recordings.

### **3 Ego-Exo4D dataset**

Next we introduce the dataset and its scope. Notably, the video capture was a distributed but coordinated effort performed by 12 research labs who worked together over nearly two years to create Ego-Exo4D. Importantly, our data collection across the sites was a coordinated effort, with common guidelines, scenarios, and camera rigs. In this way, the dataset is cohesive at the same time it is diverse.

In the following, we first introduce the ego-exo camera rig and time synchronization process (Sec. 3.1). Then we overview the domains and activities that compose the dataset (Sec. 3.2), followed by discussion of the participants’ diverse backgrounds and expertise (Sec. 3.3).

#### **3.1 Ego-exo camera rig**

To collect ego-exo data at a global scale, we developed a low-cost camera recording rig that was portable, auto-synchronized, captured a rich suite of sensor data, and attainable internationally.

Our solution consists of 1 Aria (see Figure 3), 4 GoPros<sup>1</sup>, 1 GoPro Remote, 4 Tripods, 4 SD Cards, 4 Tripod Mount Adapters, 4 Velcro’d Battery Packs, 4 USB-A to USB-C Cables, 1 Glasses Sports Strap, 1 Smartphone, and 1 Laptop or Tablet for questionnaires. The total cost excluding the Aria/phone/laptop is under \$3,000.

##### **3.1.1 Aria device and sensors**

Aria is an egocentric recording device in glasses form-factor created by Meta. It is designed as a *research tool* for egocentric machine perception and contextualized AI research, and available

---

<sup>1</sup>This represents the common core of the collection rig used in all capture settings. In certain captures, *additional* exo or ego GoPros are also used.

to researchers across the world through [projec-taria.com](http://projec-taria.com).

The Aria device emulates future AR- or smart-glasses catering to machine perception and egocentric AI rather than human consumption. It is designed to be wearable for long periods of time without obstructing or impeding the wearer, allowing for natural motion even when performing highly dynamic activities—such as playing soccer or dancing. It has a total weight of 75g (compared to over 150g for a single GoPro camera), and fits just like a pair of glasses.

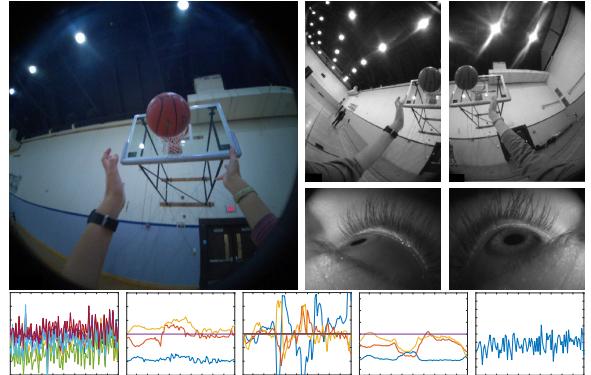
Further, the device integrates a rich sensor suite that is tightly calibrated and time-synchronized, capturing a broad range of modalities. See Figure 4. For Ego-Exo4D, the following sensor configuration is used:

- **One rolling-shutter RGB camera** recording at 30fps and  $1408 \times 1408$  resolution covering a field of view of  $110^\circ$ .
- **Two global-shutter monochrome cameras** recording at 30fps and  $640 \times 480$  resolution. They provide peripheral vision, and each cover a field of view of  $150^\circ$ .
- **Two monochrome eye-tracking cameras** recording at 10fps and  $320 \times 240$  resolution.
- **An array of seven microphones** recording spatial audio around the wearer.
- **Two IMUs** (800Hz and 1000Hz respectively), **a barometer** (50fps) and **a magnetometer** (10fps).

All sensor streams come with metadata such as timestamps and per-frame exposure times. All data is made available in raw form as part of the Ego-Exo4D dataset. For convenience, we also include pre-computed slices of data that suit specific purposes, e.g., 2D gaze points, mp4s of each camera, and smaller .vrs files with a subset of sensor streams.

### 3.1.2 Precomputed 3D spatial signals

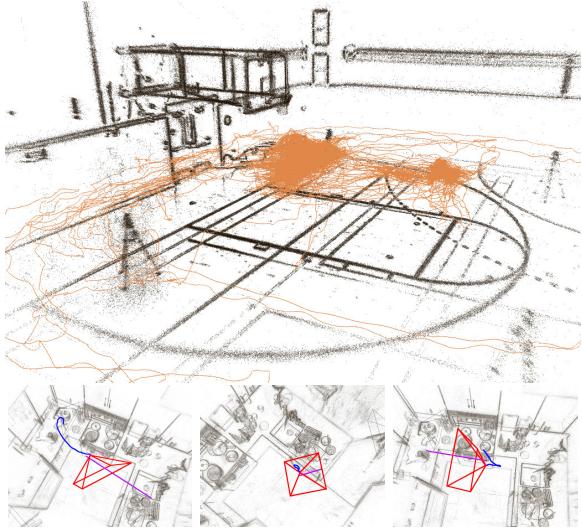
Project Aria’s machine perception service (MPS) provides software building blocks that simplify leveraging the different modalities recorded. These functionalities are likely to be available as real-time, on-device capabilities in future AR- or smart-glasses. We use the following core functionalities and include their raw output as part of the



**Fig. 4:** Sensor streams recorded by the Project Aria device. Top: RGB camera, left and right monochrome and eye cameras. Bottom: 10-second extracts from microphones, accelerometer, gyroscope, magnetometer and barometer respectively.

dataset. See Figure 5. See ([Engel et al., 2023](#)) for more details.

- **Calibration:** All sensors are intrinsically and extrinsically calibrated. MPS also provides time-varying online-calibration that corrects for tiny deformations due to temperature changes or stress applied to the glasses frame.
- **Aria 6 DoF localization:** Every recording is localized precisely and robustly in a common, metric, gravity-aligned coordinate frame, using a state-of-the-art VIO and SLAM algorithm. This provides millimeter-accurate 6 DoF poses for every captured frame, as well as high-frequent (1kHz) motion in-between camera frames.
- **Eye gaze:** The gaze direction of the user is estimated as a single outward-facing ray anchored in-between the wearer’s eyes. We use an optional eye gaze calibration procedure, where the mobile companion app directs the wearer to gaze at a pattern on the phone screen while performing specific head movements. This information was then used to generate a more accurate eye gaze direction, personalized to the particular wearer.
- **Point clouds:** A 3D point cloud of static scene elements is triangulated from the moving Aria device, using photometric stereo over consecutive frames or left/right SLAM camera. The output contains both the 3D point clouds as well



**Fig. 5:** Aria MPS output for several recordings. Top: point cloud and estimated egocentric camera trajectory for a basketball session in Chapel Hill. This single continuous recording is 60 minutes long, has a total trajectory length of 2188 m, and contains 41 distinct takes. Bottom: three screenshots of a cooking recording, visualizing the current camera pose (red), eye gaze (purple), and last second of motion (blue).

as the raw, causally computed, 2D observations of every point in the camera images.

- **GoPro 6 DoF localization:** For Ego-Exo4D, we additionally built functionality on top of the existing Aria MPS functionality, specifically to localize the static GoPro cameras. To achieve this, we use the map built with Aria’s SLAM cameras, and perform 6 DoF localization of GoPro frames on the map. To obtain the GoPro calibration, we manually calibrated one device in the lab to obtain default parameters, and then use the P4P (Kukelova et al., 2016) algorithm (with RANSAC to reject matching outliers) to estimate the 6 DoF pose, as well as re-estimate the focal length to compensate for possible calibration variation between devices.

### 3.1.3 Recording procedure

Our recording procedure involves setting up the static GoPros on tripods in locations generally consistent within each scenario, conducting a walk-around with the Aria to build a basemap for

3D reconstruction and camera localization, displaying QR codes at the start/end to assist time sync, and showing a take separation QR between each take. The camera rig was extended for certain sites with additional mounts and GoPros, discussed below.

Specifically, to sync cameras, we employ a pre-rendered sequence of QR Codes (*i.e.*, QR code video) that encode a wall-clock time. We show this QR code video using the smartphone at 29fps to all cameras in sequence and exploit the difference in frame rates to finely sync the cameras. An additional stage of manual verification ensures each GoPro camera was within 1 frame (+-16.66ms) of the Aria RGB camera. To amortize the setup and tear down time required for each recording, we record multiple ‘takes’ (*i.e.*, one instance of a certain task) back-to-back and use a ‘Take Separator’ QR code to separate takes in post-processing.

## 3.2 Domains and environments

Ego-Exo4D focuses on *skilled human activity*. This is in contrast to existing ego-only efforts like Ego4D (Grauman et al., 2022), which has a broad span of daily-life activities. We intentionally select the domains based on a few criteria: Will it illustrate skill and a variety of expertise? Do we have access to real-world settings and participants for that scenario? Is there visual variety to be expected across different instances? Will the ego and exo views offer complementary information? Will it present new challenges unaddressed by current datasets? Overall, by scoping to certain domains, we aim to build up sufficient density of data within a core set of skills for training and evaluating models.

### Physical and procedural activities

Intersecting these criteria, we arrived at two broad categories<sup>2</sup> of skilled activity: *physical* and *procedural*, together comprising eight total domains. The physical domains are soccer, basketball, dance, bouldering, and music. They emphasize body pose and movements as well as interaction with objects (e.g., a ball, musical instrument). The procedural domains are cooking, bike repair, and

---

<sup>2</sup>Note that in general physical and procedural are not mutually exclusive labels. An activity can both require physical skill and procedural steps.



**Fig. 6:** Ego-Exo4D captures skilled activity from 43 tasks and 689 keysteps in 8 domains, in a wide variety of 123 scenes in 13 cities in Japan, Colombia, Canada, India, Singapore, and 7 US states. Each domain is captured at multiple sites—from 2 to 64 unique locations. In total the dataset offers 1,286 hours of ego+exo video comprised of 5,035 takes from 740 camera wearers. An average take is 2.6 minutes.

health care. They require performing a sequence of steps to reach a goal state (e.g., a completed recipe, a repaired bike) and generally entail intricate hand-object manipulations with a variety of objects (e.g., bike repair tools; cooking utensils, appliances, and ingredients). All domains entail regular attention shifts (revealed by head pose and gaze) by the participant. Figure 6 summarizes the eight domains with example frames and data statistics for each. In Section 5 we discuss how the domains relate to the proposed benchmarks.

In total, we have 43 activities derived from the eight domains. For example, cooking is comprised of 14 recipes; soccer is comprised of 3 drills, and music is comprised of 3 instruments. See Table 2. Those 43 activities break down further into 689 total unique keysteps. The length of a take ranges from 8 seconds to 42 minutes, with procedural activities like cooking having the longest sustained captures.

#### *Distribution of activities per site*

To achieve visual diversity in the data, multiple labs across our team (typically 3-5) captured each Ego-Exo4D domain. Figure 9 shows the breakdown of which scenarios were captured by each partner institution as well as a map highlighting the locations of the 12 labs involved in data collection.<sup>3</sup> The domain selection per site is based on the lab’s own preferences and local opportunities

to capture data of these scenes at scale. Cooking is our one cross-cutting domain, collected at each site. We identified cooking as a priority domain because it resonates around the world as a human need and interest. In total, the cooking scenario of Ego-Exo4D contains more than 650 takes of cooking performed by more than 170 chefs in 60 different environments around the world, forming nearly 100 hours of ego video alone.

#### *Authentic environments for capture*

The data is collected in authentic settings—such as real-world bike shops, soccer pitches, or bouldering gyms—as opposed to lab environments. Since every domain is covered by more than one lab, the dataset exhibits visual variety from the different physical locations. For example, we have videos of chefs in New York City, Vancouver, Philadelphia, Bogota, and others; soccer players in Tokyo, Chapel Hill, Hyderabad, Singapore, and Pittsburgh. Furthermore, even within the captures done by a single lab, there are often multiple different sites used for filming (e.g., a couple different bike shops in the same city). Figure 7 and 8 shows example frames illustrating the variety of the sites and tasks.

#### *Domain-specific collection guidelines*

To ensure consistency across the dataset, we developed data collection guidelines for each domain.

<sup>3</sup>An additional four institutions not shown on the map are part of the Ego-Exo4D consortium (e.g., contributing to benchmarks) but did not collect data. They are UT Austin (USA),

KAUST (Saudi Arabia), University of Catania (Italy), and University of Bristol (UK).

Procedural	Physical
<p><i>Cooking:</i></p> <ul style="list-style-type: none"> <li>- Omelette</li> <li>- Scrambled eggs</li> <li>- Tomato and egg</li> <li>- Sesame-ginger Asian salad</li> <li>- Greek salad</li> <li>- Dumplings</li> <li>- Noodles</li> <li>- Pasta</li> <li>- Sushi roll</li> <li>- Samosa</li> <li>- Coffee latte</li> <li>- Chai tea</li> <li>- Milk</li> <li>- Cookies</li> <li>- Brownies</li> </ul> <p><i>Health:</i></p> <ul style="list-style-type: none"> <li>- COVID test</li> <li>- Cardiopulmonary Resuscitation (CPR)</li> </ul> <p><i>Bike repair:</i></p> <ul style="list-style-type: none"> <li>- Remove/install a wheel</li> <li>- Replace an inner tube</li> <li>- Clean and lubricate the chain</li> <li>- Adjust rear derailleur</li> <li>(both limit screws &amp; indexing)</li> </ul>	<p><i>Music:</i></p> <ul style="list-style-type: none"> <li>- Violin</li> <li>- Piano</li> <li>- Guitar</li> </ul> <p><i>Basketball:</i></p> <ul style="list-style-type: none"> <li>- Mikan layup drill</li> <li>- Righthand reverse layup</li> <li>- Mid-range jump shot</li> </ul> <p><i>Soccer:</i></p> <ul style="list-style-type: none"> <li>- Freestyle dribbling</li> <li>- Freestyle juggling</li> <li>- Penalty kicks</li> </ul> <p><i>Dance:</i></p> <ul style="list-style-type: none"> <li>- Easy choreography</li> <li>- Advanced choreography</li> </ul>

**Table 2:** The 43 specific activities collected for the three *procedural* and five *physical* domains

These guidelines describe the recommended camera positioning, instructions for participating camera wearers, along with important context-specific considerations. For example, given privacy concerns, our health guidelines required data collection participants to discard COVID tests before results were visible. The guidelines provide general parameters from which to collect data; however, they were not rigid steps. Indeed, to support diversity and implementation at a global scale, there are site-specific nuances, for example, differing standards for the implementation of CPR, cultural differences in the ingredients used for different targeted dishes, and location-specific bouldering routes.

The primary domain-specific design decision is the placement of the exocentric cameras. The best visibility points for a given domain depends on the general scene and objects involved (open soccer field vs. cluttered kitchen with cabinets) as well as how the person interacts with the space. For example, for soccer recordings, an exo camera placed in the goal gives great visibility of the players' shots, while in dance or music, an overhead exocentric (or downward pointing egocentric) camera captures important close-to-body detail about the participants' arms and hands. Generally the exo cameras were placed to ensure viewpoint coverage and achieve the complementary hand-object near-field interactions as well as the participants' full-body movements.

Appendix B describes the data collection details that are specific to each consortium partner site, e.g., how they recruited participants, which of the domains and activities they captured, or any modalities they added on top of the common rig.

### 3.3 Participants

Next we describe the participants who wore the egocentric cameras in Ego-Exo4D.

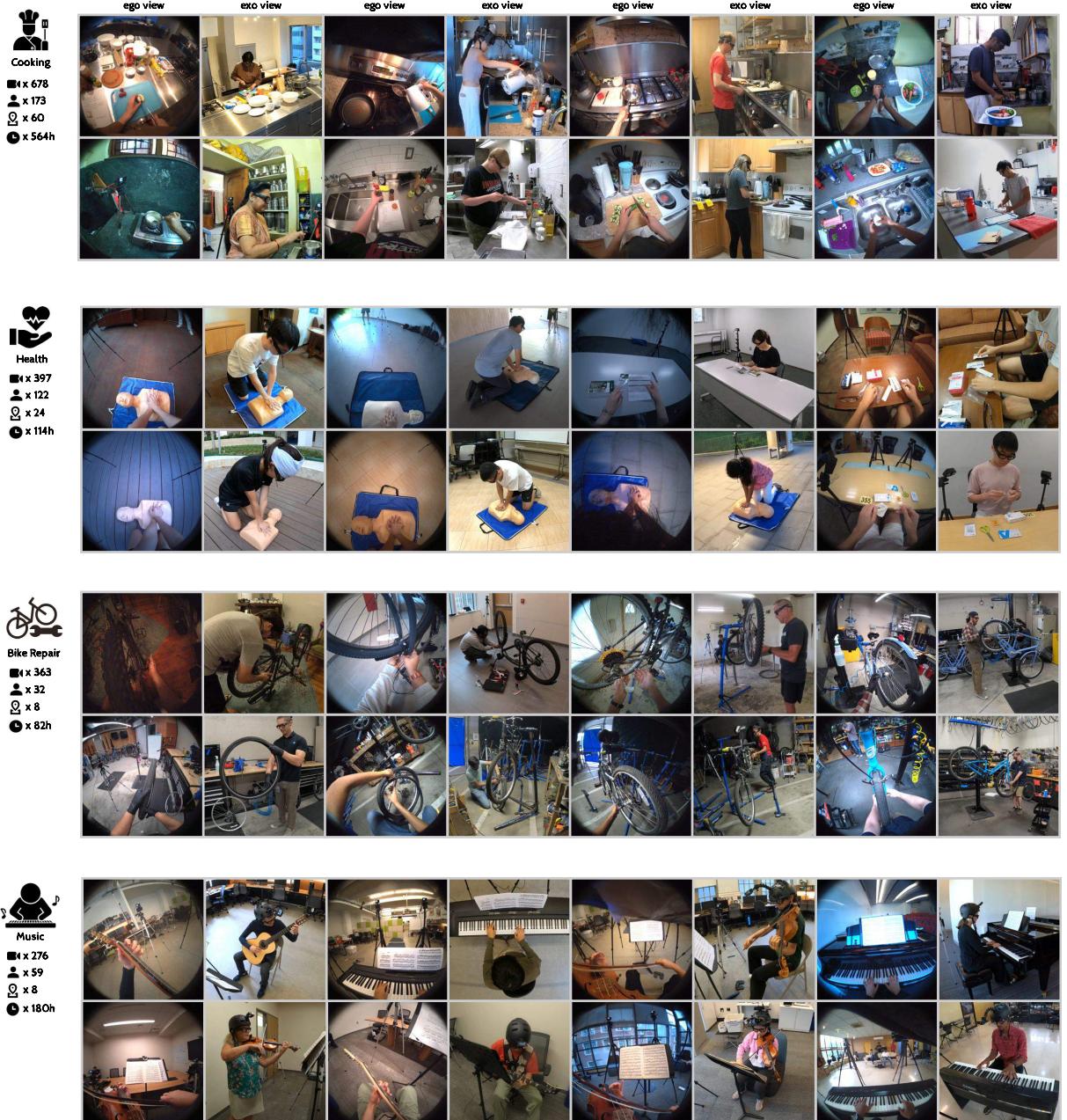
#### Credentials and expertise

We recruited 740 total participants from the local communities of 12 labs. All scenarios feature real-world experts, where the camera-wearer participant has specific credentials, training, or expertise in the skill being demonstrated. For example, among the Ego-Exo4D camera wearers are professional and college athletes; jazz, salsa, and Chinese folk dancers and instructors; competitive boulderers; professional chefs who work in industrial-scale kitchens; bike technicians who service dozens of bikes per day. Many of them have (individually) over 10 years of experience.

Experts are prioritized given they are likely to conduct activities without mistakes or distractions, providing a strong ground truth for how to approach a given task. However, we also include capture from people with varying skill levels, as well—essential for our proposed skill proficiency estimation task (Section 5.3). Notably, Ego-Exo4D represents human intelligence in a new way by capturing domain-specific expertise—both in the video as well as the accompanying expert commentary (see Section 4)—portraying the evolution of a skill from beginners to experts.

#### Demographics

The camera wearers range in age from 18 to 74 years old, with 37% self-identifying as female 60% male and 3% as non-binary or preferring not to say. See Figure 10. In total, the participants self



**Fig. 7:** Ego-Exo4D captures skilled activity from 8 domains, in a wide variety of 123 scenes in 13 different cities in Japan, Colombia, Canada, India, Singapore, and 7 US states. Every odd column shows an ego view, and the adjacent even column shows one of its paired exo views.

report more than 24 different ethnicities.<sup>4</sup> Details are in Appendix C.

---

<sup>4</sup>Sharing this information was optional for all research subjects. Ethnicity is reported based on location specific categories as defined by the relevant partner lab. No such information was gathered from research subjects participating in our collections in California, New York, and Pittsburgh, Pennsylvania.

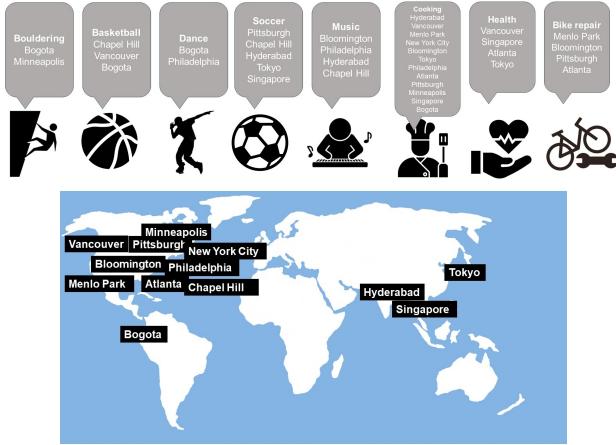


**Fig. 8:** Ego-Exo4D captures skilled activity from 8 domains, in a wide variety of 123 scenes in 13 different cities in Japan, Colombia, Canada, India, Singapore, and 7 US states. Every odd column shows an ego view, and the adjacent even column shows one of its paired exo views.

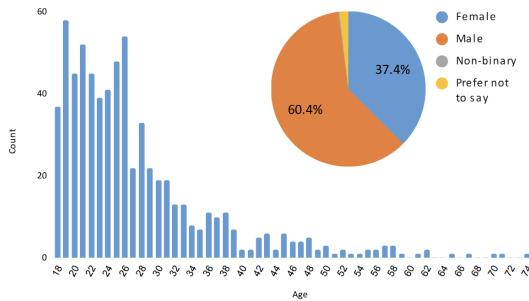
### Recruiting

To recruit participants, each partner institution chose its own approach. This included using campus email lists, flyers in coffee shops, word of mouth to family and friends, online ads, posts on

social media, university communication channels, temp hiring agencies, and connecting with schools, gyms, and teams, e.g., soccer schools, climbing gyms, professional and university athletics organizations.



**Fig. 9:** Geographic coverage of Ego-Exo4D and breakdown of which scenarios are captured in which cities. Note that even within a given city, there may be multiple sites (e.g., multiple bike repair shops or kitchens in the same city).



**Fig. 10:** Camera wearer participants' self-reported demographic information (age and gender)

### Characterizing skill levels

To ensure consistent, high quality annotations for our benchmarks (discussed below), we identified three crucial pieces of information about the participants that would be difficult to capture with third-party annotators: the participant’s *skill* level, the *objects* they are using, and the *actions* they are completing. We captured the participants’ perceived skill level and performance using pre-task and post-task surveys, respectively, which are available with the Ego-Exo4D dataset.

For the pre-task survey, we ask 10 questions like, “how many years have you been doing this

task?” and “have you taught this activity to others before?” (details in Table C1 in Appendix C). These questions are designed to be more easily quantifiable than simply asking participants to self-rate their skill level.<sup>5</sup>

For the post-task survey, we ask the participant to reflect on how well they did the task, with questions like, “what mistakes or errors did you make?” and “did it take longer or shorter than your initial expectation and why?”. Finally, to capture the *actions* and *objects* with which they interact, we ask participants to perform a round of first-person narrations called “narrate-and-act” (detailed below in Section 4).

### 3.4 Compliance with ethical standards

Ego-Exo4D was collected following rigorous privacy and ethics standards. This included undergoing formal independent review processes at each institution to establish the standards for collection, management, and informed consent. Similarly, all Ego-Exo4D data collection adhered to the [Project Aria Research Community Guidelines](#) for responsible research. Since the scenarios allow for closed environments (e.g., no passerbys) nearly all video is available without de-identification. For information about each individual partner’s protocols and restrictions, please see Appendix B. Ego-Exo4D data is gated behind a license system, which defines permitted uses, restrictions, and consequences for non-compliance.

## 4 Natural language descriptions

Ego-Exo4D also offers three kinds of paired natural language datasets, each time-indexed alongside the video: expert commentary, narrate and act, and atomic action descriptions. See Figure 11 and Figure 14 for examples of each language type highlighting their distinctions in style and point of view.

These language annotations are not steered towards any single benchmark, but rather are a general resource that will inspire new language-vision possibilities, such as grounding actions and

---

<sup>5</sup>We also obtain proficiency ratings for the participants via our expert commentators (cf. Section 4.1).



### Expert commentary

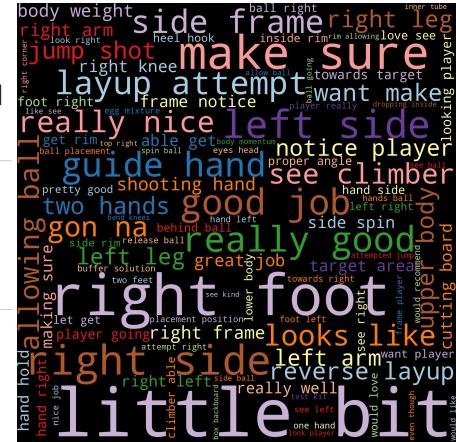
"The dancer's hand is rotated inwardly a bit. Her palm should be facing to the ground, ..."

### Atomic action descriptions

"C steps forward while shimmying her shoulders."

### Narrate and act

"Here we go, left forward, right back for one phrase 1 2 3 4 ..."



**Fig. 11:** Ego-Exo4D offers 3 paired video-language corpora. Word cloud is from expert commentary which critiques the performance.

objects, self-supervised representation learning, multimodal embeddings, video-conditioned language models, and skill assessment from video. We also anticipate the temporally grounded descriptions to be valuable for pre-training foundation models (Lin et al., 2022, Pramanick et al., 2023) or automated video captioning (Pan et al., 2020, Iashin and Rahtu, 2020, Zhao et al., 2023), both of which are rapidly growing areas of research. Furthermore, the time-anchored aspect of the three language annotations provides the opportunity to retrieve time points in specific Ego-Exo4D videos that correspond to queried moments, actions, or phrases. Finally, they are valuable to mine the dataset for the distribution of objects and activities present, e.g., for taxonomy formation.

## 4.1 Expert commentary

The first language dataset is spoken *expert commentary*. The goal is to reveal nuances of the skill that are not always visible to non-experts. We recruited 52 experts (distinct from the participants) to critique the recorded videos, call out strengths and weaknesses, explain how the specific behavior of the participant (e.g., hand/body pose, use of objects) affects the performance, and provide spatial markings to support their commentary. We provide both the transcribed speech and the raw audio (interesting for its inflection and non-word utterances), as well as the experts' spatial drawings and numeric ratings of each participant's skill.

These commentaries are quite novel: they focus on *how* the activity is executed rather than *what* it entails, capturing subtle differences in skilled execution. We believe this can unlock new fundamental problems (e.g., proficiency estimation below) and disruptive future applications (e.g., AI coaching).

In the following, we describe the qualifications and background of our experts, followed by the commentary instructions and scope.

### Experts' qualifications

The 52 experts are not only well-credentialed in their areas of expertise, but also have coaching or teaching experience. When recruiting the experts, our selection criteria focused on technical skills, communication, and performance during a live video commentating exercise.

On average, 90% of the recruited experts possess more than 10 years of professional experience and all have served during this time in the capacity of a coach, instructor or mentor. All experts further have either an advanced degree in their domain of focus or an industry certification. Certification authorities include the US Soccer Federation, the American Culinary Federation, USA Climbing, the American Red Cross, Trek Bikes, and New York State's Initial Certification in Teaching Dance, among others. Multiple individuals were recruited across each domain, with the goal of generating language and expertise diversity. Due to employment considerations,



**Fig. 12:** The 52 experts who perform the expert commentary language annotations are highly trained in the domain they are commentating, and they often have professional coaching or teaching experience.

all experts are residents of the United States. See Figure 12.

#### *Expert commentary guidelines*

Experts are provided with two time-synchronized videos of each Ego-Exo4D skills demonstration—one showing the egocentric view and another providing a single exocentric perspective specifically selected by annotators as the view that provides the best visibility on the scene (see Sec. 4.3). Experts are first asked to watch the video in full without commenting to gain an understanding of the skills demonstration and plan out important points to note in their commentary.

Then, the experts watch the video and pause every time they have a comment, typically 7 times per minute of video. The experts are encouraged to focus on critiques and teaching advice, as opposed to simply describing what the participant is doing. We record their spoken language descriptions of what is most effective or ineffective about the camera wearer’s actions, the quality of the execution, and mistakes they see. All commentary is time-anchored and retrospective, focusing on insights and perspectives relating to actions visible up until that point in the video. We choose to collect commentary as verbal recordings in order to maintain the naturalness of the performance descriptions and do so quickly. Each piece of spoken commentary is unbounded in length, and averages 4 sentences. We transcribe the commentaries automatically with OpenAI’s Whisper for automatic speech recognition (Radford et al.,



Proficiency score: 10

Commentary: Great footwork. He’s using dribble to set up his footwork and his shot. Stepping onto that left foot bringing the ball. I love that his eyes and head are up. He already knows where he’s going to go.



Proficiency score: 5

Commentary: The dancer’s hands should be a bit higher. This line should be completely straight in front of him at his shoulder length. It shouldn’t be beginning to dip lower.

**Fig. 13:** Two examples of expert commentary and proficiency scores, along with spatial drawings (red) done by the expert to augment their spoken comments.

2023). Figure 14 shows example commentaries; more are available in Table D2 in Appendix D.

Aside from the spoken commentary, the experts also provide spatial drawings and proficiency scores. During commentary, experts had the option to use a “telestrator” tool to enhance their commentary with freehand sketches to spatially localize information or otherwise help explain a point (see Figure 13). They also provide an overall proficiency rating on each video, assessing how well the task was performed with a short written justification. They score the video on a scale of 1 (least skilled) to 10 (most skilled). In many cases, experts coordinated within their domain group to calibrate this scoring.

Each video has expert commentary by 2-5 distinct experts, offering a variety of perspectives for the same content. In total, we have 117,812 pieces of time-stamped, video-aligned commentary, the result of more than 6,000 hours of work by the 52 experts. Overall, we believe the commentary is a unique window into the skilled actions that (through language) surfaces many subtleties about the actions not evident to the untrained eye.



**Fig. 14:** Examples of the three different language annotation styles: narrate and act, atomic action descriptions, and expert commentaries from four of the scenarios (bike repair, health, dance, and basketball). We also include word clouds which highlight the differences in vocabularies per scenario. In narrate and act text, we see how the participant briefly describes what they are doing and why, whereas the atomic action descriptions provide strictly a statement about the visible actions. The expert commentary offers an expert's critique of what is shown, commenting on strengths and weaknesses and explaining how the participant's actions affect their performance.

## 4.2 Narrate and act

The second language dataset consists of *narrate-and-act descriptions* provided by the participants themselves. They are in the style of a tutorial or how-to video, where the participant explains what they are doing and why. They are reminiscent of the narrations provided in Internet how-to videos, but with less stylization and without any professional post-production editing.<sup>6</sup> See Figure 14 for examples; more are in Table D2 in the Appendix.

Unlike the third-party expert commentary above, these are first-person reflections on the activity given by the people doing them. Generally the commentary is richer in constructive feedback about the quality of the activity, whereas the narrate-and-act narrations are interesting for their simultaneous nature and first-person analysis of what the participant is doing. The behavior in this extra take is expected to differ from that of the non-narrated tasks, in that it is likely that the

participant will complete the scenario more slowly than normal to concentrate on explaining what they were doing. These narrations are available for about 10% of all takes in Ego-Exo4D, since we wanted participants to execute the tasks without pausing for the bulk of the recordings. They can potentially be used for multimodal learning as is currently explored in the literature with how-to video narrations (Miech et al., 2019, 2020, Lin et al., 2022, Ashutosh et al., 2023).

## 4.3 Atomic action descriptions

The third language dataset consists of *atomic action descriptions*. Whereas the commentary and narrate-and-act language reveals spoken opinions and reasons for the actions (the “why and how”), this stream of text is specifically about the “what”. Inspired by Ego4D’s *narrations* (Grauman et al., 2022), these are short statements written by third-party (non-domain expert) annotators, timestamped for every atomic action performed by the participant for all videos in the dataset, for a total of 432K sentences. This data is valuable for mining for taxonomies of objects and actions in the data, indexing the videos with

<sup>6</sup>For some activities which were more physically intense, such as dancing or bouldering, we asked participants to instead narrate either just before or just after the action to reduce the difficulty of doing this live.

keywords for exploring the dataset, and for future research in video-language learning, as has been quite successful for the Ego4D narrations (Lin et al., 2022, Pramanick et al., 2023, Ashutosh et al., 2023).

### **Annotation description**

We present each take to the annotators as a collaged video consisting of the egocentric view, left and right grayscale SLAM, four or five fixed-position exocentric cameras, and single-track composite audio; for a subset of videos, a helmet-mounted GoPro view is also available. Annotators are asked to provide a play-by-play description of what happens, as seen across *any* of the views. Potential contents include actions by the camera wearer, other individuals interacting with the camera wearer, and relevant environmental events.

Each narration is atomic and time-anchored: as much as possible, each narration should be limited to one verb and have a single associated timestamp, roughly within a second of its occurrence in the video. For consistency across narrations, and consistency with Ego4D’s narrations, the camera wearer in each take is referred to as “C” (e.g., “C picks up a wrench.”). Other individuals are referred to by other letters (e.g., “Man X kicks the soccer ball back to C.”); these letter labels are not necessarily consistent across takes, but refer to the same individual within a take. Many videos are narrated by two independent annotators, and we make both available. Figure 14 shows examples, and more are in Table D2 of the Appendix. See Table D3 in the Appendix for atomic action descriptions summary statistics.

### **Visibility labels**

Because of the multi-view nature of the Ego-Exo4D capture rig, certain actions or events may not be visible across all camera feeds. While we hope Ego-Exo4D leads to increased attention toward multi-view learning, many existing systems fundamentally assume a single view at a time; if a camera does not have a view of the narrated action or event, this may lead to a confusing learning signal, or pose an impossible ask for a model to infer.

Thus, we also ask that the annotators answer two additional question per narration: 1) an indicator of whether the narration is visible from the

egocentric camera, and 2) which (if any) of the static exocentric cameras provide the best view. If there are multiple equally good views, annotators are free to pick any. In particular, we found this *best exocentric view* helpful for other Ego-Exo4D annotation efforts: the narration visibility tags played a role in exocentric view selection for both the correspondences benchmark (Section 5.1) and expert commentary (above), and frame selection for hand and body pose (Section 5.4).

### **Comparing the language annotations**

How do the statistics from all three forms of language differ? Overall, expert commentary tends to use a much larger vocabulary and more lengthy statements, since commentators are giving more elaborate statements of advice and explanation. The temporal density of the atomic action descriptions is greater than the other two forms, since the annotators are pausing to describe every single action of the camera wearer. Narrate-and-act comments use a vocabulary size in between the other two, reflecting the more free-form speech (compared to the written atomic actions) is used. Across the different scenarios, the trends are mostly similar, with the most noticeable differences being the temporal density; it is particularly high for both cooking and soccer. In the former, there are many procedural steps, whereas in the latter there are many instances of the drill being executed.

See Figure D8 in the Appendix for the detailed statistics, and Figure D7 in the Appendix for word clouds per scenario and annotation type highlighting the differences in vocabulary and word frequency.

## **5 Ego-Exo4D benchmark tasks**

Our second major contribution is to define the core research challenges in the domain of egocentric perception of skilled activity, particularly when ego-exo data is available for training (if not testing). To that end, we devise a suite of foundational benchmark tasks organized into four task families: relation (Sec. 5.1), recognition (Sec. 5.2), proficiency (Sec. 5.3), and ego-pose (Sec. 5.4).

Benchmark	Annotation Type	Ego-Exo4D (v2)	
		Num Takes	Annotations
Relations	Manual	1335	5566 objects 742K ego masks 1.1M exo masks
Keystep recognition	Manual	1088	17 activities, 664 keysteps 27.6K ego segments (87h) 143K ego+exo segments (454h)
Procedure understanding	Manual	628	6 activities, 186 keysteps 8.6K segments (30h)
Proficiency estimation	Semi-automatic	2987	2987 proficiency scores (demonstrator)
	Manual	912	19K “good” segments 20K “tips” segments (demonstration)
Ego pose (Body)	Automatic	2559	9.2M 3D / 46.87M 2D
	Manual	1358	376K 3D / 2M 2D
Ego pose (Hand)	Automatic	976	4.3M 3D / 21M 2D
	Manual	458	68K 3D / 340K 2D

**Table 3:** Summary of annotation statistics for the different benchmark tasks of Ego-Exo4D.

For each task, we provide not only suitable multimodal data, but also high quality annotations that allow training and evaluating models, as well as baselines that provide a starting point from which the research community can build. Table 3 overviews the annotations provided with Ego-Exo4D, and Table 1 summarizes key distinctions with existing datasets.

We ran the first formal teaser Ego-Exo4D challenge in 2024, and will launch the full suite of challenges and leaderboards next year.

In the following sections, for each benchmark task we provide 1) the motivation and applications of solving that task, 2) the formal task definition, 3) key related work, 4) a description of the annotations, 5) an overview of the metrics we use to evaluate the task, 6) the design the baseline models, and 7) their results on the released dataset.<sup>7</sup>

Important: To ensure fair comparisons in any future work using Ego-Exo4D, researchers need to account for 1) the precise task input-output definitions and 2) the train/test/val splits available with the annotations. Specifically, for each task below, when formally defining the inputs and outputs, we also explicitly specify which inputs are *excluded* from use, if any.

## 5.1 Ego-exo relation

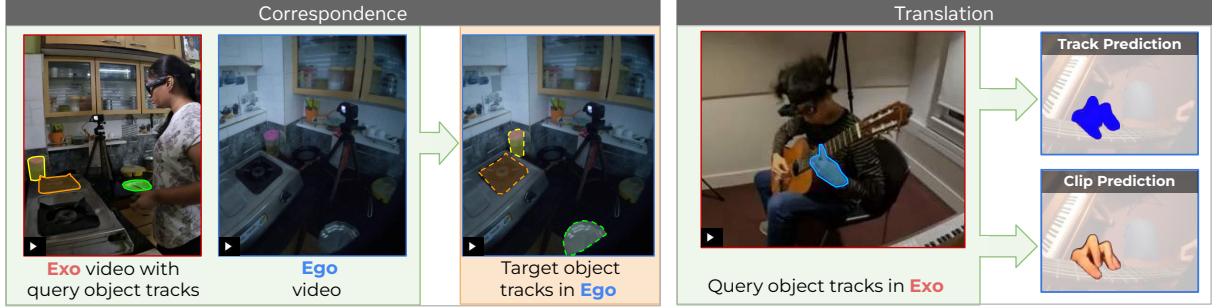
Our ego-exo *relation* tasks deal with relating the video content across the extreme ego-exo viewpoint changes. They take the form of object-level matching (correspondence) and synthesis of one view from the other (translation). See Figure 15.

### 5.1.1 Ego-exo correspondence

#### Motivation

Establishing object-level correspondences between ego and exo viewpoints would allow AI assistants to provide visual instructions by matching third-person observations of objects from instructional

<sup>7</sup>Note that all results are from “v2” of Ego-Exo4D released in March 2024. The smaller v1 is now considered obsolete and should not be used for future publications and comparisons.



**Fig. 15:** The ego-exo relation family consists of the tasks of correspondence (left) and translation (right).

videos to those in the user’s first-person view. Compared to the general correspondence problem, our setting requires tackling a number of challenges: extreme viewpoint differences, high degrees of object occlusion, and many small objects (e.g., cooking utensils and bike repair tools).

### Task definition

Given a pair of synchronized ego-exo videos and a sequence of query masks of an object of interest in one of the videos, the task is to identify the corresponding mask for the same object in each synchronized frame of the other view, if visible. See Figure 15, left. The task can be posed with query objects in either the ego or exo video, with both directions presenting interesting challenges (e.g., high degree of occlusion in ego views, and small object size in exo views). This task is especially challenging in our dataset, since we have to handle long videos with an average length of 3 minutes, as well as very small objects with areas of only a few pixels.

Importantly, the input to the model *excludes* semantic labels or names for the objects, camera pose information relating the two views, and IMU or active range sensor measurements. We do not use such information as we want to encourage the development of methods for open-world correspondence, not relying on predefined sets of objects or inputs that require non-consumer camera devices.

### Related work

Related tasks are image-level sparse correspondence given query points (instead of object masks) (Jiang et al., 2021) and image-level object co-segmentation (Vicente et al., 2011) for jointly segmenting semantically similar objects. Our task

goes beyond static object correspondence, since the interplay between human pose and object state changes during manipulation necessitate using temporal context and tracking as the query object can be highly occluded or blurry (Tang et al., 2023).

### Annotations

We annotate pairs of temporally synchronized egocentric and exocentric videos with segmentation masks for selected object instances from six scenarios: *Cooking*, *Bike Repair*, *Health*, *Music*, *Basketball* and *Soccer*. We exclude *Bouldering* and *Dance* from this benchmark as they have limited diversity of objects. We focus on objects used by the camera-wearer at any point during the execution of the activity and that are visible in both views for at least some frames of the sequence. These masks allow us to define object-level correspondence between the views. We used a multi-stage annotation process for annotating paired ego-exo videos. There are 1.8M masks annotated at 1fps covering 5.6k objects from 1335 takes. Overall, an average of 5.5 objects are annotated with correspondences between the two views in each take, with each object tracked for an average of 173 frames (excluding frames with occlusions). See Appendix E.1 for details and statistics.

### Metrics

We adopt the following metrics in our evaluation:

1. *Location Error* (LE), which we define as the normalized distance between the centroids of the predicted and ground-truth masks.
2. *Intersection Over Union* (IoU) between the predicted and ground-truth masks.

3. *Contour Accuracy* (CA) (Perazzi et al., 2016), which measures mask shape similarity after translation is applied to register the centroids of the predicted and ground-truth masks.
4. *Visibility Accuracy* (VA) (Brodersen et al., 2010), which evaluates the ability of the method to estimate the visibility of the object in the target view, as in practice it may often be occluded or outside the field of view. We measure this performance using balanced accuracy. Note that, in contrast to the previous metrics that compare segmentation masks at frames where the object is visible in both views, this metric is computed based on all frames with query masks.

### Baselines

Finding object mask correspondences across pairs of videos is an under-explored area in video understanding. Therefore, we investigate two diverse baseline approaches for our ego-exo correspondence task: (a) a *spatial model* that tackles the correspondence problem independently at each time point, and (b) a *spatio-temporal model* that takes into account the history of predicted correspondences.

- *Spatial baseline model.* This model receives as inputs an egocentric frame, the associated exocentric frame, and a query object segmentation mask in one of the views. It then outputs the mask in the other view (if the object is visible in that view). It can be thought of as a generalization of query-point correspondence approaches proposed for sparse image correspondence (Jiang et al., 2021). We implement this baseline in the form of a Transformer-based image correspondence model, *XSegTx* (Cross View Segmentation Transformer), which extends SegSwap (Shen et al., 2022), a method originally proposed for image co-segmentation, i.e., for segmenting common objects in a pair of images.
- *Spatio-temporal baseline model.* The spatio-temporal model receives as input the pair of ego-exo video clips as well as an object segmentation track in one of the views, and outputs segmentation masks in the other view for the frames that the object is visible in both views. It can be thought of as performing generalized

tracking across views. We build our baseline model on top of XMem (Cheng and Schwing, 2022), a model originally proposed for tracking a specific target object given its segmentation mask in the first frame. In particular, our baseline model, called *XView-XMem*, adapts XMem to track the object across different views given ground-truth segmentation masks for one of the views in each frame.

See Appendix E.1.1 for implementation details.

### Results

We benchmark our XSegTx and XView-XMem baseline models on the test set in Table 4. We experiment with two settings: providing the ground-truth object track in the exo view (exo query mask) and predicting it in the ego view, and vice versa.

First, we observe that exploiting temporal cues helps with tackling the object correspondence task as shown by the significant increase in performance achieved by the spatio-temporal baselines (ST type) compared to the spatial ones (for example, IoU improves from 13.88% to 22.14% in the Ego→Exo setting).

Second, we can see a big difference in performance between the Ego→Exo and Exo→Ego settings for all the baselines. In particular, models perform worse when the sequence of query masks is provided for the egocentric video and the model needs to predict query masks in exocentric video. This might be due to the heavy occlusion and very small size of objects in the exocentric views, making segmentation very challenging. While predicting a very tiny mask in the exo view can be very difficult, models can reason about the type and rough location of the object from a tiny mask in the exo view and thus accurately detect and segment it in the ego view, where it is much larger.

However, all our baselines achieve a performance inferior to 23% IoU in the Ego→Exo setting and inferior to 24% IoU in the Exo→Ego setting. This shows the challenging nature of the task and the dataset. We note that our dataset includes a great degree of object shape variation and high number of very small objects which are very difficult to model.

We also show some qualitative results in Fig. 16. As we can see, the spatial baseline (XSegTx) struggles to track the same object

throughout the video. For example, in the bottom example, XSegTx alternates between predicting one and two object masks whereas the spatiotemporal baseline (XView-XMem) reliably tracks a single object throughout the sequence, showing the importance of exploiting temporal cues in the data. See Appendix E.1.1 for more analysis and visualizations of the results.

### 5.1.2 Ego-exo translation

#### Motivation

The second of the two ego-exo relation tasks is *ego-exo translation*. Our translation task entails synthesizing a target ego clip from a given exo clip. This problem may be viewed as a form of ego-exo correspondence with missing information: given the masks of an object in the exo clip, its corresponding masks and pixel values must be generated in the *unobserved* ego clip. We note that this problem cannot be solved by image-based rendering or via geometric transformations, since the dramatic differences in ego-exo viewpoints cause the two cameras to capture different regions of the same objects.

We believe this problem will drive novel research for combining recognition and object synthesis. For example, in Figure 15 (right), the approach must perceive the hand as the generation target and make effective use of the hand’s object-specific shape and appearance priors in order to synthesize the ego view of the fingertips—which are not visible in the exo clip. Furthermore, this task will stimulate advances in visual odometry, as the method must be able to infer the ego camera pose from the third-person clip.

Ego-exo translation also holds strong application potential, as it may unlock the ability to generate first-person renderings of videos that were originally captured from a third-person perspective. For example, it may enable AR coaching, with objects and interactions lifted from a third-person instructional videos and re-synthesized from the camera-wearer perspective to better guide the user in the execution of complex activities. Ego-exo translation may also be used to generate abundant first-person training data from existing large-scale collections of third-person videos in order to train robot perception models.

#### Task definition

The translation benchmark focuses on generating information in the egocentric view given the exocentric view. We decompose ego-exo translation into two separate tasks: *ego track prediction* and *ego clip generation* (Figure 15, right). Ego track prediction estimates the segmentation mask of an object in the *unobserved* ego frames given the object masks in the observed exo clip. Ego clip generation entails generating the image values (i.e., RGB) within the given ground-truth ego mask by making use of the exo clip and the object masks in those frames. This decomposition effectively splits the problem into two tasks: 1) predicting the location and shape of the object in the ego clip, and 2) synthesizing its appearance given the ground-truth position. For both sub-tasks, the input exo clip consists of 5 frames evenly sampled from a time span of 5 seconds.

We believe that the decoupling of these two tasks will promote faster progress on the individual sub-problems and facilitate understanding of the key challenges in each of them. For each, we consider a variant where the pose of the ego camera with respect to the exo camera is available to use at inference time. This simplifies the problem but reduces the applicability of the method, since this information is typically not available for arbitrary third-person videos. Finally, we note that while the opposite direction of translation (i.e., ego-to-exo) could be considered, here we focus on the task of ego generation because of its higher value for robotics and AR applications.

Note that we restrict the input to include only the exo view and the object masks in order to promote the design of methods that can translate arbitrary third-person video into an egocentric one. Thus, the input *excludes* depth maps, 3D point clouds, IMU, or SLAM, which would simplify the task at the expense of general applicability, since these signals are typically not available for in-the-wild video. The only exception is a variant of the task where the ego camera pose for all frames of the clips is given as input. We consider this formulation in order to estimate a sort of “upper bound” on translation performance under the unrealistic assumption of known ego-exo camera relation.



**Fig. 16:** Qualitative results for the different ego-exo correspondence baselines.

Query Mask	Method	Type	Bal. Acc. $\uparrow$	IoU $\uparrow$	Location Score $\downarrow$	Contour Acc. $\uparrow$
Ego	XSegTx (random weights)	S	50.00	0.48	0.118	0.014
Ego	XSegTx	S	<u>66.31</u>	18.99	<u>0.070</u>	<u>0.386</u>
Ego	XMem (w/o finetuning)	ST	64.39	<u>19.28</u>	0.151	0.262
Ego	XView-Xmem (w/ finetuning)	ST	61.24	14.84	0.115	0.242
Ego	XView-Xmem (+ XSegTx)	ST	<b>66.79</b>	<b>34.90</b>	<b>0.038</b>	<b>0.559</b>
Exo	XSegTx (random weights)	S	50.00	1.08	0.203	0.024
Exo	XSegTx	S	<b>82.01</b>	<b>27.14</b>	<b>0.104</b>	<b>0.358</b>
Exo	XMem (w/o finetuning)	ST	60.35	16.56	0.160	0.204
Exo	XView-Xmem (w/ finetuning)	ST	<u>61.72</u>	21.37	0.139	0.269
Exo	XView-Xmem (+ XSegTx)	ST	59.71	<u>25.00</u>	<u>0.117</u>	<u>0.327</u>

**Table 4:** Baseline evaluation for the ego-exo correspondence benchmark on test set. Best results are reported in bold, whereas the second best results are underlined.

### Related work

Ego-exo translation relates to cross-view image synthesis ([Regmi and Borji, 2018](#), [Tang et al., 2019](#), [Lu et al., 2020](#)). Within this genre, the problem of exo-to-ego generation was recently introduced for both images ([Liu et al., 2020](#)) and video ([Cheng et al., 2024](#), [Luo et al., 2024](#), [Liu et al., 2021](#)), and approached using GANs or diffusion conditioned on the input view. Our work not only formalizes this task with ample data, but its formulation also draws attention to the need for a *semantic* basis to new view synthesis across extreme view changes.

### Annotations

Translation uses the same annotations as the correspondence task discussed above in Section 5.1.

### Metrics

We adopt a diverse set of metrics to assess the different aspects of the generated translation. As for the task of correspondence, we use *Visibility Accuracy* (VA) to evaluate the ability of the method to predict the visibility of the target object in the ego view but this time given only exo frames as input. We consider the visibility prediction correct if and only if either of these two conditions are met: (1) the predicted mask is empty when the object is invisible in the ego view, or (2) the predicted mask is non-empty when the object is visible in the ego view. Furthermore, we adopt the following metrics defined for correspondence to gauge the performance of *Ego Track Prediction*: 1) *Location Error* (LA) 2) *Intersection Over Union* (IoU) and 3) *Contour Accuracy* (CA) ([Perazzi et al., 2016](#)).

The IoU and CA are calculated after registering the centroids of the predicted mask and the ground-truth ego mask, in order to gauge mask prediction independent of location error. To evaluate *Ego Clip Generation* we use two popular image quality metrics (SSIM, PSNR ([Hore and Ziou, 2010](#))) and three perceptual metrics (DISTS ([Ding et al., 2020](#)), LPIPS ([Zhang et al., 2018](#)) and CLIP similarity ([Radford et al., 2021](#))).

### Baselines

For track prediction, we implement the GAN-based method pix2pix ([Isola et al., 2017](#)) and the NeRF-based method GNT ([Varma et al., 2023](#)). For clip generation, we employ the GAN-based method pix2pix ([Isola et al., 2017](#)) and the diffusion model DiT ([Peebles and Xie, 2022](#)). It is worth noting that, as discussed below, we introduce specific modifications to adapt these methods to our task requirements. All baselines utilize exo images and masks, with only the GNT model making use of the extra input of ego camera pose.

*Ego Track Prediction* involves generating segmentation masks for the egocentric view based on the exocentric video clip and the exocentric object masks. We consider the following two baselines for this task:

- **pix2pix-mask.** We modify the generator of pix2pix to have inputs and outputs of 4 channels. Specifically, the exo frame and the exo mask are concatenated as the inputs while the 4-channel outputs are ego frame (3 channels) and ego mask (1 channel). The ego frame is

supervised with the losses used in pix2pix. We use the bootstrapped cross-entropy loss (Reed et al., 2014) and the dice loss (Sudre et al., 2017) for mask prediction.

- **GNT-mask.** We adopt the Generalizable NeRF Transformer (GNT) (Varma et al., 2023) as another baseline leveraging the camera poses. In our adapted version, the image encoder takes a 4-channel image (exo frame and mask) as inputs to predict the ego frame and ego mask. Formally, during the training of our GNT-mask, for each point  $x$  and viewing direction unit vector  $d \in \mathbb{R}^3$ , the ray transformer  $f$  in GNT predicts two key attributes: RGB Color ( $c$ ) and Object Existence Score ( $e$ ), in which  $e$  signifies the probability of an object being present at point  $x$ . During rendering, the volumetric radiance field encoded by the ray transformer can then be rendered into a 2D image as well as a 2D object mask.

*Ego Clip Generation* requires producing pixel values representing the target object in the egocentric view. To achieve this, we leverage 6 different input images for each frame: exo frame, exo mask, exo object crop, cropped exo mask, ego mask and cropped ego mask. The cropped exocentric and egocentric masks are generated by considering a bounding box to isolate the relevant portions of the exocentric and egocentric masks, respectively. The “exo object crop” refers to the RGB image obtained by cropping out the relevant region using the cropped exocentric mask. We resize these 6 images to the same size ( $256 \times 256$ ). We evaluate two baselines for this task:

- **DiT-pix.** We adopt the Transformer-based diffusion model DiT (Peebles and Xie, 2022). We predict the ego object crop by conditioning the DiT on the 6 input images in two manners. Initially, these six images are concatenated along the channel dimension and subsequently combined with the noisy ego object crop, forming the input to DiT. Additionally, two ResNet-50 architectures encode the six images into low-dimensional features, which are then incorporated into each layer of DiT via AdaLN (Perez et al., 2018).

- **pix2pix-pix.** We adopt pix2pix (Isola et al., 2017) for clip generation as well by concatenating the 6 images along the channel dimension as inputs to the pix2pix model.

All of the above-mentioned baselines perform image-to-image generation. We implement also clip-to-clip variants of these methods by taking multiple frames as inputs and predicting results for all frames jointly. For pix2pix, we achieve this by replacing the original 2D-Conv with 3D-Conv, and 2D-BatchNorm with 3D-BatchNorm. For DiT, we use space-time divided attention as in TimeSformer (Bertasius et al., 2021).

## Results

We employ the validation set for the purpose of selecting optimal checkpoints and hyperparameters, which are subsequently evaluated on the test set.

For the task of Ego Track Prediction, both pix2pix-mask and GNT-mask perform poorly in estimating the object visibility, achieving Visibility accuracy around 50%, i.e., same as random guess (50.0% for GNT-mask and 56.2% for pix2pix-mask, on the v2 test set). However, the ResNet-50 trained exclusively to attend to this binary classification achieves a VA of 82.9% on the v2 test set. We assess mask quality by considering distance (Location Error) and similarity metrics (IoU and Contour Accuracy) between predicted and ground-truth masks after registration. The 3D-aware NeRF-based baseline, GNT-mask, outperforms the implicit baseline, pix2pix-mask, overall. However, it does so by exploiting the ego camera pose as additional input. It is noteworthy that both baselines perform poorly on this task, likely due to the inherent challenges in correctly predicting the location and shape of the target object in the ego view, probably due to the fact that it often has diminutive size in the exo view.

In the case of Ego Clip Generation (Table 5b), the Diffusion model DiT-pix demonstrates superior performance across all metrics compared to the GAN-based pix2pix-pix. Qualitative results (Figure 17a) illustrate that DiT-pix can generate highly photorealistic images, aligning closely with the ground-truth in most instances. However, there are occasional cases (the last 2 rows) where

**Table 5:** Results on exo to ego translation task.

(a) Evaluation of translation baselines for the sub-task of ego track prediction.

Method	Ego Cam.	Location Pose	Contour Error↓	IoU ↑ Acc.↑
pix2pix-mask	No	21.6	4.5	5.3
+multi-frame	No	20.1	3.1	3.5
GNL-mask	Yes	<b>19.6</b>	<b>15.5</b>	<b>10.3</b>

the shape of the object is accurately generated, but the texture deviates slightly.

We further verify the importance of each input in Figure 17b. Without exo object crop as input, the model fails to correctly infer the color and texture of the target object in the ego view. This result is expected as the source objects often represent a very small region of the entire exo frame. Additionally, without the ego crop mask as input, the model predicts the orientation of the object incorrectly. These observations highlight the importance of the cropped inputs.

We can observe in Table 5a that multi-frame (i.e., clip-to-clip) prediction does not provide a quantitative advantage over frame-to-frame prediction. Yet, we noticed that the multi-frame variant often yields generations that are more consistent across frames, even for frames where the exo view is heavily occluded, as can be seen in Figure 17c. This is reasonable as a clip-level model can more effectively learn about the target object from multiple frames and fill-in information that is missing in individual exo frames.

Please see Appendix E.1.2 for a break down of ego-exo translation results across different scenarios.

## 5.2 Ego-exo keystep recognition

This family of tasks centers around recognizing the keysteps of a procedural activity and modeling their dependencies. Specifically, there are three tasks: fine-grained keystep recognition (Sec. 5.2.1), efficient multimodal keystep recognition (Sec. 5.2.2), and procedure understanding (Sec. 5.2.3). We refer to the family of tasks as “ego-(exo)” since exo may be available at the time of training but not inference. See Figure 18.

(b) Evaluation of translation baselines for the sub-task of ego clip generation.

Method	SSIM ↑	PSNR ↑	DISTS ↓	LPIPS ↓	CLIP ↑
pix2pix-pix	0.42	16.4	0.36	0.50	79.8
DiT-pix	0.59	16.1	0.31	0.46	81.9

### 5.2.1 Fine-grained keystep recognition

#### Motivation

Recognizing the step a camera wearer is performing is non-trivial: keysteps in the same activity may look similar and may involve differentiating subtle differences in hand-object interactions with heavy occlusions and head motion. Models with access to multiple views during training can leverage their complementarity to account for the deficiencies of each one, by learning viewpoint invariant representations or distilling multi-view signals into a single model (e.g., human hands from ego; body pose from exo).

#### Task definition

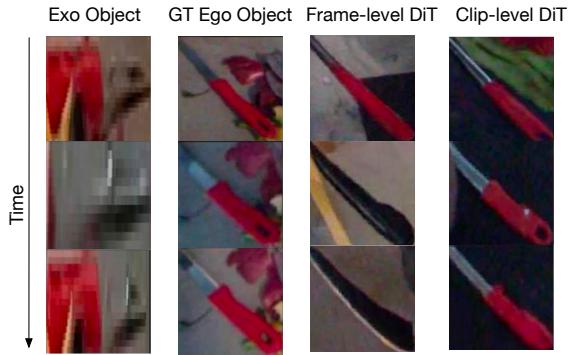
We consider trimmed video clip classification as the keystep recognition task. At training time we are given a labeled collection  $\mathcal{D}$  of ego-exo video clips:  $\mathcal{D} = \{(\mathcal{V}_{ego}^{(1)}, \mathcal{V}_{exo^{1-M}}^{(1)}, y^{(1)}), \dots, (\mathcal{V}_{ego}^{(N)}, \mathcal{V}_{exo^{1-M}}^{(N)}, y^{(N)})\}$  where  $y^{(n)}$  denotes the keystep label of the  $n$ -th sample. The video clips are manually trimmed from long procedural videos to contain only the keysteps to recognize. At test time, given *just* the ego view of a trimmed clip  $\mathcal{V}_{ego}$ , the model must predict its keystep label  $y$ .

Classification of trimmed video clips is a problem formulation commonly adopted in action recognition benchmarks (Kay et al., 2017, Soomro et al., 2012, Goyal et al., 2017). However, our task differs from action recognition in three fundamental aspects. First, it targets fine-grained keystep recognition rather than classification of coarse activities. We note that this adds significant complexity, since different keysteps of an activity often involve manipulating the same objects in the scene (e.g., folding the bedsheet and smoothing



(a) Qualitative ego clip generation by DiT-pix on the test set. The model takes 6 input images (exo frame, exo crop, exo mask, exo crop mask, ego mask, and ego crop mask). The ego frame, serving solely as a reference, does not constitute either an input or an output element.

(b) Qualitative demonstration of the importance of the different inputs given to DiT-pix. The exo crop image and ego crop mask are critical for good performance.



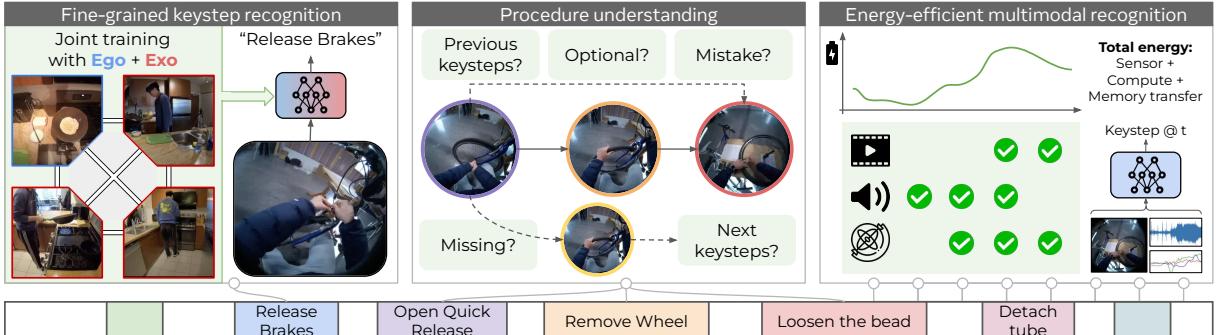
(c) Comparison of ego clip generations using frame-to-frame vs clip-to-clip variants of DiT-pix. The clip-to-clip version of the model produces outputs that are more coherent across the frames of the clip, even for frames where the exo view is heavily occluded.

**Fig. 17:** Qualitative results for exo-to-ego translation task.

out the bedsheet) and are consequently difficult to tell apart. Second, different keysteps may be represented over largely different time spans (e.g., the average time span for “kneading dough” is 87.3 seconds, in stark contrast with “getting salt”, which averages at 3.6 seconds), thus requiring analysis at different levels of temporal granularity. The third key difference is the potential to leverage contextual cues available in exocentric videos during training to improve the prediction accuracy on egocentric videos.

Note that at test time, the input to the model includes just the ego-view videos (RGB only). Exo-view videos, activity and scenario names, narrations, audio and associated metadata such

as eye gaze, 3D point clouds, camera pose, and IMU information are *excluded* as inputs for inference. Intuitively, these additional modalities could provide valuable contextual cues, such as environmental awareness from exo-view videos, semantic meaning from narrations, or attentional signals from eye gaze, which could help the model better understand the visual content of the ego-view videos and improve its keystep recognition performance. We encourage exploring their potential utility in training to leverage these benefits, but for the purpose of evaluation, we restrict the input to RGB video only at test time to ensure our approach remains vision-centric.



**Fig. 18: Ego-exo keystep recognition.** This family of tasks consists of fine-grained recognition (left, Section 5.2.1), procedure understanding (center, Section 5.2.3), and energy-efficient multimodal recognition (right, Section 5.2.2).

### Related work

Keystep recognition has been studied in first-person (Sigurdsson et al., 2018, Ragusa et al., 2021, Bansal et al., 2022, Song et al., 2023) or third-person (Tang et al., 2020, Zhukov et al., 2019, Zhou et al., 2018, Ashutosh et al., 2023, Mavroudi et al., 2022) videos; however, limited work considers both views together. Prior work considers cross-view learning with unpaired videos (Ardesir and Borji, 2018, Xue and Grauman, 2023, Li et al., 2021) and view-invariant feature learning on paired videos (Sigurdsson et al., 2018). In contrast, we explore keystep recognition in large-scale, procedural activities with fully synchronized training videos.

### Annotations

We annotate videos featuring any of the three procedural activities (i.e., cooking, bike repair, health) with temporal segments of *keysteps*, i.e., actions that contribute towards the completion of a procedural task. Each keystep annotation contains the start and end timestamps, a category label, a natural language description (e.g., “add dried herbs” or “fit the tire onto the bike”), and a flag indicating whether the keystep is essential or optional for task completion. To accurately model the hierarchical nature of the activities, we also develop a hierarchical keystep taxonomy concurrently with the annotation process, in an iterative, data-driven manner. In total, we annotate 143,442 segments, spanning 664 keysteps across 17 activities. Figure 20 shows example keystep

annotations, highlighting the challenges of fine-grained keystep recognition where subtle differences in hand-object interactions and contextual cues are crucial for distinguishing between activities. Complete details on the annotation interface and taxonomy development are in Appendix E.2.

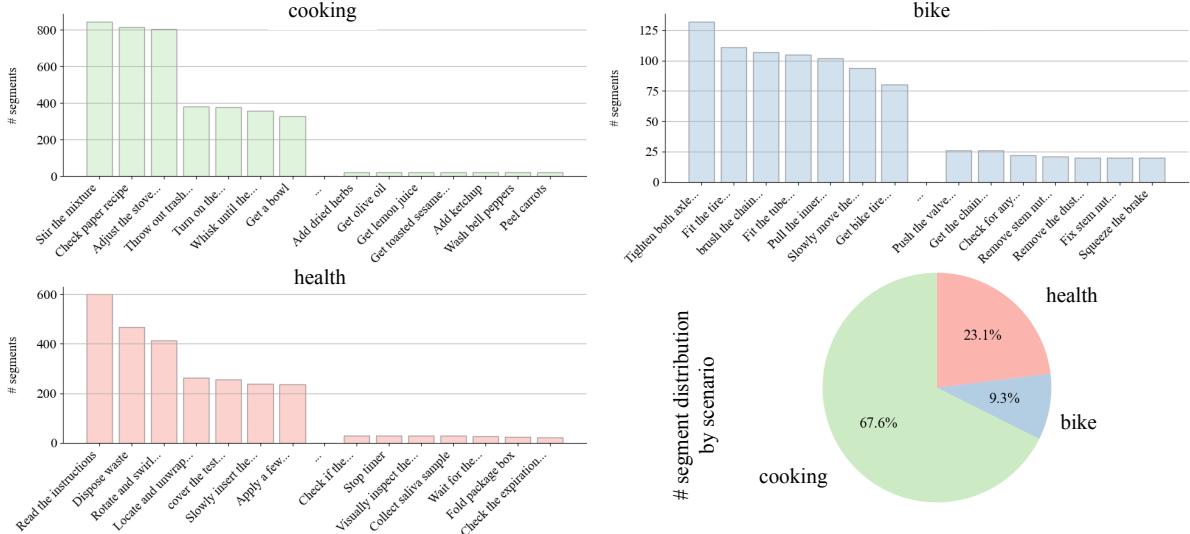
### Metrics

We report top-1 accuracy for evaluation. Since the keysteps in our dataset exhibit a very long-tailed distribution, we set a cutoff threshold at 20 samples per keystep, limiting our analysis to 278 unique keysteps as shown in Figure 19. Some of these keysteps are illustrated in Figure 20. Dataset split details are in Appendix E.2.

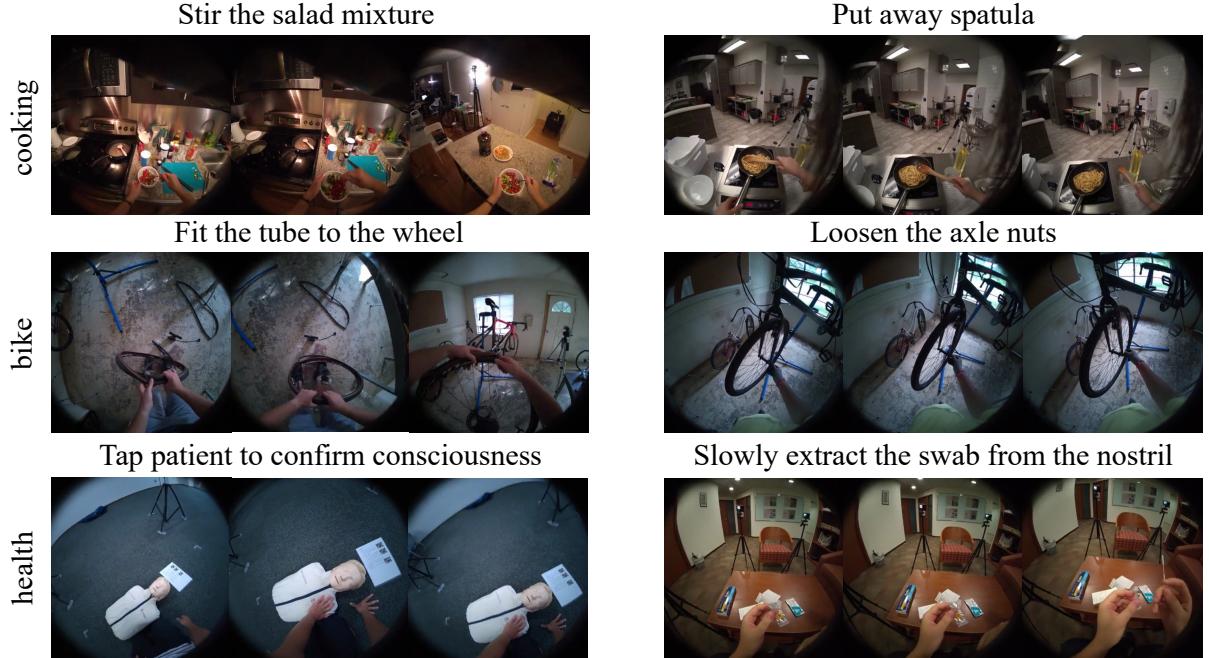
### Baselines

To understand the best strategy for egocentric keystep recognition with paired ego-exo training data, we consider a diverse set of baselines approaches, including methods for action classification, video representation learning, and ego-exo transfer.

- **Action classification.** As a prototypical example of this classic genre, we select a TimeSformer (Bertasius et al., 2021) model initialized with the checkpoint pretrained on the large-scale third-person action dataset Kinetics-600 (Kay et al., 2017) due to its strong performance in various video understanding tasks.
- **Video-language pretraining.** We adopt the EgoVLPv2 framework (Pramanick et al.,



**Fig. 19:** Keystep distribution in our dataset for each procedural scenario: cooking, bike repair, and health.



**Fig. 20:** Example keysteps from cooking, bike repair, and health scenarios. Keystep labels are displayed above each frame sequence.

2023) and pre-train the model jointly on the Ego4D (Grauman et al., 2022) (which contains only ego views) and the Ego-Exo4D datasets (which encompasses both ego and exo views).

We balance the number of samples between these two datasets by augmenting Ego-Exo4D with LaViLa-style (Zhao et al., 2023) narration

rephrasing.

- **View-invariant learning.** A two-stage training approach is employed. In the first stage, we utilize all available (ego, exo) video pairs in the dataset for training a view-invariant (VI) encoder. The training objective is a clip-level contrastive loss (Oord et al., 2018), aiming at identifying the synchronized (ego, exo) pairs as positive, and the non-synchronized pairs as negative. In the second stage, this pretrained model is further trained with a classification loss, aligning with the clip-level classification nature of the downstream task. Note that to align with the clip-level classification task, our contrastive loss operates at the clip-level, rather than at the frame-level as was done in view-invariant loss proposed in (Sigurdsson et al., 2018, Sermanet et al., 2017).
- **Viewpoint distillation.** This also adopts a two-stage training approach. In the first stage, we train a multi-view teacher that takes both ego and exo views as input. In the second stage, a single-view ego student is trained, distilling knowledge (Hinton et al., 2015) from the multi-view teacher to encapsulate information from both views.
- **Ego-exo transfer.** Here we follow the methodology proposed in Ego-Exo (Li et al., 2021) which uses egocentric pseudo-labels to pre-train the network. We employ a masked autoencoder (MAE) (Tong et al., 2022) backbone, initialized from a Kinetics checkpoint, and the pseudo-labels provided from the Ego-Exo checkpoint to fine-tune with two auxiliary heads (Object-Score and Interaction-Map). We then further finetune the model with a classification loss for fine-grained keystep recognition.

For the first two baselines (which utilize pre-trained checkpoints from well-established benchmarks), two training settings are further examined: one using only the ego view for the classification loss and the other utilizing both ego and exo view videos, with the training objective being the sum of ego view and exo view classification losses. Implementation details are in Appendix E.2.

## Results

Table 6 reports the Top-1 accuracy for ego classification on both validation and test sets. Among all the baselines, the VI Encoder emerges as the top performer, achieving a test accuracy of 41.53%. It is closely followed by Viewpoint Distillation and EgoVLPv2 pretrained on EgoExo4D, which attain test accuracies of 39.49% and 38.76% respectively. These results open discussion on how to effectively utilize exo videos during training to enhance ego keystep recognition during test time.

First, we note that different approaches respond differently to the addition of exo-view videos during training. Specifically, while the TimeSFormer (K600) exhibits a degradation when the exo classification loss is integrated into the objective (i.e., test accuracy drops from 35.93% to 31.04%, EgoVLPv2 pretrained on EgoExo4D benefits from the introduction of exo-view videos (i.e., test accuracy improves from 37.72% to 38.76%). This enhancement is also evident in the VI encoder and viewpoint distillation when compared to TimeSFormer (K600) that only utilizes ego-view videos for training. These observations suggest that certain baselines are better equipped at leveraging exo information during training to improve ego keystep recognition.

Finally, in Appendix Figure E16, we show a breakdown of performance by viewpoint. In short, we find that ego views are more informative for steps involving manipulation of small objects, like ‘cut carrots’ and ‘unpack the new tube’, while exo views show advantages for keysteps like ‘have a conversation asking different questions’. Overall, we posit that the endeavor to enhance view-invariant learning and to more effectively harness the complementary information from exo views for ego keystep recognition remains an open avenue. Our findings underscore the need for further investigation and innovation in this domain.

### 5.2.2 Energy-efficient multimodal keystep recognition

#### Motivation

Current activity detection models assume access to densely sampled clips from the full video and ample computational resources to process them. These assumptions are incompatible with real-world devices (e.g., mobile phones, AR glasses)

Method	Train data	Ego Accuracy (%)	
		Val	Test
TimeSFormer (Bertasius et al., 2021) (K600)	ego	35.13	35.93
TimeSFormer (Bertasius et al., 2021) (K600)	ego,exo	32.68	31.04
EgoVLPv2 (Pramanick et al., 2023) (Ego4D)	ego	36.51	37.55
EgoVLPv2 (Pramanick et al., 2023) (Ego4D)	ego,exo	35.84	36.59
EgoVLPv2 (Pramanick et al., 2023) (EgoExo4D)	ego	36.04	37.72
EgoVLPv2 (Pramanick et al., 2023) (EgoExo4D)	ego,exo	<u>39.10</u>	38.76
VI Encoder (Oord et al., 2018) (EgoExo4D)	ego,exo	<b>40.34</b>	<b>41.53</b>
Viewpoint Distillation (Hinton et al., 2015)	ego,exo	38.19	<u>39.49</u>
Ego-Exo Transfer MAE (Li et al., 2021)	ego,exo	37.17	36.58

**Table 6:** Top-1 accuracy of keystep recognition on val and test data. The pre-training dataset is denoted in parentheses. VI is short for view-invariant.

where the camera is not always on and the compute budget is limited by battery life. This task focuses on building energy-efficient video models to pave the way for feasibility on real-world hardware.

### Task definition

Whereas the keystep recognition task (presented in Sec. 5.2.1) entails classifying keystep video clips in batch without regard for energy costs, in this task, the goal is to perform *online* classification of keysteps in a streaming egocentric multi-modal video, within an energy budget. We consider an ego video  $\mathcal{V}_{ego}$  of arbitrary length  $T$  comprising a stream of  $K$  different sensory modalities (e.g., RGB images, audio, IMU, etc.). At each time step  $t$ , where  $1 \leq t \leq T$ , the video consists of samples for each available modality, such that  $\mathcal{V}_{ego}^t = \{S_1^t, \dots, S_K^t\}$ , where  $S_j^t$  denotes the sample at time  $t$  for the  $j^{th}$  modality.

Given  $\mathcal{V}_{ego}$  and an energy budget  $B$ , our task is to learn a model  $\mathcal{F}$  that maximizes the overall keystep recognition performance across the full video while also ensuring that the combined energy for sensing and running model inference does not exceed  $B$ .  $\mathcal{F}$  consists of a sensor triggering policy  $\mathcal{F}^P$  and a keystep prediction model  $\mathcal{F}^K$ . At every step  $t$ , the policy  $\mathcal{F}^P$  decides which sensor(s) to activate and sample from, in order to produce the model’s current observation  $O^t$ , such that  $O^t \subseteq \{S_1^t, \dots, S_K^t\}$ . Given  $O^t$ , the keystep predictor  $\mathcal{F}^K$  outputs its estimate of the ground truth keystep for the current step.

The energy budget accounts for the cost of operations in each model forward pass, the cost of moving intermediate activations in and out of

memory and the cost of the continuous operation of sensors, each having different cost profiles (e.g., IMU and audio sensors are relatively cheaper to operate than camera sensors). Note that the sensor triggering policy may be static (e.g., sample video at 4 frames per second (fps), keep audio/IMU off; sample 1 fps video, keep IMU always on) or dynamic (e.g., depending on the audio, decide whether to trigger video capture). We keep our task definition general, allowing the challenge to admit a wide variety of recent approaches ranging from pure video-based efficient backbone architectures (Feichtenhofer, 2020) to multi-modal triggering approaches and, naturally, a combination of them.

Note that at test time, the input to the model can only include current and past observations as our task is strictly an online recognition task. However, we encourage exploring modalities beyond those considered in our experiments, e.g., IMU or camera poses inferred from video, audio, and IMU.

### Related work

Prior work on efficient models considers light-weight architectures (Feichtenhofer, 2020, Vasu et al., 2023, Howard et al., 2017, Zhang et al., 2018, Tan et al., 2019, Mehta and Rastegari, 2021), efficient input processing (Gao et al., 2020, Korbar et al., 2019, Ghodrati et al., 2021, Meng et al., 2020, Tan et al., 2023), or inference optimizations (Iandola et al., 2016, Esser et al., 2019, Polino et al., 2018, Zhu and Gupta, 2017, Wu et al., 2018). In all cases, they optimize computation (FLOPs), parameter count, or prediction throughput (FPS), which in isolation are

insufficient to characterize running on real-world devices. To address this, we propose the first benchmark for *energy-efficient* video recognition that is tied to real-world, on-device constraints, and measures total power consumed.

### Annotations

This task uses the same egocentric videos and annotations as keystep recognition. However, in addition to the raw RGB video, it uses the audio stream (and potentially other sensors) as another sensor modality.

### Measuring energy consumption

Accurately measuring energy consumption of models is crucial for their use in AR/VR devices (Abrash, 2021, Chen et al., 2019). The energy used comes from a complex interplay of sources including sensors, compute, communication, data processing, memory transfer (SRAM and DRAM), and leakage – many of which are typically ignored when building *efficient* computer vision models, despite their large energy consumption (e.g., memory transfer accounts for over 50% of the total power (Yang et al., 2022)).

We consider three key factors when modeling energy consumption following prior work (Sze et al., 2020). (1) Compute energy: the cost of each model forward pass as a function of the number of operations (MACs). (2) Memory transfer energy: the cost associated with memory read-write operations for storing intermediate activations and model outputs. (3) Sensor triggering energy: the cost associated with turning on / off and continuous operation of sensors (camera, audio, IMU). For a model that processes an observation  $O^t$ , the total energy consumed can then be formulated as:

$$E(O^t) = \alpha * C(O^t) + \beta * M(O^t) + \sum_{j=1...K} \gamma_j * \mathbb{1}(S_j \in O^t) \quad (1)$$

where  $C(O^t)$  corresponds to the total number of multiply-add operations computed during the forward pass (in MAC/s),  $M(O^t)$  corresponds to the total memory transferred to/from DRAM (in MB/s), and  $S_j \in O^t$  corresponds to whether the  $j$ -th sensor is active. Finally,  $\alpha, \beta, \gamma_j$  are weighting factors that measure the contribution of each

energy source. We select these weighting parameters to reflect real-world AR/VR hardware capabilities. Namely,  $\alpha = 4.6 \text{ pJ/MAC}$  (Sze et al., 2020, Desislavov et al., 2023);  $\beta = 80 \text{ pJ/byte}$  (Horowitz, 2014);  $\gamma_{rgb} = 15 \text{ mW}$  and  $\gamma_{audio} = 0.5 \text{ mW}$  (Liu et al., 2020).

We adapt off-the-shelf profiler software built for PyTorch to compute the quantities in Eqn. 1 – the energy consumption expressed as power (mW). Details about the profiler are in Appendix E.3.

### Metrics

Following prior work (De Geest et al., 2016), we evaluate online keystep detection performance using per-frame calibrated mean average precision (mcAP), which accounts for the imbalance in the keystep labels in our dataset. We measure energy consumption in mW as described above. There is a natural trade-off between efficiency and better performance. Thus, we evaluate models in two tiers by setting a budget for the power consumption in each tier, namely 20 mW for the *high-efficiency* tier and 2.8W for the *high-performance* tier, selected based on existing efficient architectures. More details about the tiers are in Appendix E.3.

### Baselines

We provide a family of (less/more expensive) keystep prediction models for solving the task. Each model has a unimodal or audio-visual feature encoder followed by a keystep classification head. Experimental setup and implementation details are in Appendix E.3.

- **X3D-XS (Feichtenhofer, 2020).** This is a vision-only model comprising the X3D-XS feature encoder, which progressively expands the feature size and representational capacity of its layers, and later contracts them for achieving better performance-efficiency trade-off.
- **LaViLa (Zhao et al., 2023).** This is another vision-only model where the visual feature encoder is trained through CLIP-style video-language pre-training.
- **Light-ASDNet (Liao et al., 2023).** This is an audio-only model that represents audio as spectrograms and efficiently encodes them

Method	Modality	mcAP (%) ↑	Power (mW) ↓
Light-ASDNet (Liao et al., 2023) + $s = 5$	A	65.18	19.67
X3D-XS (Feichtenhofer, 2020) + $s = 10$	V	76.85	<b>19.14</b>

(a) *High-efficiency tier* (budget = 20 mW).

Method	Modality	mcAP (%) ↑	Power (mW) ↓
Lavila (Zhao et al., 2023) + $s = 5$	V	<b>93.24</b>	<b>2245.66</b>
AV-LF w/ Lavila + $s = 5$	AV	92.18	2274.40

(b) *High-performance tier* (budget = 2.8W).

**Table 7:** Keystep prediction results.

by splitting 2D convolutions into 1D convolutions along the spectrogram temporal dimension (Liao et al., 2023).

- **Audio-Visual Late Fusion (AV-LF).** This is an audio-visual model that does late fusion of visual features (encoded with X3D-XS or LaViLa) and audio features from Light-ASDNet by using linear layers.

To improve the energy efficiency of the aforementioned keystep predictors, we employ the following baseline policies for determining when to sample or skip each modality:

- **Fixed stride.** This is a policy that samples the input (video or audio) every  $s$  prediction steps. We evaluate different  $s$  values, where  $s$  ranges from 2-150 steps.
- **AV-LF + greedy.** This is a policy that greedily uses up the budget by sampling both audio and vision as early as possible, and uses the AV-LF backbone for keystep prediction.
- **AV-LF + random.** This is a policy that randomly samples or skips the audio and/or visual inputs until it runs out of budget, and uses the AV-LF backbone for prediction.
- **Audio-Visual (AV) Cascade.** This is a policy that initially uses the Light-ASDNet model to predict the keystep, and switches over to the LaViLa model if the audio-based prediction confidence is below a confidence threshold of 0.5.

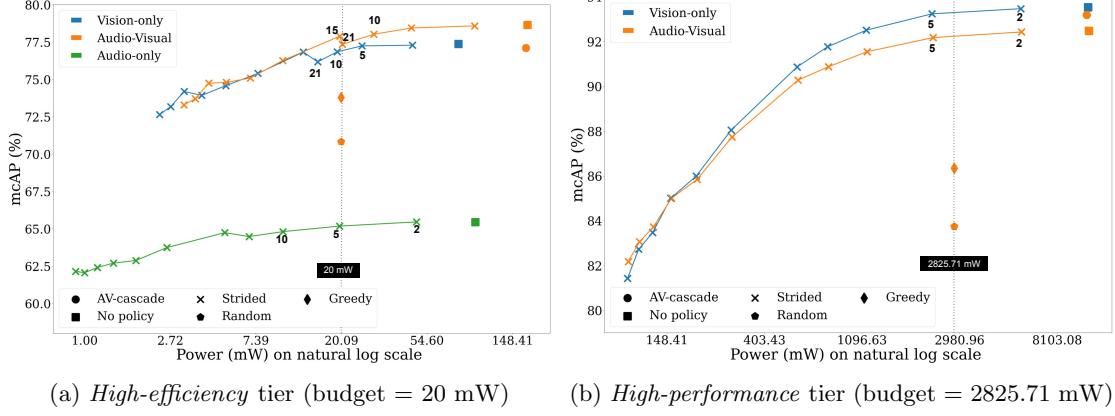
## Results

In Fig. 21a, we plot the recognition mcAP of all models against their total power consumption for the *high-efficiency tier*. We can see that combining vision and audio is better than using only vision or audio. Thus suggests that the two modalities carry complementary cues that are useful for the task. However, all vision-only models outperform their audio-only counterparts, which indicates that vision is the most critical modality for the task. The raw backbones generally perform better than the models using a sampling policy, but at the cost of requiring higher energy, making them impractical to use in online settings. Among the models that use a fixed stride, a lower stride generally improves the performance while hurting energy efficiency. Using the greedy or random policy with AV-LF leads to a sharp decline in performance compared to using a fixed stride, showing that sampling very early or randomly in the episode is suboptimal for our online recognition task. AV-cascade also performs worse than most audio-visual models while also requiring more energy, possibly because the audio backbone often outputs wrong but over-confident predictions that prevent switching over to the more reliable vision backbone when required.

For easy reference, in Table 7a we report the recognition performance and total power consumption of our best uni-modal and audio-visual models within budget for the *high-efficiency tier*.

In Fig. 21b, we plot the recognition mcAP of all models against their total power consumption for the *high-performance tier*. Different from the high-efficiency tier, the audio-visual backbone generally performs worse than the vision-only backbone, possibly because the LaViLa features are strong enough by themselves, and fusing them with audio features through the simple mechanism of linear late fusion reduces their expressivity. Otherwise, the overall behavior of different sampling policies is similar across the two tiers. We report the recognition performance and total power consumption of the best uni-modal and audio-visual models within the *high-performance* budget in Table 7b.

Finally, in Appendix Figure E17, we present a breakdown of performance by keystep labels and the behavior of audio- and vision-only models across them. In short, we find audio-only models



**Fig. 21:** Keystep prediction performance (mcAP) vs. total power consumption with different prediction backbones and sampling policies for both *high-efficiency* (left) and *high-performance* (right) tiers. For the models using a fixed stride, we show their stride value in text if their total energy consumption is close to the budget.

have an affinity for sounding actions like *stir fry egg mixture* and *cut butter*.

### 5.2.3 Procedure understanding

#### Motivation

The procedure understanding task consists in inferring the underlying structure of a procedure from the observation of natural videos of subjects performing the procedure.

The real-world motivation for our procedure understanding task has basis in augmented reality (AR), robotics, and more in general in assistive systems. Indeed, automatically understanding the *structure* of a procedure from video, e.g., inferring keystep orderings and preconditions, will allow to assist or guide users carrying out the procedure through AR or to allow robots to learn from human demonstrations. Beyond recognizing the current keystep, an assistive system could verify that some mandatory keysteps are missing, suggest possible future ones, and detect procedural mistakes. Similarly, robots could learn the structure of a procedure from human demonstrations. Mining the structure of procedures has been shown useful for planning (Chang et al., 2020, Bi et al., 2021) and improving keystep recognition (Ashutosh et al., 2023, Zhou et al., 2023) keystep discovery (Bansal et al., 2022), and for procedural mistake detection (Seminara et al., 2024).

#### Task definition

Figure 18 (center) illustrates the proposed task. Given a video segment  $s_i$  and its segment history  $S_{i-1} = \{s_1, \dots, s_{i-1}\}$ , models have to 1) determine *previous keysteps* (i.e., keysteps which should be performed before  $s_i$ ); infer if  $s_i$  is 2) *optional* (i.e., it can be omitted without compromising the correct execution of the procedure) or 3) a *procedural mistake* (a keystep which should not have been performed in that moment due to missing pre-conditions); 4) predict *missing keysteps* (i.e., key-steps which should have been performed before  $s_i$ ); and 5) forecast *next keysteps* (i.e., keysteps for which dependencies are satisfied and hence which could be executed next).

The task is weakly supervised, with two versions based on the level of supervision: 1) instance-level: video segments and their keystep labels are available during training and inference, similar to an action recognition task; 2) procedure-level: unlabeled video segments and a taxonomy of procedure-specific keystep names are given for training and inference. Note that, being weakly supervised, in both cases, explicit information on the structure of the procedure—such as the occurrence of mistakes or lists of pre-conditions—are not available for training. Also note that, when the procedure-level supervision is considered, the input to the model *excludes* keystep labels both at training and test time. At both the procedure and instance levels of supervision, models are required

to process the video in a causal fashion, meaning that predictions made at time  $t$  only depend on observations made at time  $t' < t$ .

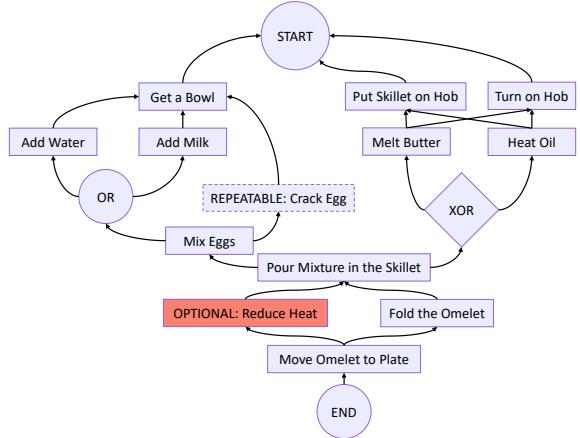
### Related work

Prior work focusing on procedural understanding learns an explicit graph (Jang et al., 2023, Xu et al., 2020, Soran et al., 2015) as ground truth or uses a task graph for representation learning (Narasimhan et al., 2023, Ashutosh et al., 2023, Zhou et al., 2023) and short-term step understanding (Dvornik et al., 2022, Ashutosh et al., 2023, Zhou et al., 2023). Other work (Sener et al., 2022, Ding et al., 2023) studies mistake detection in a supervised setting. We are the first to propose procedural understanding to evaluate the long-term structure of the task in a weakly-supervised setting.

### Annotations

For this task, we considered the following procedures: *i.e.*, *Covid-19 Rapid Antigen Test*, *Fix a Flat Tire - Replace a Bike Tube*, *Remove a Wheel*, *Install a Wheel*, *Clean and Lubricate the Chain and First Aid - CPR*. These scenarios represent structured activities with clear procedural constraint, yet allow a certain degree of variability in correct task executions. For each of the considered procedures, we manually labeled task-graphs as structures encoding the keystep orderings leading to a correct execution of the procedure (detailed below). A task graph is meant as a way to encode all orders of keysteps which lead to a correct execution of the task.

*Task-graphs.* We define a task-graph as a directed graph in which nodes represent keysteps and directed edges represent dependencies. For instance, in the example task-graph reported in Figure 22, the “Add Milk → Get a Bowl” structure denotes that keystep “Get a Bowl” has to be executed before keystep “Add Milk”. If a keystep has more than one dependency, all of them need to be satisfied. For instance, both “Put Skillet on Hob” and “Turn on Hob” need to be executed before “Heat Oil”. Besides directed edges, task-graphs also contain “OR” and “XOR” structures, which combine dependencies logically, as well as “optional” and “repeatable” node attributes. For instance “Mix Eggs” can be performed if either



**Fig. 22:** Example task-graph of a ”Cooking Omelet” procedure.

“Add Water” or “Add Milk” (or both) are executed, whereas “Pour Mixture in the Skillet” requires either “Melt Butter” or “Heat Oil” to be executed, but not both. Repeatable nodes (e.g., “Crack Egg”) can be repeated as long as their outgoing nodes (pre-conditions) are satisfied and incoming nodes (future nodes) are not executed. For instance, one could keep cracking eggs as long as the bowl is in place, but not after mixing the eggs. An optional node (e.g., “Reduce Heat”) can be omitted, but, if included, it needs its preconditions to be correctly satisfied.

*Task-graph construction.* We first familiarized ourselves with the procedural tasks by watching videos with annotated keysteps. We then initialized task graphs with procedural dependencies obtained from keystep annotations through the following procedure: a) a directed graph is first generated from the observed keystep transition frequencies; b) edges of the transition graph are filtered based on transition probabilities using a threshold parameter which is manually tuned for each scenario; c) edge directions are inverted to convert frequent transitions into dependencies. These initial graphs were then refined and manually corrected.

*Segment-level annotations.* A task graph is a global representation of a procedure including information on dependencies and partial orderings of keysteps. Since our task is defined at the keystep level, we need to “project” the constraints expressed by the task graph onto keystep

video segments, which we do with the following procedure.

Let  $S = \{s_1, \dots, s_n\}$  be a labeled sequence of keysteps in a given video. We denote with  $y_i$  the annotated keystore label of segment  $s_i$  and with  $Y_{:i} = \{y_1, \dots, y_i\}$  the sequence of labels up to the  $i$ -th keystore. Using these keystore annotations, each segment  $s_i$  is automatically matched to a task-graph and augmented with the following attributes: 1) a list of *previous keysteps*—these are the in-neighbors of the matched node, 2) *optional* labels—directly derived from the optional node attribute, 3) a *procedural mistake* label—this is set to “true” if the in-neighbors of the matched node do not correspond to segments in the history  $Y_{:i}$ , 4) the list of *missing keysteps*—the in-neighbors of the matched node not listed in  $Y_{:i}$ , and 5) the list of *next steps*—nodes for which in-neighbors appear in  $Y_{:i}$ . Non-repeatable nodes are listed only if they do not appear in  $Y_{:i}$ .

Given the weakly supervised nature of the task, we only release keystore level annotations on the validation set, while annotations on the training set are not shared nor used for the development of the baselines, and test labels are private, with evaluations on the test set possible by submitting predictions to a server. Also note that we do not release the labeled task graphs to avoid leaking test and training labels.

### Metrics

We consider the task of determining lists of keysteps as a detection task with an imbalance between positives (the keysteps to be detected, e.g., preconditions) and negatives (the keysteps which are not to be detected, e.g., keysteps which are not preconditions) and evaluate all methods using the calibrated Average Precision (cAP) (De Geest et al., 2016). Note that, according to this measure, a random baseline would on average achieve a performance of 50%.

### Baselines

We consider graph-based and end-to-end baselines. Graph-based baselines (see Figure 23(a)) include two main components: a segment-keystore assignment module (A1) which provides a pseudo-labeling of video segments based on a pre-trained video-language model, and a procedural reasoning module (B) which makes predictions based on

a transition graph built either from ground truth (for instance-level supervision) or pseudo-labels (for procedure-level supervision). Additionally, we provide a baseline where the keystore assignment step is replaced with label predictions from the Keystore Recognition task (A2). Note that in the training set, segments with a confidence score below 20% have been discarded.

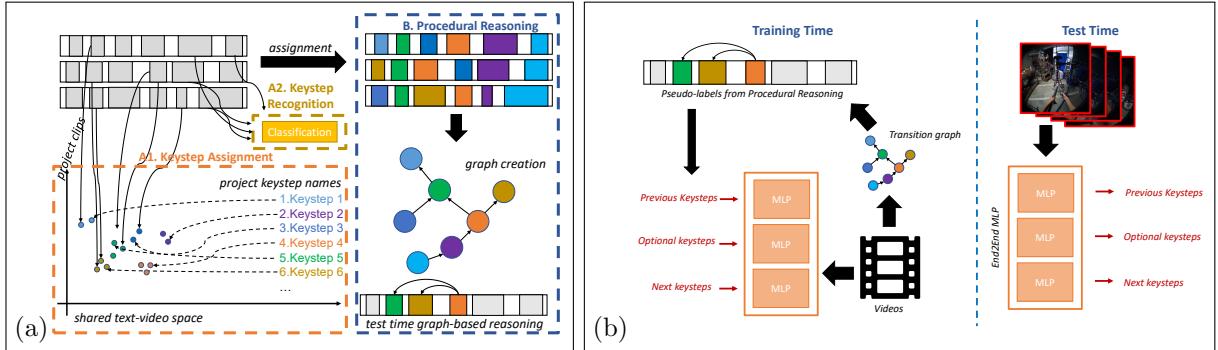
End-to-end baselines (see Figure 23(b)) are trained to predict the same results as graph-based baselines directly from video, with the aim to obtain a compact algorithm which does not explicitly make use of a graph. The end-to-end architecture consists of three MLPs designed to predict previous, optional, and future keysteps. Each MLP has six heads, one for each considered scenario.

Additional baseline implementation details are provided in Appendix E.4.

### Results

Table 8 reports the results obtained by our baselines and compares them against those produced by a “uniform” baseline, predicting previous/optional/mistakes/missing/next keysteps with equal probabilities. Results show that the graph-based baseline relying on ground truth annotations significantly outperforms the uniform baseline for most of the tasks, excluding future keystore predictions. This suggests that even simple keystore co-occurrences are informative to some degree of the overall structure of the procedure.

The limited performance gains on future keystore prediction highlight the complexity of the task and the need for further research. The end-to-end model trained with instance-level supervision achieves lower or similar performance, trading accuracy for test-time efficiency, due to the absence of an explicit graph. Procedure-level baselines achieve lower results because they do not rely on ground truth labels. The keystore prediction approach achieves better results compared to the keystore assignment mechanism for all tasks, except for optional keysteps. Despite our efforts, performance is below the uniform baseline, indicating that there is room for future investigations.



**Fig. 23:** Overview of the two procedure understanding approaches considered in our evaluation: (a) graph-based baselines for procedure understanding rely on a Keystep Assignment or a Keystep Recognition and a Procedural Reasoning component; (b) the architecture of our end-to-end baseline.

Supervision	Baseline	Keystep Labels	Inf. Set	Prev. Keysteps	Opt. Keysteps	Proc. Mistakes	Miss. Keysteps	Fut. Keysteps
-	Uniform Baseline	-	Val/Test	59.18/59.13	56.71/56.73	60.54/60.66	65.58/65.64	65.65/65.65
Instance-Level	Graph-Based	Ground Truth	Val/Test	82.49/82.32	58.95/62.10	73.19/73.06	84.29/82.63	63.48/62.82
Instance-Level	End-to-End	Ground Truth	Val/Test	62.05/62.05	51.85/61.39	56.75/52.07	60.11/61.77	60.35/59.25
Procedure-Level	Graph-Based	Keystep Assignment	Val/Test	54.26/53.43	49.86/52.36	56.46/57.81	60.97/53.92	52.50/53.54
Procedure-Level	End-to-End	Keystep Assignment	Val/Test	55.37/54.82	<b>52.12/60.78</b>	52.84/54.73	56.11/53.75	<b>58.88/57.47</b>
Procedure-Level	Graph-Based	Keystep Prediction	Val/Test	<b>64.56/66.22</b>	49.51/49.00	<b>61.15/58.59</b>	<b>61.50/64.18</b>	<b>57.87/58.34</b>
Procedure-Level	End-to-End	Keystep Prediction	Val/Test	57.43/57.92	<u>51.54/61.01</u>	51.68/54.92	54.99/55.15	57.35/56.92

**Table 8:** Results for the procedure understanding task. Best results are reported in bold, the second best results are underlined. All results are in percentage.

### 5.3 Ego-exo proficiency estimation

#### Motivation

Going beyond recognizing what a person is doing, this task aims to infer the user’s skill level. Such an ability could lead to novel coaching tools that let people learn new skills more effectively, or new ways to *evaluate* human performance in domains like sports or music.

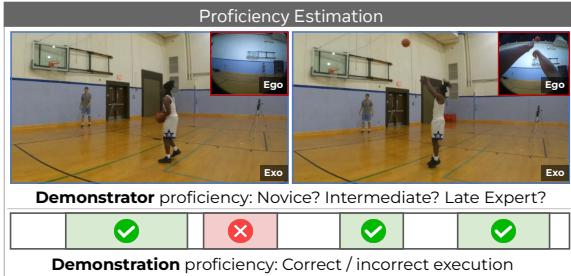
#### Task definition

We consider two variants: (1) *demonstrator* and (2) *demonstration* proficiency estimation. Both tasks consider one egocentric and (optionally)  $M$  exocentric videos of a demonstrator performing a task, which are synchronized in time, as their inputs:  $\mathcal{V} = (\mathcal{V}_{ego}, \mathcal{V}_{exo}^1, \dots, \mathcal{V}_{exo}^M)$ . See Figure 24 for an illustration. We provide more details for each variant below.

*Demonstrator proficiency estimation:* The goal is to estimate the demonstrator’s skill level from one or more task demonstrations. It is formulated as a video classification task with the following classes: (novice, early expert, intermediate expert, late expert).

*Demonstration proficiency estimation:* Given a single task demonstration, the goal is to identify parts of the video where the task execution was good (i.e., ‘good executions’) or needs further improvement (i.e., ‘needs improvement’). It is formulated as a temporal localization task, where we localize instances of ‘good executions’ and ‘needs improvement’ throughout the task demonstration. Formally, we can express the demonstration proficiency estimation function  $h$  as  $\hat{G}, \hat{I} = h(\mathcal{V})$ , where  $\hat{G} = \{t_1^g, t_2^g, \dots, t_{|\hat{G}|}^g\}$  are the timestamps where the participant shows good task execution, and  $\hat{I} = \{t_1^i, t_2^i, \dots, t_{|\hat{I}|}^i\}$  are the timestamps where the participant needs to improve their skill level. Note that parts of the video that do not reveal the participant’s skill are left unlabeled.

Both tasks inherently benefit from multi-view data. Egocentric video captures fine-grained information about the hand pose and object interactions, which can be critical in tasks such as cooking (e.g., chopping vegetables) and music (e.g., placement of fingers on the guitar). On the other hand, the exocentric videos provide broader



**Fig. 24:** Demonstrator and demonstration proficiency estimation.

information about the demonstrator’s body pose, which can be highly indicative of proficiency in tasks that require extensive physical motion such as basketball, soccer, and dancing.

Note that the input to the model excludes textual descriptions/narrations of the activity, audio, gaze sensor readings, and any subject information, which would simplify the task significantly at the expense of usability since these signals are typically not available for in-the-wild video. Our formulation encourages the development of proficiency estimation methods from visual cues.

### Related work

Prior work uses egocentric (Bertasius et al., 2017, Doughty et al., 2019) or exocentric (Parmar and Morris, 2017, Parmar and Tran Morris, 2019, Ismail Fawaz et al., 2018) views for proficiency estimation in sports (Pirsavash et al., 2014, Bertasius et al., 2017, Parmar and Tran Morris, 2019), health (Ismail Fawaz et al., 2018, Liu et al., 2021, Zhang and Li, 2013, Zia et al., 2017), and others (Doughty et al., 2019, Yu et al., 2021). We propose the first multi-view egocentric and exocentric proficiency estimation benchmark. Unlike prior work, our benchmark spans diverse, day-to-day physical and procedural scenarios and includes temporally localized annotations of (in)correct executions.

### Annotations

We now describe the annotation procedure for the two proficiency estimation tasks.

*Demonstrator proficiency estimation.* We assign four proficiency labels (novice, early expert, intermediate expert, late expert) to each person performing activity demonstrations (one label per person). Most levels correspond to experts

	Demonstrator			Demonstration		
	Train	Val	Test	Train	Val	Test
Basketball	575	143	167	146	47	19
Bike repair	-	-	-	41	9	15
Cooking	200	53	86	80	24	39
Dance	380	127	148	80	35	27
Health	-	-	-	42	12	16
Music	138	36	71	94	32	35
Rock Climbing	561	159	230	65	15	22
Soccer	129	43	66	8	3	6
Total	1983	561	768	556	177	179

**Table 9:** Distribution over video takes in proficiency estimation benchmark.

since Ego-Exo4D videos are dominantly targeted towards expert participants who can perform the task successfully (see Appendix B). Four proficiency classes makes the task challenging but still approachable.<sup>8</sup> We derive annotations for this task from participant surveys and expert commentary. Please see Appendix E.5 for more details, including a visualization of the proficiency score distribution for each scenario (Figure E18).

We split our dataset into train/val/test splits based on the common split shared across benchmarks. The dataset statistics are shown in Table 9. Note that we exclude the bike repair and health scenarios from the demonstrator proficiency task. The distribution of participants for bike repair is heavily skewed towards late experts. The predominant activity in the health collection is COVID testing, where skill levels are hard to determine due to the simplicity of the task.

*Demonstration proficiency estimation.* We leverage temporally localized annotations that include the timestamps of steps demonstrated in the video as well as the proficiency category for each demonstrated step instance (i.e., good execution or needs improvement). For this task, we consider all 8 scenarios, as shown in Table 9. We derive annotations for this task from expert commentary, where task experts carefully analyze videos and provide timestamped commentary on the participant’s performance (see Section 4.1). In particular, given a single timestamped comment from an expert, we annotate whether the comment describes a good execution and/or provides tips for improving the participant’s skill level. See Table E8

<sup>8</sup>Subtle variations between five or more levels of proficiency can be insufficiently observable from vision alone, and even difficult for expert annotators to reach consensus.

in Appendix E.5 for example annotations. These annotations are then associated with the timestamp provided with each comment to obtain a list of timestamps for good executions  $\{t_1^g, t_2^g, \dots\}$  and tips for improvement  $\{t_1^i, t_2^i, \dots\}$  in each video. Overall, the demonstration proficiency estimation task consists of 556 train / 177 val / 179 test videos (see Table 9 for a breakdown per scenario).

### Metrics

For demonstrator proficiency estimation, we measure performance using top-1 classification accuracy. For demonstration proficiency estimation, we measure the temporal localization performance using a modified mean average precision (mAP). Unlike prior temporal action localization methods which use temporal IoU between segments, we use the  $L_1$ -distance between the predicted and ground-truth timestamps to measure mAP. Therefore, we define mAP based on  $L_1$ -distance (in seconds) thresholds.

### Baselines

Next, we define the baselines for each task.

*Demonstrator proficiency estimation:* We adopt TimeSformer (Bertasius et al., 2021) for our experiments. We train one model on the egocentric view (“ego model”), and a separate model on all 4 exocentric views (“exocentric model”). The models are trained to classify individual clips using the cross-entropy loss. At inference time, we perform late fusion to incorporate information from both egocentric and exocentric video streams. We average the softmax predictions across both egocentric and exocentric models to obtain the final video label prediction. We also average results over three spatial crops during inference following prior work (Bertasius et al., 2021).

*Demonstration proficiency estimation:* We adopt ActionFormer (Zhang et al., 2022), a video action localization model for our experiments. Unlike traditional action localization, we infer only a single timestamp since our annotations contain only a single point in time for each good execution or tip for improvement. We accordingly adapt ActionFormer for timestamp regression and define  $L_1$ -distance based mAP metrics. We

Method	Pretraining	Accuracy					
		Ego		Exos		Ego + Exos	
		Val	Test	Val	Test	Val	Test
Random	-	26.4	26.4	26.4	26.4	26.4	26.4
Majority-class	-	32.3	42.4	32.3	42.4	32.3	42.4
TimeSFormer	-	40.6	33.9	39.0	<b>47.5</b>	39.9	45.7
TimeSFormer	K400	<b>47.2</b>	44.1	37.8	47.0	40.3	46.1
TimeSFormer	HowTo100M	45.1	36.7	39.8	46.6	<b>43.7</b>	47.0
TimeSFormer	EgoVLP	44.7	43.8	<b>40.5</b>	44.5	39.4	43.5
TimeSFormer	EgoVLPv2	46.7	<b>50.4</b>	37.0	47.0	37.1	<b>48.7</b>

Inference with multiple takes per demonstrator  
TimeSFormer EgoVLPv2 48.3 51.0 36.0 47.3 43.1 49.1

**Table 10: Demonstrator proficiency estimation.**  
We report top-1 accuracies for various baselines on the demonstrator proficiency estimation task. Our learned models use the TimeSFormer architecture (Bertasius et al., 2021).

train our models with Omnivore features (Girdhar et al., 2022) extracted from overlapping time intervals in the video. For the experiments involving multiple views (i.e., multiple exo views or ego + exo views), we simply concatenate the features for all views at each time step.

Please see Appendix E.5 for additional implementation details about the baselines.

### Results

Our Ego-Exo4D dataset has 5 views (1 egocentric view, and  $M = 4$  exocentric views). We run the proficiency tasks in two settings: one where the exo view is available at test time, and one where it is not. For the latter, benchmarking baseline models with only the egocentric view is important when the target is augmented reality applications like wearable headsets and mobile robotics. For the former, results with benchmarking with both egocentric and exocentric views are helpful to capture the multi-view aspect of the problem.

*Demonstrator proficiency estimation.* We present results for demonstrator proficiency estimation in Table 10. We include two naïve baselines to account for biases in the dataset. The random baseline uniformly samples one skill level at random. The majority-class baseline predicts the majority class within each scenario. TimeSFormer trained from random initialization outperforms the naïve baselines by a significant margin, demonstrating the ability of learned methods to quantify skill levels from videos.

Ego videos are sufficient to achieve good performance in most cases, while the exo videos are beneficial in tasks such as bouldering, highlighting the complementary nature of the ego and exo viewpoints. Initializing TimeSFormer using pre-trained weights improves over random initialization, particularly on ego videos. Furthermore, fusing the predictions from the ego view and exo views does not improve performance, likely due to the simplicity of late fusion. In the last row of Table 10, we further report results when providing multiple demonstrations from a participant for evaluating TimeSFormer. This matches or outperforms evaluating TimeSFormer on a single demonstration, highlighting the potential for obtaining more accurate skill estimates by studying multiple demonstrations. We further study scenario-specific performance of the baselines in Appendix E.5. We find that egocentric views are beneficial for scenarios such as cooking that require close-up views of hands and objects, whereas exocentric views are more useful for scenarios such as bouldering that require body-pose information. Overall, our benchmark presents new challenges for video-based skill understanding and our results highlight the difficulty of the task, suggesting good scope for improvement in future work.

*Demonstration proficiency estimation.* We present results for the demonstration proficiency estimation task in Table 11. We include three naïve baselines along with ActionFormer (Zhang et al., 2022). The “Random tips/good exec.” baseline randomly predicts a tip or a good execution label every 5.97 seconds, i.e., the average temporal span between adjacent annotations in our dataset. The “Uniform tips” baseline predicts a tip for improvement label every 5.97 seconds. The “uniform good exec.” baseline predicts a good execution label every 5.97 seconds. We evaluate ActionFormer models trained on ego only, exo only and ego + exo views. All naïve baselines perform poorly on this task. The learned ActionFormer baseline outperforms the naïve baselines by a good margin. However, the absolute mAP scores are fairly low, suggesting that the task is very challenging and has a significant scope for improvement in methods.

## 5.4 Ego pose

### Motivation

Having presented benchmark tasks about ego-exo relation, recognition, and proficiency assessment, we now define the final family of tasks centered on body and hand pose. This family of tasks is motivated by recovering the skilled body movements of participants, even in the extreme setting of monocular ego-video input in dynamic environments, as shown in Figure 25. Estimating the physical state of a person’s body—the 3D positions of the arms, legs, hands—from the ego view is essential for wearable AI systems that can support human activity. Challenges include subtle and flexible movements, frequent occlusion, and body/hand parts out of view.

For each scenario, we invite experts as the participants to enhance the complexity and variety of the motions captured. As an example, expert musicians typically demonstrate more advanced and varied finger techniques (300+) compared to beginners or intermediate players (< 100) during the recordings. Such complexity enables the model to (1) extract more representative latent features, and (2) learn subtle patterns and relationships that might be missed in a more homogeneous dataset, for both estimation and prediction tasks.

### Task definition

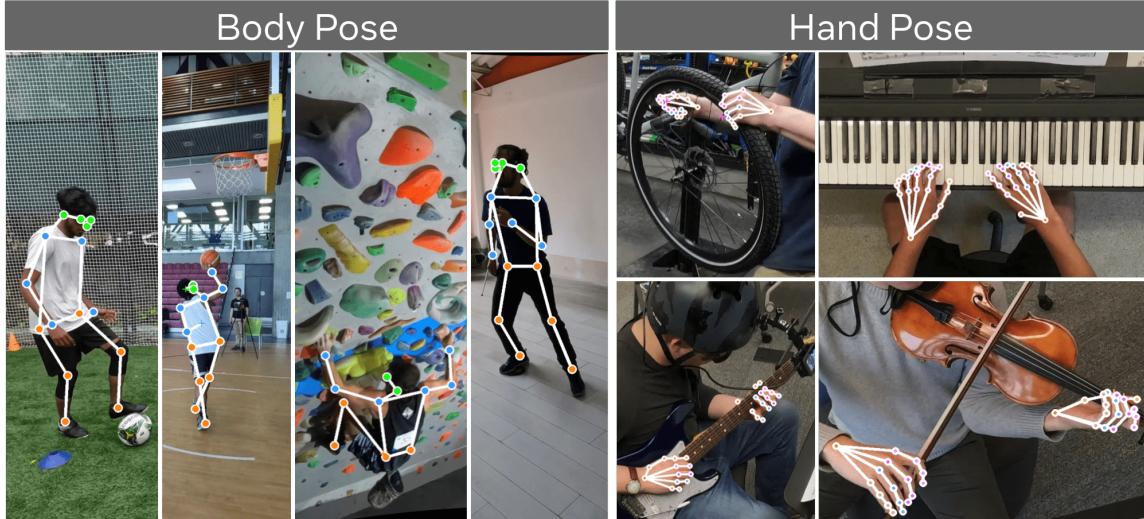
The ego pose benchmark is divided into two separate tasks: *body pose estimation* and *hand pose estimation*.

In our *body pose estimation* task, the goal is to estimate the 3D human pose sequence  $\mathcal{P} = \{\mathcal{P}_1, \dots, \mathcal{P}_T\}$  using either an egocentric RGB video input sequence  $\mathcal{V}_{ego} = \{\mathcal{V}_1, \dots, \mathcal{V}_T\}$ , an IMU sensor sequence  $\mathcal{H}_{imu} = \{\mathcal{H}_1, \dots, \mathcal{H}_T\}$ , or both, where  $1 \leq t \leq T$  is the given time window remapped to make the starting timestamp to be 1, and  $P_t \in \mathcal{R}^{17 \times 3}$  correspond to the 17 joints following the MS COCO convention in 3D.  $T$  can have different values depending on the length of the particular annotated clip. Note that at test time, we only estimate the error across the visible annotated joints at frame  $t$ .

The *ego hand pose* task entails predicting the three-dimensional coordinates of the camera wearer’s hands. Given an egocentric frame, the goal is to estimate the 3D joint location for the hands that are (at least) partially visible in the ego

Method	Val/test results								
	mAP <sub>0.25</sub>	Ego mAP <sub>1.0</sub>	Avg.	mAP <sub>0.25</sub>	Exos mAP <sub>1.0</sub>	Avg.	mAP <sub>0.25</sub>	Ego + Exos mAP <sub>1.0</sub>	Avg.
Random	0.48/0.45	5.23/4.72	2.46/2.20	0.48/0.45	5.23/4.72	2.46/2.20	0.48/0.45	5.23/4.72	2.46/2.20
Uniform tips	0.49/0.45	5.28/5.18	2.48/2.39	0.49/0.45	5.28/5.18	2.48/2.39	0.49/0.45	5.28/5.18	2.48/2.39
Uniform good exec.	0.43/0.46	4.79/4.62	2.27/2.17	0.43/0.46	4.79/4.62	2.27/2.17	0.43/0.46	4.79/4.62	2.27/2.17
ActionFormer (Zhang et al., 2022)	<b>0.95/1.04</b>	<b>6.33/7.56</b>	<b>3.27/3.87</b>	<b>1.08/1.14</b>	<b>7.50/7.36</b>	<b>3.84/3.87</b>	<b>0.97/1.14</b>	<b>7.03/7.90</b>	<b>3.57/4.04</b>

**Table 11: Demonstration proficiency estimation benchmark.** We report the mean average precision (%) for various baselines on the demonstration proficiency estimation task for the val and test splits. mAP<sub>k</sub> is measured at an  $L_1$ -distance threshold of  $k$  seconds. The average mAP (Avg.) measures the mAP averaged across  $k = \{0.25, 0.5, 1.0\}$  seconds.



**Fig. 25:** Hand and body keypoints for ego-pose estimation.

view. The output is parameterized as 21 3D joints per hand following the MS COCO dataset convention (Lin et al., 2014). Frames from the ego view are extracted and undistorted for both training and evaluation. The 2D hand bounding boxes are generated by projecting the 3D hand joints onto the 2D image planes and subsequently enclosing these projections.

Since the Ego Pose benchmark is aimed at promoting the development of methods that perform body pose estimation solely from first-person raw video or IMU data, the input *excludes* egocentric modalities that would unfairly simplify the task (e.g., audio captured from a wearable camera, eye gaze), as well as exocentric video or any signals that can be extracted from it.

### Related work

Limited prior work explores 3D body pose from a wearable camera. Some methods assume no body visibility (Li et al., 2023, Jiang and Grauman, 2017, Yuan and Kitani, 2018, 2019, Luo et al., 2021), while others assume partial observability by modifying cameras to capture the body (Rhodin et al., 2016, Tome et al., 2019, Xu et al., 2019, Ahuja et al., 2019, Hwang et al., 2020). Our dataset can be used for both paradigms.

Existing hand pose datasets use constrained environments (Simon et al., 2017, Moon et al., 2020) with simple hand motion (Kwon et al., 2021, Ohkawa et al., 2023, Hampali et al., 2020), whereas we include diverse real-world scenarios with skilled hand motions, e.g., with expert musicians and bike mechanics.

### Annotations

The *3D human body pose* annotation process consists of two main stages: (1) automatic ground truth generation, and (2) manual multi-view keypoint annotation/correction. Through this process we derive 3D keypoint annotations for approximately 14M frames.

In the automatic ground truth generation phase, we use off-the-shelf models ([MMPoseContributors, 2020](#)) to predict the 2D bounding boxes from each of the exocentric views. Since there could be multiple people in the scene and we only want to consider the one wearing the egocentric camera, we project the 3D headset location from the MPS output to select which box corresponds to the camera wearer. Then, we run an off-the-shelf 2D human keypoint detector ([MMPoseContributors, 2020](#)) for each bounding box to obtain the 2D keypoints. Finally, we run 3D triangulation with RANSAC to minimize the reprojection errors to obtain the 3D keypoints for the camera wearer. In the manual annotation phase, we import the undistorted frames and the reprojected 2D keypoints into our multi-view annotation interface.

The *3D human hand pose* annotation process also consists of two stages, i.e., the automatic ground truth generation and the manual multi-view keypoint annotation. Compared to the body pose, the main difference in automatic ground truth generation is that we also detect hand keypoints from the egocentric frame, and we use the result from the whole body pose estimation to infer the hand locations when there are multiple people in the scene. Similarly, for manual annotation, besides the exocentric frames, we also show the annotators the egocentric frames to allow them to annotate/correct hand keypoints. For each annotated joint in manual annotations, we provide the number of views used for triangulation as the indicator of the confidence for the provided ground truth data. Meanwhile, the correction the annotators make for hand joints on ego images can serve as the indicator to understand the difficulty for hand reconstruction from the given ego view.

Ego-Exo4D offers the largest available manually annotated body pose (376K 3D/2M 2D) and hand pose (68K 3D/340K 2D) annotations. Along with this, we also provide 9.2M/47M (body) and 4.3M/21M (hand) automatically generated

groundtruth 3D and 2D poses, totaling about 13.M frames. In total, we have approximately 14M frames of 3D ground truth (GT) and pseudo-GT combined across body and hands. To our knowledge, this represents the largest collection of body pose annotations in the literature, whether for ego or exo video.

How good is the auto GT? Between manual and automatic annotations, the body and hand MPJPEs are 3.33 cm and 1.87 cm, respectively, much smaller than the best baseline methods. It is important to note that Ego-Exo4D tackles real-world scenarios with five or fewer cameras rather than controlled environments. This introduces challenges like increased occlusions from body and objects along with limited view and resolution of hands from distant cameras. Despite this, our auto generation pipeline surpasses baselines, showcasing robustness and efficacy. Experiments below further show performance boosts across baselines when using automatic ground truth, demonstrating its effectiveness. Note that automatic GT and manual GT are not mutually exclusive, and people can choose whether/how automatic GT is used for training.

### Metrics

To evaluate the performance of *body pose estimation* approaches we calculate the Mean Per Joint Position Error (MPJPE) in centimeters (cm), and the Mean Per Joint Velocity Error (MPJVE) in meters per second (m/s).

The *ego hand pose* baselines are evaluated according to both the MPJPE and the PA-MPJPE metrics. The MPJPE measures absolute Mean Per Joint Position Error, while the PA-MPJPE calculates the average 3D joint errors after performing Procrustes Alignment on hand poses. Both metrics are reported in millimeter (mm) unit.

### Baselines

We evaluate three state-of-the-art baseline methods for the *body pose estimation* task. Moreover, to gauge the performance of deep-learning-based methods, we create a static pose baseline, which consists of fixing the 3D human body pose prediction to be the average pose in the training set and translating it according to the IMU sensor. Thus, the fixed prediction matches the camera location at each frame.

- **Kinpoly.** Kinpoly (Luo et al., 2021) proposes to use a simulated humanoid to track head pose and create full-body motion based on action types. Based on the input head-pose and action type, Kinpoly synthesizes realistic human pose and human-object interactions inside a physics simulator. Different from kinematic-based methods that directly output joint angles or positions for pose estimation, Kinpoly outputs joint torques as the final product and controls a simulated humanoid for pose estimation.
- **EgoEgo.** EgoEgo (Li et al., 2023) uses a two-step approach for egocentric body pose estimation, by estimating the head pose from the egocentric video first, and then using a diffusion model to generate the full body motion sequence based on the head pose sequence. For head pose estimation, it obtains the initial head pose trajectory using DROID-SLAM (Teed and Deng, 2021), and then uses learning-based methods to correct the head pose, including a GravityNet to estimate the additional rotation and a HeadNet with optical flow features as input to estimate the scaling factor to the trajectory. The full body pose is generated with a modified version of DDPM (Ho et al., 2020) that is conditioned on head pose and trained on AMASS (Mahmood et al., 2019). We show the evaluation of the conditional diffusion part here.
- **Location-based.** This baseline is inspired by state-of-the-art methods that use transformer-based models for body pose estimation from sparse inputs (Castillo et al., 2023, Jiang et al., 2022). We adapt these methods to utilize 3D positions as opposed to the traditional parametric body model. During the training phase, the model was subjected to 40,000 iterations, using the Adam optimizer with a learning rate of  $1e^{-4}$ . The window size for temporal analysis was set at 40 frames, and we minimized the Mean Squared Error (MSE) loss between predicted poses and ground truths. As for the input, our model receives a sequence of head poses captured by the device.

We implemented and/or trained four baseline models for the *ego hand pose estimation*.

To estimate the 3D hand joint from monocular 2D ego view images, 2D heatmaps can be explicitly estimated and lifted to 3D space, or 3D joints can be directly estimated from extracted 2D features. The feature extractor backbone could be CNN-based or transformer-based. The proposed baseline methods cover different choices of model designs. All the baseline models work on single frame images without temporal information. The baseline models are trained on manual or manual+automatic annotations, and are only evaluated on manual annotations.

Notably, most baseline methods generate hand mesh as final results in their original paper. We modified them to be trained and supervised only on 2D/3D hand joints (not on hand mesh) to fit the benchmark.

- **THOR-net.** THOR-net (Aboukhadra et al., 2023) uses Keypoint-RCNN as the feature extractor to obtain 2D information and derive 2D hand keypoints heatmaps explicitly. The method then lifts 2D estimates to the 3D space using GraFormer (Zhao et al., 2022), which is a model consisting of Graph Convolutional layers and Attention layers. We use only the 2D-to-3D pose GraFormer branch in THOR-net to adapt the method to our task. The training takes around 4 hours for the manual dataset on a GeForce RTX 4090 Graphics Card, and around 10 hours for the dataset combining manual and automatic annotations.
- **HandOccNet.** HandOccNet (Park et al., 2022) uses a ResNet50(He et al., 2016)-based FPN (Lin et al., 2017) to extract 2D features. The method then uses two Transformer-based modules: Feature Injecting Transformer (FIT) to inject hand information into occluded region, and Self-Enhancing Transformer (SET) to further refine the 2D features. The method proposes a regressor based architecture to produce 2D keypoints, MANO (Romero et al., 2017) pose, and MANO shape parameters to predict joints and vertices. To accommodate our baseline, only 2D keypoints and 3D joints location losses are used in the training phase. The training takes around 2 hours for the manual dataset on 8 NVIDIA V100 Graphics Cards.

- **POTTER.** POTTER (Zheng et al., 2023) proposes Pooling Attention Transformer (PAT) to extract 2D visual features, which significantly reduces the memory and computational cost without sacrificing performances. The method then applies a mesh regression head HybrIK (Li et al., 2021) to generate 3D joint and mesh results. The training takes around 43 minutes for manual dataset, and around 4 hours for manual+auto dataset on a GeForce RTX 4090 Graphics Card.
- **METRO.** METRO (Lin et al., 2021) extracts a CNN-based global image features. The method then uses a transformer encoder to jointly model vertex-vertex and vertex-joint interactions, and outputs 3D joint coordinates and mesh vertices simultaneously. Since the training of METRO strongly depends on hand mesh supervision, which is not present in the annotations, we borrowed the checkpoint trained on FreiHand (Kolotouros et al., 2019) dataset and run the inference only, without training it on our benchmark.

## Results

Table 12 shows the evaluation results of all the baseline approaches for the *body pose estimation* task. First, note that the static pose baseline obtains a significantly higher MPJPE than all the other approaches. This finding suggests that the poses across different scenarios in the dataset are extremely diverse. Thus, attempting to have the same static pose for all test cases is unfeasible. In contrast, the proposed baseline implementations achieve notable enhancements in performance. Table 13 shows the performance of each method per scenario. While these developments are promising, we believe that further refinement is possible, especially in lower body pose estimation and to ensure temporal consistency in predictions.

We report the MPJPE and PA-MPJPE of the baseline models for the *body pose estimation* task in Table 14, and their corresponding parameter numbers and multiply-accumulate operations (MACs) in Table 15a. We further analyze the error distribution across different hand joints. Figure 26 shows that the thumb finger and finger tips tend to have larger errors, most likely because they are occluded or invisible more often.

Method	Validation		Test	
	MPJPE	MPJVE	MPJPE	MPJVE
Static pose	254.29	-	215.87	-
EgoEgo	24.53	0.78	26.38	0.66
Kinpoly	21.66	0.86	24.36	0.65
Location-based	20.73	0.74	18.51	0.64

**Table 12: Results for the 3D human body pose benchmark.** We report the Mean Per Joint Position Error in cm and the Mean Per Joint Velocity Error in m/s for all the baseline approaches.

Scenario	EgoEgo	Kinpoly	Location-based
Basketball	21.36	24.98	19.89
Soccer	23.08	19.09	16.62
Bike repair	30.18	25.19	20.61
Cooking	23.71	20.80	12.65
Health	32.57	29.23	11.63
Dance	20.93	18.03	21.15
Music	33.81	30.30	15.00

**Table 13: Body pose estimation Test results per scenario.** We report the Mean Per Joint Position Error in cm.

	Manual		Manual+Auto	
	MPJPE	PA-MPJPE	MPJPE	PA-MPJPE
METRO*	-	20.61	-	20.61
THOR-net	51.24	17.99	47.64	17.61
HandOccNet	-	17.22	-	13.56
POTTER	30.57	11.14	28.94	11.07

**Table 14:** MPJPE and PA-MPJPE in mm for ego hand pose baseline models. \* denotes methods *not* trained on the benchmark.

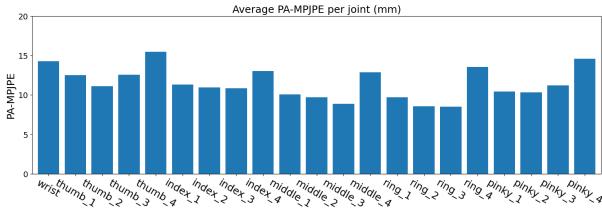
	THOR-net (Aboukhadra et al., 2023)	HandOccNet (Park et al., 2022)	POTTER (Zheng et al., 2023)
Params (M)	59.5	37.22	14.5
MACs (G)	123.6	15.5	5.2

(a) Number of parameters and MACs for the different ego hand pose baselines. [consider joining up the smaller tables into a singl figure with subparts for flow.](#)

# visible views	3	4	5	6
PA-MPJPE (mm)	14.01	12.15	11.03	10.02

(b) PA-MPJPE for joints that are visible in different number of views (including ego and exo views). Results generated from POTTER (Zheng et al., 2023) evaluation.

**Table 15:** Analysis for the hand pose benchmark.



**Fig. 26:** Average PA-MPJPE for each joint. Results generated from POTTER (Zheng et al., 2023) evaluation.

For each annotated joint, the manual annotations keep record of the number of views where the joint is visible. The visible 2D observation is then used for triangulation in 3D ground truth generation. This can be taken as an indicator of the uncertainty of the ground truth, and the difficulty level for the estimation of the joint (usually, a joint visible by fewer views indicates that it is more entangled with objects or other part of the hand). Table 15b shows that the PA-MPJPE decreases as the visible number of views increases. To guarantee the ground truth accuracy, all experiments are performed only on joints at least visible in 3 cameras.

## 6 Conclusions

Ego-Exo4D provides a robust data collection pipeline, a dataset of unprecedented scale and realism, and a benchmark suite for ego-exo skilled activity understanding and video learning. We propose a replicable pipeline to collect synchronized multi-exo cameras along with egocentric data in diverse settings indoors and outdoors. The setup was replicated across country boundaries to collect a homogeneous dataset. This offers, for the first time, collection of ego-exo data outside mocap suits or lab settings, capturing the participants where they naturally carry out their skilled activities—e.g. chefs in their kitchens, dancers in their studios, and football players on the pitch.

Eight compelling domains were selected with diverse skilled activities. We divide these domains into physical skills—those that particularly require strengthening, flexing or training the human body to carry out a skill, e.g. dancing; and procedural skills—those that one masters through the usage of tools to manipulate the surrounding environment, e.g. cooking. Both types of

skilled activities (physical and procedural) have never been explored jointly. By bringing the two types of skilled activities to a common dataset and benchmark, Ego-Exo4D will address the everlasting promise of assistive technologies, beyond a single domain or application.

The dataset comes with a suite of benchmarks, models, evaluation scripts, web-based visualizer, and baselines, to assist the research community in exploring and building on the challenges posed by Ego-Exo4D.

One of the challenges of Ego-Exo4D, and consequently its limitations, is the difficulty in optimizing the positions of exo cameras in the various settings. Often, the action is occluded by the person in the majority of the exo camera due to the standard and static positioning of these cameras. This impacted the annotations at times, and we opt to manually select suitable exo cameras during annotations. Additionally, the data is long-tailed due to the natural durations of activities. For example we have 9x more hours of cooking than soccer, as preparing a meal takes much longer than a soccer drill. The tasks also differ in their skill challenge; for example, learning to shoot the basketball into a hoop requires a lot more training and expertise than learning to carry out a COVID test. This diversity, while part of daily activities, could introduce challenges to current model training.

In addition to the three foundations of this dataset: pipeline, dataset, and accompanying benchmarks, two additional research gems should be highlighted for interesting future directions.

First, the videos are accompanied by three levels of linguistic descriptions: (i) fine-grained narrations of ongoing actions, (ii) descriptions of the activity by the participants themselves reflecting on their expertise and why they perform in a certain manner—we refer to this as ‘act and narrate’, as well as (iii) commentary from expert tutors. These are people trained to teach or evaluate the skill of others. The commentary is temporally synced with the action and enriched by spatial highlights to pay attention to particular ways in performing skills, both showcasing excellence as well as points for improvement—we refer to this as ‘expert commentary’. Ego-Exo4D thus offers the first resource of its kind to compare how actors and observers reflect, similarly or distinctly, on their skills.

Second, Ego-Exo4D offers for the first time the chance to study detailed hand-pose including hand-object interactions and full body pose in one dataset. A few recent works have showcased the potential of combining both in synthetic (Tendulkar et al., 2023) or controlled (Taheri et al., 2020) settings. Ego-Exo4D can be used as a base to take these directions further into real-world recordings.

Ego-Exo4D is a massive step towards a holistic understanding of the individual camera wearer—particularly their personal goals to advance their skill levels for their job or hobby. Models that understand one’s skill will offer the ultimate assistive companion, as noted by the survey paper (Plizzari et al., 2024). Such a companion can offer actionable personalised feedback (Ashutosh et al., 2024), then continuously monitor and quantify that feedback’s impact on the camera wearer’s skill over time, towards life-long skill mastering. Beyond skill understanding in video, Ego-Exo4D serves as a resource to deepen research in general 3D vision (environment reconstruction, camera relocalization, and others), video-language learning (grounding actions and objects, multimodal representation learning, language generation), and traditional exocentric activity understanding.

## Contribution statement

This project is the result of a large collaboration between many institutions over the last two years. Initial authors represent the leadership team of the project. Kristen Grauman initiated the project, served as the technical lead, initiated the recognition and proficiency benchmarks and expert commentary, and coordinated their working groups. Andrew Westbury served as the program manager and operations lead for all aspects of the project. Lorenzo Torresani led development of the capture domains, initiated the relation and ego-pose benchmarks, and coordinated their working groups. Kris Kitani led development of the multi-camera rig and supported the Ego-Exo4D engineering team on all aspects of the data annotation and organization. Jitendra Malik served as a scientific advisor. Authors with stars (\*) were key drivers of implementation, collection, and/or annotation development throughout the project. Authors with daggers (†) are faculty and senior researcher PIs

for the project. The Appendices detail the contributions of individual authors for the various benchmarks, data collection, and annotation pipelines.

## Acknowledgements

We gratefully acknowledge the following colleagues for valuable discussions and support of our project: Vittorio Caggiano, Sarah Carroll, Ilé Danza, Ahmad Darkhalil, Zona de Bloque, Alex Dinh, Rene Martinez Doehner, Ivan Cruz, Matt Feiszli, Vance Feutz, Kelly Forbes, Rohit Girdhar, Pierre Gleize, Andrés Hernández, Shun Iwase, Hanxiao Jiang, Armin Kavian, Bolin Lai, Vivian Lee, Brighid Meredith, Ashley Massie, Natalia Neverova, Manohar Paluri, Joelle Pineau, Artsiom Sanakoyeu, Paresh Shenoy, Jiaray Shi, Jiasheng Shi, Gaurav Shrivastava, Mitesh Singh, Manasi Swaminathan, Arjang Talattof, Ali Thabet, Laurens van der Maaten, Andrea Vedaldi, and Tobby Zhu.

We also sincerely thank the 52 experts who contributed to the expert commentary for their expertise and support; they are listed individually in Appendix D. Thank you to the Common Visual Data Foundation (CVDF) for hosting the Ego-Exo4D dataset. Finally, thank you to the 740 participants who contributed to this dataset and shared their skills in video.

UT Austin is supported in part by the IFML NSF AI Institute. University of Catania is supported in part by the project Future Artificial Intelligence Research (FAIR) – PNRR MUR Cod. PE0000013 - CUP: E63C22001940006. Luigi Seminara is supported by PNRR PhD scholarship “Digital Innovation: Models, Systems and Applications” DM 118/2023. Simon Fraser University is supported in part by the Canada Research Chairs Program (CRC-2019-00298) and NSERC Discovery (2019-06489) Georgia Tech is supported in part by NSF (#2144194) and NSF-GRFP Indiana University is supported in part by NSF DRL-2112635 (AI Institute for Engaged Learning). Univ. of Bristol is supported in part by EPSRC UMPIRE (EP/T004991/1) and EPSRC PG Visual AI (EP/T028572/1). Z. Zhu is supported by UoB-CSC Scholarship. University of Tokyo is supported in part by JST ASPIRE Grant Number JPMJAP2303 and JSPS KAKENHI Grant Numbers JP24K02956.

## A Camera setup and recording details

### A.1 Time sync

To sync cameras, we employ a pre-rendered sequence of QR Codes (*i.e.*, QR code video) that encode a wall-clock time. We show this QR code video using the smartphone at 29fps to all cameras in sequence and exploit the difference in frame rates to finely sync the cameras. In theory, the QR code decoded on a frame that captures a QR change is likely the one that was visible during that frame’s center of exposure. With a single QR, the camera’s center of exposure time could be anywhere within the 34.48ms that the QR is shown. However, with two consecutive frames with the same QRs, we can localize that time down to  $\pm 0.574\text{ms}$ . The same approach yields  $\pm 0.558\text{ms}$  for the 59fps GoPros given 3 consecutive frames (see Figure A1), providing sub-frame synchronization accuracy.

We manually verified that each GoPro camera was within 1 frame ( $\pm 16.66\text{ms}$ ) of the Aria RGB camera by visually comparing them at single-frame moments (*e.g.*, contact frames) using a synced video collage at the start and end of each capture. We checked points near the start and end of each capture under the logic that sync is a linear mapping and camera clock speed is mostly constant, so if the error is  $\pm 1$  frame at the start and  $\pm 1$  frame at the end, it will be  $\pm 1$  frame throughout.

An ‘audio sync’ fingerprint was played at the start and end of each capture to synchronize audio streams but has not been used.

#### Challenges and workarounds

In practice,  $\sim 70\%$  of recorded captures yielded frame-accurate sync through our automated pipeline. Inaccurate sync causes included observed issues (*e.g.*, phone changing orientation mid-playback, video playback interruptions) and suspected ones (*e.g.*, videos not playing back at precisely 29fps, center exposure times not being evenly spaced). To recover these captures, we employed a manual sync procedure wherein people manually selected frame timestamps that should be aligned based on precisely time-localizable events, *e.g.*, a lighter first sparking, a soccer ball



**Fig. A1:** With a QR code timer playing back at exactly 29fps, cameras with evenly spaced center-exposures can be precisely time-localized to the QR timer with these multi-QR patterns.

making contact with a cleat, or a hand beginning a fast slide down the neck of a guitar. This unblocked the remaining  $\sim 30\%$  of captures at the cost of less accurate sync.

#### Alternatives

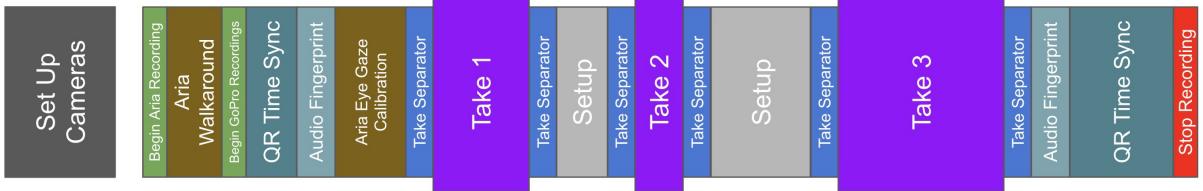
We explored and disqualified other sync options—notably using Timecode with TentacleSync or Ultrasync. Both of these solutions use LTC to encode a 1fps timestamp into the audio channel of a connected device. Using them with GoPros would cost us the stereo-audio modality, which we opted to keep to support audio-based research areas. We additionally lacked an ergonomic input solution for Aria to use while recording, so that mandated non-intrusive sync solutions.

### A.2 Take separation

To amortize the setup and tear down time required for each recording, we record multiple ‘takes’ (*i.e.*, one instance of a certain task) back-to-back and use a ‘Take Separator’ QR code (different from the time sync QR code video) that is identified in post-processing to auto-separate each take. This enables us to scale up recording—particularly for the physical scenarios where a single take can be less than a minute long. Data collectors track metadata for each take, identifying them by index and marking data such as participant ID (anonymous unique identifier), task (*e.g.*, making tea, making cucumber salad, performing CPR), and whether the take should be dropped (*i.e.*, if it is just setup time between activity enactments).

### A.3 Recording procedure

Our rig setup procedure entails setting up the stationary exo cameras in the recording environment and displaying QR codes to perform time sync and



**Fig. A2:** Overview of the recording procedure

then take separations. Figure A2 overviews our recording procedure.

1. Position tripods, power on GoPros, and set camera angles to ensure maximum human coverage.
2. Begin Aria recording via smartphone.
3. Conduct a walk-around with the Aria glasses to build a basemap for 3D reconstruction and camera localization. Match the viewpoint of each GoPro camera by positioning the Aria directly in front of its lens.
4. Start QR Timesync Video off-screen. Show QR video to Aria RGB camera.
5. Use GoPro Remote to begin GoPro recording. Show QR video to each GoPro camera. Play Audio Sync fingerprint from the center of the space.
6. Pass Aria glasses to (new) participant. Perform Eye Gaze Calibration via the Aria app. Show ‘Take Separator QR’ to one GoPro and begin the take. Show ‘Take Separator QR’ to one GoPro after the take is complete and repeat this step for each participant/take. Do not repeat gaze calibration if the participant has not changed.
7. Play Audio fingerprint from the center of the space. Restart the QR Timesync Video off-screen. Show it to the Aria RGB camera, then each GoPro. Stop recording on all cameras.

The core camera rig was extended to handle onsite requirements and regional challenges. The team at Universidad de los Andes introduced a top-down (ceiling mounted) GoPro for dance, which was adopted by the team at the University of Pennsylvania with an overhead pole mount. The teams at University of Pennsylvania, IIT-Hyderabad, and Indiana University added an additional egocentric, head-mounted GoPro.

#### A.4 Aria post-processing

First, the Aria Machine Perception Services (MPS) pipeline is invoked for *each full Aria recording*—these typically are about 20 minutes to 1 hour long and can include several takes, the hand-over in-between takes, as well as some other set-up steps. This is followed by localizing all GoPro videos of that scene as described above, and finally followed by time-synchronization across Aria and the GoPro cameras, as well as take-separation, as described below.

There are total of 783 Aria recordings processed by MPS—containing the total 5,035 takes in the dataset. 95.9% of these recordings have successful Aria localization throughout the whole recording, with only 3.5% containing a partial tracking failure (leading to short gaps in the 6DoF trajectory). Three (0.6%) recordings failed completely. The most common failure reason is physical shock on the glasses, for example when the glasses are accidentally dropped on the ground or the table.

Furthermore we attempted to localize a total of 3,724 GoPro recordings, 91.4% of which are successfully localized. Similar to the Aria recordings, GoPro’s are localized on a *recording* level rather than on a *take* level. This helps in particular with very short takes as are common during physical activities—as there otherwise would not be sufficient visual overlap across Aria and GoPro perspectives. The most dominant reason for GoPro localization failure occurs when the GoPro is pointed to an texture-less area (e.g. a white table) which lacks the necessary visual features to perform localization. As the GoPro’s are static, this cannot be compensated for by device motion as is the case for the moving Aria device.

Technical documentation and open-source tooling for Aria recordings and MPS output is

available on Github<sup>9</sup> and the associated documentation page<sup>10</sup>. It includes both python and C++ tools to convert, load, and visualize data; as well as sample code for common machine perception and 3D computer vision tasks.

## B Data collection

Twelve research labs came together for nearly two years to create Ego-Exo4D. Importantly, our collection across the sites was a coordinated effort, with common guidelines, scenarios, and camera rigs. In this way, the dataset is cohesive at the same time it is diverse. In this section we describe the data collection details that are specific to each partner site, e.g., how they recruited participants, which of the 8 scenarios they captured, or any modalities they added on top of the common rig.

Figure 9 shows the breakdown of which scenarios were captured by each partner institution as well as a map highlighting the locations of the 12 labs involved in data collection. Note that an additional four institutions not shown on the map are part of the consortium (e.g., contributing to benchmarks) but did not collect data. They are UT Austin (USA), KAUST (Saudi Arabia), University of Catania (Italy), and University of Bristol (UK).

### B.1 Carnegie Mellon University

Carnegie Mellon University focused on three skill-based activity scenarios: (1) soccer, (2) bike-repairs, (3) cooking. The exocentric cameras for our collections, four in total, were arranged approximately in a square configuration at a consistent height to capture the full range of the activity. Notably, for the soccer activities, an additional exocentric viewpoint was positioned inside the goal post to offer a more comprehensive perspective on the participants.

**Soccer** In the soccer scenario, we collaborated with professional players from the Pittsburgh Riverhounds team, representing the experts, and students from Carnegie Mellon University (CMU) as the beginners. We captured the soccer scenario across 4 different locations. The drills featured a variety of movements such as dribbling, goal kicks,



(a) Cooking Scene



(b) COVID Test Scene



(c) Bike Repair Scene

**Fig. B3:** Views from two different cameras for each scenario collected in Atlanta, GA, USA.

and juggling, with each participant performing for a minimum of 3 minutes. This scenario resulted in roughly 4 hours of egocentric footage and 18 hours from exocentric perspectives, encompassing 32 participants in total.

**Bike repair** In the bike-repair segment, our experts were seasoned mechanics with over a decade of experience from Allegheny county. To ensure authenticity, we visited each mechanic in their respective shops to allow usage of their own tools and setup. Four tasks were captured for each bicycle, and we ensured bicycle diversity by selecting different sizes, shapes, colors, and makes. The tasks include tire removal, tube change/ inflation, tire reassembly, and clean/lube chain. This yielded 3 hours of egocentric recordings and 12 hours of exocentric footage, encompassing 22 different bicycles.

**Cooking** For the cooking section, we documented a professional chef in his traditional kitchen environment. Our dish of choice was scrambled eggs, and to inject variety, the chef prepared it using different techniques. This segment summed up to an hour of egocentric recordings and 4 hours from exocentric viewpoints.

<sup>9</sup>[https://github.com/facebookresearch/projectaria\\_tools](https://github.com/facebookresearch/projectaria_tools)

<sup>10</sup>[https://facebookresearch.github.io/projectaria\\_tools/docs/intro](https://facebookresearch.github.io/projectaria_tools/docs/intro)

All recordings were conducted in Pittsburgh, PA, USA, strictly adhering to CMU’s Institutional Review Board (IRB) guidelines. Every participant was briefed about the recording process, and prior to their involvement, a signed consent form was obtained.

## B.2 FAIR, Meta

We collected 119 total takes of skills demonstrations in New York and three different locations in California. We focused on cooking and bike repair, bringing in a skilled workforce of chefs and bike technicians that serve major kitchens and repair shops in the area. We used the unified camera rig of 1 Aria and 4 GoPros without any additional sensors.

**Bike Repair** Our skilled mechanics performed four different bike repairs for a total of 102 takes. We focused specifically on wheel repairs (removing and installing the wheel & flat repairs). While we strive for diversity in terms of the model of bikes, a majority of those in the dataset are drawn from standard fleet bike models, which contain identical parts and components. The location featured in the dataset is a well-equipped, industrial scale bike shop.

**Cooking** Our chefs recorded five different recipes as part of 17 unique takes, including salads, egg dishes, and Asian garlic noodles. Locations featured in the dataset are three different professional kitchens used to prepare and serve hundreds of people each day.

Internal documentation and processes ensured all participants provided informed consent to appear in the dataset and participation was strictly voluntary.

In total, we were able to mobilize five chefs and four bike mechanics. Participating chefs and bike technicians are highly skilled, with all research subjects reporting that they do the activity shown in the dataset daily or weekly. Similarly, eight research subjects have more than 10 years professional experience.

## B.3 Georgia Tech

Collection at Georgia Tech focused on the Health, Cooking, and Bike Repair scenarios. Across these 3 scenarios, 279 takes were captured with 34 unique participants. For all scenarios, the unified camera rig was positioned such that 2 exocentric

cameras would ensure capture of the participant’s hands, and the other 2 exocentric cameras would capture the participant’s full body and the full environment.

Participants were recruited from different sources including flyers, campus organizations, email lists, and word of mouth. Five of these participants completed data collections for 2 scenarios (4 participated in Health and Cooking, and 1 participated in Cooking and Bike Repair). Potential participants were provided with the study description and consent form prior to scheduling a recording session. At the beginning of each session, study personnel walked through the consent form with the participant, and answered any questions. The participant then reviewed and signed the consent form to confirm participation in the study.

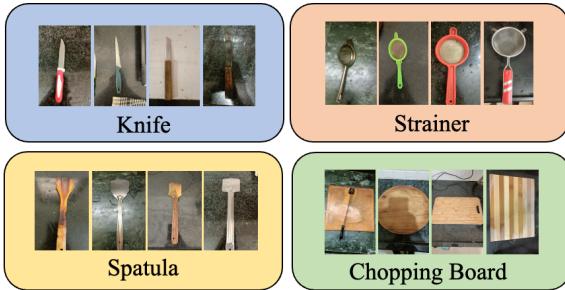
The recording environment differed by scenario and included participants’ homes, campus meeting rooms, and an on-campus bike shop. Fig B3 shows a sample environment and camera setup for each of the Health, Cooking, and Bike Repair scenarios. Further details of the data collection specific to each scenario is provided below.

### *Health*

Participants for the Health Scenario took COVID rapid test kits while seated at a table. Recordings were captured in 2 different on-campus meeting rooms. Participants were recruited through campus email lists and flyers in local coffee shops. Each recording session lasted approximately 40-60 minutes and consisted of a participant completing 5-7 test kits, using 2-4 different types of test kits. 7 different types of COVID test kits were used across the full collection. In total, 96 takes were recorded from 16 unique participants.

### *Cooking*

Participants for the Cooking Scenario prepared dishes from three recipes: Asian Salad, Tomato & Eggs, and Garlic Noodles in their home kitchens. Participants were recruited via mailing lists of local apartment complexes, contacting participants from prior research studies, and word of mouth. Each recording session lasted 2-3 hours, capturing 3-6 takes of a recipe being cooked from start to finish. Participants cooked 2-3 of the recipes during their session, depending on



**Fig. B4:** In Hyderabad, India, cooking was captured in different kitchens with socio-economic diversity. We observe that the same kitchen tools appeared in different shapes.

dietary restrictions and preferences. Participants were provided with ingredients and a paper copy of the recipe, and used equipment from their own kitchen to prepare the food. In total, 71 takes from 15 unique participants were captured. The takes were about evenly distributed among the three recipes. Recordings were completed in 10 unique kitchen environments.

#### Bike Repair

Participants for the Bike Repair Scenario performed repairs including taking off a wheel, putting on a wheel, replacing a tube, and cleaning a dirty chain. We recruited skilled participants from a campus bike repair organization. There were 8 unique participants, who each completed 1-3 recording sessions. Each session lasted 40-60 minutes and captured 5-7 takes of individual bike repairs. In total, 112 takes were recorded, showing the distribution across repair tasks. One session of 6 takes was recorded in a participant's home, while the rest were recorded in the campus organization's bike shop space, which is shown in Fig B3c. Due to the organization's access to a large quantity of used bicycles, there is large diversity in the make and model of bicycles across takes.

The study protocol was reviewed and approved by our Institutional Review Board (IRB).

#### B.4 IIIT-Hyderabad

In Hyderabad, we contributed to three scenarios - (a) cooking, (b) soccer, and (c) music. We formulated a data collection strategy tailored to the specific scenario, as outlined below.

Our primary objective was to comprehensively capture body and hand movements, along with their interactions with objects, during the execution of the activities. In general, we adhered to the standard camera setup instructions. Nonetheless, we incorporated an extra exo-camera for capturing soccer activities in order to enhance the overall coverage of the event. Additionally, for music activities, we introduced a head-mounted Go-Pro camera.

This decision stemmed from the observation that expert musicians frequently do not directly look toward their instruments while playing. Consequently, the head-mounted camera guarantees continuous visibility of both the hands and the musical instruments, providing an ego-view perspective.

The collection in India was done during the peak summer, and this led us to a challenging situation where the cameras frequently shut down due to overheating. To address this, we mostly avoided capturing multiple takes in one capture and placed the cameras into an ice chest box in between the captures to cool them down.

**Cooking** For cooking, we reached out to people located in Hyderabad with varying socio-economic backgrounds and explained the data collection plan, and goals. We also requested them to engage their family members as well as friends in this data capturing process. Finally, we recorded the videos with the 41 informed participants capturing in a diverse set of 19 kitchens, geographically well-apart in and around Hyderabad, resulting in a rich dictionary of kitchen utensils (see Figure B4), narrations in four different languages. Additionally, we made an effort to ensure a balanced representation of genders in our overall data collection process.

**Soccer** For soccer, we reached out to three different soccer training schools in Hyderabad with the overall recording plan and process. They helped us in recruiting local soccer teams who play professional tournaments and practice almost everyday. We also recruited few players from our university soccer teams. In total, we recorded 49 participants, 'performing dribbling, juggling, and penalty-kick activities.

**Music** For the music scenario, we contacted one music school and recruited 4 musicians from them having at least 3 years of experience of playing either the piano, guitar, or both instruments.

To add diversity, the musicians were asked to play western as well as Indian pieces.

Our collection protocol was reviewed and approved by our university's Institutional Review Board (IRB). The primary conditions set forth by the IRB encompass the following aspects: (a) participants with 18+ age are deemed suitable for inclusion in the project, (b) participants have provided explicit consent for their facial and vocal presence to be featured in the released videos, (c) participants have willingly agreed to take part without receiving any immediate financial incentives from the videos, and (d) participants have the autonomy to engage in the activities in an environment of their choice. The participants were given the detailed descriptions of the project beforehand and requested to sign the consent form. Each participant received compensation as part of the process.

We selected participants from a wide range of age groups, spanning from 18 to 61 years old, to introduce an additional layer of diversity. Moreover, the participants were from diverse professional backgrounds (e.g., coach, software engineer, data annotators, project managers etc.). Before sharing, we carefully examined each video to ensure there was no sensitive content.

## B.5 Indiana University

We focused on cooking, bicycle maintenance, and music scenarios. All activities were collected using the unified camera rig, including additional sensors in specific scenarios. For cooking, the 4 GoPros were placed 90 degrees apart from each other, with 2 placed close to the participants to capture hands and objects and 2 placed further to capture the overall scene. In music, 4 GoPros were placed in front of the player, approximately 45 degrees apart from each other. In addition, we attached an additional GoPro HERO10 camera to the participant's head (using a helmet), tilted down roughly 80 degrees to capture hand movements. In bike repair environments, the 4 GoPros were placed 90 degrees apart from each other, of which 1 GoPro was placed close to the bike, 1 GoPro was placed close to the workbench and tools, and 2 GoPros were placed further away from the participants to capture the overall scene.

**Cooking** For cooking, we had a total of 18 participants collect 72 takes and 20.5 hours of

video. For 15 of the participants, we used a commercial test kitchen at our university. We purchased all of the ingredients and kitchen equipment ahead of time and had them ready when each participant arrived. We asked them to make four dishes (chai tea, sesame-ginger salad, tomato and eggs, and noodles) and provided printed recipes for these dishes. The remaining three participants chose to record in home kitchens, and the four dishes they made varied based on their preferences (one participant made omelet, cucumber salad, noodles, and chai tea, another made scrambled eggs, sesame-ginger salad, sushi rolls, and brownies, and the third made scrambled eggs, cucumber salad, noodles, and milk tea). Due to concerns about food safety, we discarded (composted) the cooked dishes instead of allowing the participants to eat them.

**Music** For music, we had a total of 17 participants collect 60 takes and 6.5 hours of video. Participants were recruited based on their self-assessed proficiency in one of three instruments: piano, violin, or guitar. We recorded in 4 different locations including two studios, an office, and an auditorium that had a piano. Participants were instructed to play scales and arpeggios (2 mins), sheet music provided by us (3 mins), freeplaying (10 mins), and then recall and talk about any mistakes that were made during the playing and what could be improved (2 mins).

**Bike repair** For bike repair, a total of 13 participants recorded 108 takes and about 8 hours of video. We initially planned to hire professional bike technicians, but it was very difficult to recruit them in our relatively small city. Instead, we recruited more generally, looking for participants with (self-assessed) proficiency to do four basic bike maintenance tasks: removing a wheel, changing an inner tube, reinstalling a wheel, and cleaning and lubricating the chain. Most of the takes were recorded in a small house that is used for storage by our university's landscaping staff, and provided a realistic garage-like environment. We provided participants with a bike rack and supplies including bike tubes, pumps, tools, chain cleaner and lubricant, and gloves. To achieve diversity in different bikes and bike types, we asked participants to bring their own bike when possible, and we also provided 4 bikes (one of which belonged to one of the authors and the other

three which we bought at a salvage shop). Most participants performed takes on about 3 bikes. One participant chose to record in an apartment, and one recorded in a hallway in a university building instead of the garage due to scheduling conflicts.

Our protocol was reviewed and approved by our university’s Institutional Review Board. For each potential participant in each scenario, we first scheduled an online introduction meeting to tell them about the study and answer their questions and concerns. If they were interested, we agreed on the activity they would perform and when and where to meet for recording. We also sent them the informed consent form to give them sufficient time to review. On the recording day, we first asked them to sign the consent form, and then started recording their activities. All activities were recorded in an enclosed space to make sure that no one else accidentally entered the field of view of the cameras. We also ensured that the space did not have privacy-sensitive content, and we instructed participants not to use their phones or other devices that might show private content.

Within a few days, we securely sent the videos to the participant so that they could review the video and ensure that they were comfortable sharing it with others. They also completed a brief online demographic study, and then were sent an incentive payment in the form of an electronic Amazon.com gift card. We made clear to participants that if they were not comfortable sharing their video, we would destroy it and they would still receive their incentive payment, although none of the participants chose this option. We gave the participants US\$20 in gift cards for each hour of their time spent recording (with a minimum of \$20, and partial hours rounded up to the nearest \$5). We gave an additional \$20 gift card to reimburse travel costs for those who came to our facilities to record (e.g. in our kitchen, bike repair shop, or on-campus studio or auditorium). For cooking and bike repair, we gave an additional \$20 gift card to participants who provided their own ingredients or bike maintenance supplies, to defray these expenses.

We recruited participants in the Bloomington, Indiana, USA area through online email advertisement, word of mouth, physical flyers, and posting on social media. We recruited participants who were 18 years of age or older, had self-assessed

expertise in the activities as described above and could perform the tasks without wearing prescription glasses (which could interfere with the Aria’s gaze tracking).

## B.6 National University Singapore

In Singapore we focus on the following scenarios: soccer, health-related activities including COVID-19 ART testing and Cardiopulmonary Resuscitation (CPR), and cooking. In total, our collected data encompasses around 26 hours of egocentric videos and 117 hours of exocentric videos. These videos spread across 327 takes. In general, we adhered to the standard camera placement guidelines; however, for each scenario, we fine-tuned the position of the exo cameras based on practical considerations. For instance, in a small kitchen for cooking, we positioned the camera on the table to broaden its field of view.

**Soccer** For soccer, we conduct recordings at a university sports field. Our participants were primarily sourced through referrals provided by skilled participants recruited through online calls for participants. Additionally, during outdoor recording sessions, we occasionally invited surrounding bystanders to participate.

**Health** For health activities, we recorded in vacant classrooms, meeting rooms, and outdoor fields. CPR sessions are captured either in a yoga classroom or in a quiet, empty outdoor field. For recruitment, we circulated online calls for participants and then, for skilled activities like CPR, we collaborated with experts to organize training courses. Participants would participate in these courses and were trained to be proficient and then conducted recording afterwards.

**Cooking** As for cooking, which requires a kitchen, we used the kitchen in our lab mates’ apartments and arrange other participants to go there.

Our data capture has been approved by our university’s Institutional Review Board (IRB). The main requirements include that participants: (1) agreed to take part in the study, (2) agreed to donate their speech, image, video, IMU, and 3D scan data for the purposes of this research, (3) agreed that their face, tattoos, and voice may appear in the data, (4) have the right to withdraw their recorded data at any time.

In Singapore, high temperatures often pose the challenge of camera overheating, particularly for GoPro cameras, which can lead to protective shutdowns and interrupt data collection. To mitigate this, we place small ice cubes wrapped in wet wipes on the GoPro cameras to help cool it down during recording. Furthermore, we attempted to schedule our participants' recordings in the evening or during an overcast day.

Our data pool comprises contributions from about 93 meticulously selected participants, ensuring a proficient completion of the recordings. Particularly in soccer, most participants have extensive experience and were members of their school or college soccer teams.

## B.7 Simon Fraser University

We captured three types of scenarios in a variety of environments: cooking, basketball, and COVID-19 testing. In total, 88 participants carried out activities in the three scenarios we collected in a total of 61 data capture sessions, resulting in 519 activity takes.

We used the unified camera rig and followed the general collection guidelines with a number of small adjustments to facilitate scenario-specific capture. In kitchen and health scenarios where the participant interacts with small objects in tabletop height settings, the placement of exocentric cameras was optimized in a “two near, two far” setup to provide for visibility of the small objects and hands while also capturing the overall human pose during the activity.

**Cooking** The cooking scenario was captured in a decentralized fashion by going to the participants’ own residences and asking them to cook in their kitchen. This allowed for diversity in the environment as well as in the participant during data capture. Our data capture sessions resulted in 112 cooking takes.

**Basketball** Collection for the basketball scenario was done in a “round robin” fashion to reduce player-to-player overhead. We targeted a spectrum of experience levels, for example going from university basketball team players who compete at the national level to more amateur basketball players who only have played basketball occasionally. We collected 355 takes of basketball activities.

**Health** Following the standard data collection guidelines for health activities, we gathered 52 takes of health activities.

We followed the institutional research board (IRB) process at our institution to acquire approval for the participant recruitment strategy, study setup, and participant consent acquisition forms. All participants consented to their data being collected and distributed for research purposes. Participants have the right to request that their data be withheld from inclusion in the dataset.

We recruited participants by word of mouth, reaching out to specific clubs and groups for some of the activities, and more generally through advertisement using university-affiliated communication channels.

## B.8 Universidad de los Andes

We collected around 40 hours of video spanning four distinct scenarios that encompassed three physical activities (basketball, bouldering, and dancing) and one procedural activity (cooking). Figure B5 shows examples of the diverse scenarios that we collected. In total, we collected 2062 takes across all the activities. We used the unified camera rig with additional activity-specific sensors as described below.

### *Bouldering*

We partnered with a local climbing gym, which serves as a teaching and competition center in Colombia. We used the gym as the recording location and recruited participants who practice or teach bouldering there. Our focus was to recruit participants with four different levels of expertise: beginner, intermediate, advanced, and professional climbers. We hired expert route setters to design 33 climbing routes. These routes varied from beginner (V1) to expert level (V7). For data collection, each participant attempted to complete seven routes, having 3 minutes to make as many attempts as possible for each route. The routes were selected considering the expertise level of each participant. We located four exo cameras to capture each take; two horizontal cameras were facing the climbing wall, and the other two vertical cameras were on each side of the wall. Thus, the four cameras captured a complete view

of the climbing wall and the participant's movements at every moment. We gathered 1251 takes for the bouldering scenario from 40 participants. We ensured ethnic, age, and expert-level diversity across the takes.

### **Dancing**

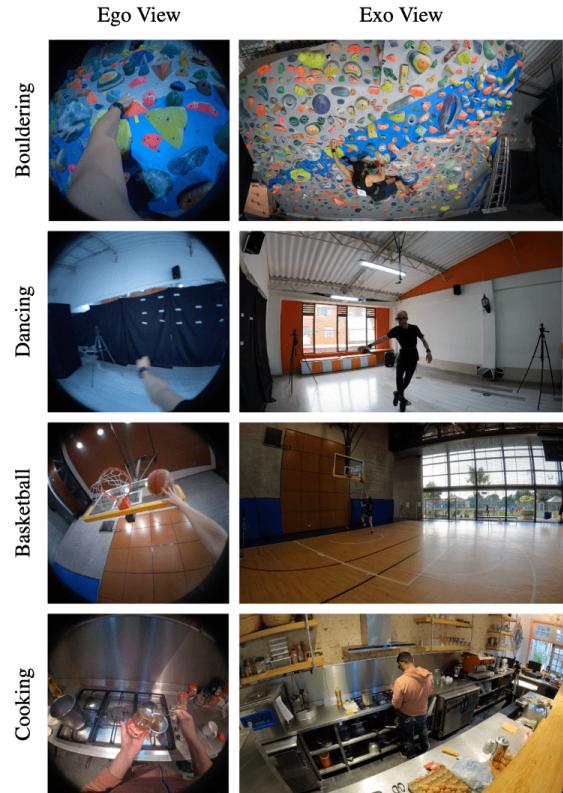
We collaborated with a salsa dance academy to use as a recording location and to help with participant recruitment. We recruited students from three expertise levels: beginners, intermediate, and advanced. According to the expertise level, each dancer performed different choreographies. Beginners recorded a single choreography, while intermediate and advanced participants recorded an additional one according to their expertise. Each attempt lasted one minute, and the dancer performed from six to ten attempts. The choreographies were designed by professional dancers who teach at the academy. We used five exo cameras: four forming a square, defining the dancing area, and the fifth camera placed on the ceiling. Given the salsa dance's velocity and the movements' complexity, this fifth exo camera gave a crucial point of view for further analysis. We gathered 600 takes from 40 participants across the three expertise levels.

### **Basketball**

We collected data from the professional women's team and students from a basketball class at our University. Each participant performed six to ten attempts for each basketball exercise. We collected all captures at the basketball court at our University's Sports Center. For this setup, we used four exo cameras around the basketball ring, ensuring a complete view of each exercise. For this scenario, we collected 167 takes from 38 participants.

### **Cooking**

We rented a professional kitchen equipped with all the necessary utensils to perform the captures. We focused on collecting data from two types of recipes: a dish with egg and a drink. Each participant could choose between cooking an omelet, scrambled eggs, tomato and eggs, and coffee latte or tea for the drink. Each participant was free to choose how to complete each recipe. Thus, our takes show diverse ways to prepare each recipe. For this setup, we placed four exo cameras around



**Fig. B5:** Egocentric and one exocentric view for each of the recorded scenarios in Bogota, Colombia.

the kitchen, all facing the user, to capture the whole kitchen without losing any detail of the person making the recipe. We placed two cameras on a counter facing the kitchen and the other two on each side of the kitchen. We collected 44 takes for the cooking scenario from 20 participants.

The Institutional Review Board (IRB) of our university reviewed and approved our study protocol. All participants signed a consent form before participating in the study.

We partnered with professional training centers for physical activities that helped us recruit volunteers with different expertise levels. These volunteers were previously familiarized with the activities and the environment where the captures occurred. In addition, we recruited family members, friends, and acquaintances of students and faculty members of our research group for cooking.

## B.9 University of Minnesota

Collection at the University of Minnesota focused on two main scenarios: Bouldering and Cooking. A total of 249 takes with 53 unique participants were collected. We collected all data using the unified camera rig with no additions.

### *Bouldering*

The bouldering activity was collected at a local bouldering gym, focusing on a wall with 14 different routes ranging in difficulty from beginner to expert. We collected 210 takes from 42 unique participants. Participants were asked to climb four to five routes of their choice, with the ability to take breaks within or between takes. Expert climbers who felt comfortable with the routes were able to narrate their approach and climb in real time. As participants were able to choose routes freely, our five exo-cameras were set up to accommodate the entire wall.

### *Cooking*

Cooking activity was collected on-site at each individual's home kitchen. Five exo cameras were set up in each kitchen to maximize coverage of both the participant and the environment. We captured 9 unique kitchen environments with 14 unique participants whose skill levels ranged from cooking novice to commercial chef. Participants focused on three recipes each (scrambled eggs, Greek salad, and pasta noodles from scratch), which were performed back-to-back on the day of recording.

Our data collection protocol was reviewed and approved by the Institutional Review Board at our university. At every take, the study personnel provides a guidance to a participant through the consent form prior to participation, ensuring the participant understands the purpose of the study and all risks involved, with each participant receiving payment proportional to their contribution.

Participants were recruited via word of mouth, campus organizations, and digital flyers which were distributed via local social media (Facebook) communities.

## B.10 University of North Carolina

Throughout our data collection at UNC, we focused on three skill-based activity scenarios: (1) basketball, (2) soccer, and (3) music drills. We used three unique environments (i.e., a basketball gym, a soccer field, and a music studio) to capture the data for each scenario. All recordings took place on the UNC campus. UNC's Institutional Review Board (IRB) reviewed and approved our study protocol. All participants signed a consent form before participating in the study.

To recruit participants, we used an online research study database, where participants from the local area could sign up to perform our study. We recruited participants willing to perform skill-based activities such as basketball, soccer, or music drills regardless of their skill level. Additionally, to recruit a more skilled group of participants, we contacted expert musicians from UNC's School of Music and athletes from UNC's basketball and soccer teams.

In total, we collected approximately 19 hours of egocentric and 76 hours of exocentric video data spanning approximately 548 takes of activity demonstrations from 56 participants (41 male, 15 female). Among the 56 participants, 44 were aged 18-25, 10 aged 25-50, and 2 aged 50-75. Furthermore, 26 participants had more than 10 years of experience in the scenario they chose to perform (e.g., basketball, soccer, music), 13 participants had 1-10 years of experience, and 17 had less than 1 year of experience. We used standard camera placement guidelines and the same recording devices described above.

### *Basketball*

All participants performed three basketball drills: Mikan Layup, Reverse Layup, and Mid-range Jumpshooting, for 388 takes. We recruited 11 expert players from the university team with 10+ years of experience. To improve the participant skill diversity in the dataset, we also recruited novice players with less than 1 year of playing experience. The location of data collection was a university basketball gym.

### *Music*

For the music scenario, we asked all 9 participants to play 5 minutes of scales and arpeggios and 10

minutes of free play for 27 takes. All of our participants were recruited from the university music club and considered themselves as experts at playing their respective instruments. The instruments featured in our collected dataset were piano, trombone, trumpet, and saxophone. The Ego-Exo4D music guidelines called for just piano, violin, and guitar, but we found it necessary to expand this list in order to gather data for this domain. All data was recorded in a university music room.

### *Soccer*

For soccer, we focused on three drills: dribbling, juggling, and penalty shots for 133 takes across 12 unique participants. 7 of these participants were experts with 10+ years of experience, whereas the remaining 5 participants were casual soccer players. All videos were collected at a university soccer field.

Our study protocol was approved by the Institutional Review Board (IRB). All participants signed a consent form before participating in the study.

### **B.11 University of Pennsylvania**

The University of Pennsylvania focused on capturing videos of experts of various levels playing musical instruments, dancing, and cooking. Over the spring and summer of 2023, UPenn captured 521 usable takes across 95 participants for the consortium’s collections with up to 7 views.

One primary goal of this project is to capture detailed body movement, especially hands, across ego view and exo views. We work to ensure highly engaged experts enjoy demonstrating their full skill capability.

The hand-object/instrument interaction region is the key to understanding human activities and evaluating their skills. Comprehensive hand pose information is especially important for the full analyses of scenarios collected at UPenn, especially the music scenario, where slight differences in finger motion result in entirely different performances.

We also observed that experts had a tendency to not need to look at their hands during play. Thus, we found the initial data capture using the general camera setup to lack crucial visual information in such scenarios due to:

(1) (in ego view) limited field of view of Aria

glasses, and skilled experts don’t need to look at their hands,

(2) (in exo views) frequent occlusion and self-occlusion caused by participants’ motion.

We added two cameras to maximize the view coverage.

**Head-mounted Camera:** The head-mounted camera on a helmet angled downwards to capture the hand/body region: (1) (ego) it follows the subject’s body motion faithfully, and (2) (exo) it is designed to focus on the hand-object interaction region with much less self-occlusion. Empirically, we found this additional camera is crucial for capturing guitar, violin, and cooking scenarios.

**Overhead Camera:** We replace the head-mounted camera with an overhead camera in (1) piano scenarios, where the overhead camera can have similar performance, and (2) dance scenarios, where the helmet can dramatically worsen the experience and performance of the participants.

We believe the goal is not to maximize the number of hours captured but to have the participants show (1) diverse techniques to build models for the scenarios, and (2) unique techniques to demonstrate their skill levels.

To get the most representative recordings of the participants, we aim to maximize their engagement during the data capture. Specifically, we (1) walk through the whole process with the participants before the data capture to familiarize them with the setup (2) let them choose their favourite music piece to play or dance with in music and dance scenarios; and (3) have a narrate and act section for the musicians to demonstrate how they feel about their performance.

### *Music*

For musical instrument playing, classified as a “physical” activity, we captured takes of musicians (1) warming up (scales and or Etudes) (2) sight-reading simple sheet akin to Suzuki Practice books or Etudes exercises, and (3) freplaying. We captured takes of violin, piano, and guitar, with a duet between a cello and a violin for one trial. Participants were recruited from a diverse pool of musicians, spanning the Penn Orchestra, local music schools, and independent music students. This pool’s experience ranged from professional instructors and performers to complete



**Fig. B6:** The Aria Glass ego view, head-mounted semi-ego view, overhead view and other static exo views in playing guitar in Philadelphia, PA, USA.

newbies. We totaled 275 takes over 37 participants. Notably, during the shooting, we observed that participants were particularly uncomfortable with the helmet used to mount the GoPro; it interfered with their head movements and the bow sometimes ended up knocking against the mounted GoPro. To combat this, we added additional cushioning to depending on the subject's head shape and broke sessions into chunks to allow for breaks.

### Cooking

Cooking, categorized as "procedural", consisted of preparing four dishes: an egg dish, a salad, a noodle dish, and a dessert. The group of participants consisted primarily of Penn students with experience ranging from amateurs to hobbyists. Professionals were unavailable due to scheduling conflicts. We totaled 81 takes over 20 participants. The entire filming process was undertaken within a three-week span, primarily at the apartment of one of the team's participants. This location expedited our data collection for this task by providing a stove and fridge for regular use.

### Dancing

Dance captures, classified as a "physical", consist of four takes of dancers performing dance routines to a song. The dance types recorded included Lindy-Hop Jazz, Bollywood, Latin, and Chinese Folk Dance; across these genres, we totaled 165 takes over 38 participants. The Lindy-Hop Jazz dancers came from the Jazz Swing Attacks, a dance club in Philadelphia. Contact was established via Instagram, and data, collected weekly over a month. This group contained a balanced mix of experienced instructors and beginner dancers. The Bollywood dancers, the Drexel's Philly Maza, were recorded in the Drexel Engineering Building. They compete nationally but routinely train beginner recruits. The Chinese

Folk Dancers were members of the local Great Wall Chinese School's dance club and independent student volunteers with prior competitive dance training. These were captured in the SIG Lab for collections.

All participants were confirmed to be at least 18 years of age by the time of participation and gave written consent for participating in these data collection trials. The consent form, in compliance with IRB guidelines, but gives participants the choice to back out. The collected information on basic demographic information should not be used to identify participants individually. All other data collected per participant (prior experience with task, average times spent per session, etc) could not be used to identify participants.

## B.12 University of Tokyo

In Tokyo, we collected video data for three scenarios: cooking and health for procedural activities and soccer for physical activities. We followed the standard camera configuration and calibration procedure of the Ego-Exo4D dataset for all scenarios. In the following paragraphs, we will describe the specifics of each scenario, particularly the unique aspects of our data gathering.

### Cooking

We recruited 12 Japanese participants living in the Tokyo area through a temporary employment agency. The gender and age of the participants were balanced to collect diverse behavior patterns. The participants cook three days or more each week in their daily lives. Each participant prepared three dishes: an omelet, a white radish & lettuce & tomato & cucumber salad, and a sushi roll. We recorded both versions with and without narrations for each dish and participant. A one-page summary of each recipe was provided before data collection and was shown during

video recording so that the participants could prepare the dishes smoothly, and the procedure of each recipe should be consistent between the participants.

All video recording of the cooking scenario was done in a rental kitchen studio equipped with an island kitchen and all the necessary kitchenware for four consecutive days. The studio is situated in a busy location in downtown Tokyo, and some external noises, like ambulance sirens, were audible during the recordings. We collected 68 takes from 12 participants. Of the 68 takes, 34 takes are with narrations, and 34 takes are without narrations. The length of each participant's making each dish twice with and without narrations is about 35 minutes, ranging from 24 min 27 sec (Sushi roll) to 55 min 4 sec (Salad). The length includes time for the camera synchronization procedure.

During the recording of the cooking scenario, we discovered a flickering issue in some of the video data due to the incompatibility between the Aria Glass sampling rate and the power frequency in Tokyo. To overcome this issue, we attempted to shoot as much as possible in daylight and adjusted the fps when using artificial lighting. While processing the videos, we discovered some exo videos had decoding errors due to damaged frames. Each corrupted video contained one to three damaged frames for an unknown reason. To address this, we re-encoded these videos by replacing the damaged frames preventing decoding with the nearest good frame. Note that some videos still contain damaged frames as long as those frames did not influence decoding. In addition, the original MP4 files recorded by GoPro contain 4 streams: video, audio, data0, and data1, but the re-encoded videos only contain a video stream and an audio stream.

### ***Health***

We recruited 17 Japanese participants living around Tokyo, Japan, through a temporary employment agency. The gender and age of the participants were balanced to collect diverse behavior patterns. We recorded videos of the 17 participants performing two tasks: COVID-19 rapid antigen testing using three test kits and performing CPR on a mannequin. We conducted all recordings in the same meeting room on campus

over two days. For the COVID test, an instruction manual of each test kit was provided to the participants before and during the recording. Also, we did not show the participants or record any COVID test results for privacy protection. For CPR, the participants took an introductory lifesaving course provided by the Tokyo Fire Department before recording. Besides, a one-page summary of the CPR procedure was provided before the data collection and shown during the video recording. This is so that the participants can perform CPR smoothly and the procedure is consistent among the participants. We collected 73 takes from 17 participants. All of the CPR procedures (17 takes) were recorded without narration. For the COVID test, we recorded the videos of the 17 participants using the three test kits (51 takes) without narrations. The video length of each participant's performing CPR is 9 min 48 sec on average, including the camera synchronization procedure. Similarly, the video length of each participant's using the three COVID test kits is 29 min 22 sec on average. Additionally, we recorded extra takes of 5 participants out of the 17 using a test kit with narrations (5 takes in total).

### ***Soccer***

We gathered videos of 14 Japanese participants, each performing three fundamental soccer drills: dribbling and juggling for two minutes each and penalty kicks ten times. Of the 14 participants, 13 are soccer players from a university football club. We recruited them through the staff of the club. The remaining one participant is not from the club but is an expert with over ten years of soccer experience. All the participants are male, and their age ranges from 18 to 30s. We recorded the videos on an outdoor soccer field at a local university over four days, with three to four participants participating each day. For juggling, we instructed the participants to include various movements such as juggling with thigh, inside and outside of feet, and alternating feet. For penalty kicks, we instructed them to shoot to the right side of the goal 5 times and to the left side five times. During penalty kicks, a helper aids the participant in retrieving the ball. This helper stands within the goal area and might be recorded by some cameras. We collected 42 takes from 14 participants. All the takes were recorded without narrations.

Our university's institutional review board reviewed and approved our study protocol. We explained the objective and the range of use of the videos through documents and took consent from each participant before the recording. In particular, we took the consent not to blur their faces to keep the naturalness of the videos.

## C Participants

We provide self-declared information on ethnic groups by the participants. Sharing this information was optional for all research subjects. Ethnicity is reported based on location specific categories as defined by the relevant partner lab. No such information was gathered from research subjects participating in our collections in California, New York, and Pittsburgh, Pennsylvania.

### *Atlanta, Georgia, USA*

100% of participants that reside in Fulton County, Georgia self-reported their ethnic group membership as follows:

Ethnicity	Number of participants
Asian	23
White	8
Hispanic/Latino	3

### *Bloomington, Indiana, USA*

100% of participants that reside in Monroe County, Indiana self-reported their ethnic group membership as follows:

Ethnicity	Number of participants
Asian	22
Black	1
Middle Eastern	1
White	18
Prefer not to say	1

### *Minneapolis, Minnesota, USA*

100% of participants that reside in Hennepin County, Minnesota self-reported their ethnic group membership as follows:

Ethnicity	Number of participants
White	41
Hispanic/Latino	4
Asian	8
Black	1

### *Tokyo, Japan*

100% of participants that reside in Tokyo self-reported their ethnic group membership as follows:

Ethnicity	Number of participants
Asian (Japanese)	45

### *Hyderabad, India*

100% of participants that reside in Hyderabad self-reported their ethnic group membership as follows:

Ethnicity	Number of participants
Asian (Indian)	95

### *Chapel Hill, North Carolina, USA*

100% of participants that reside in Orange County, North Carolina self-reported their ethnic group membership as follows:

Ethnicity	Number of participants
White	20
Indian	1
Asian	13
African American	9
Hispanic/Latino	3
Prefer not to say	3

### *Vancouver, British Columbia, Canada*

100% of participants that reside in Vancouver self-reported their ethnic group membership as follows. Please note that research subjects in this case opted not to use any assigned category and independently described their identity.

Ethnicity	Number of participants
African/Nigerian	4
Asian	9
White/Caucasian	10
Chinese	26
European	1
Iranian/Persian	14
Italian	1
Jamaican	2
Kazakh	1
Kyrgyz	2
Middle Eastern	1
Mixed	3
South Asian	2

### *Philadelphia, Pennsylvania, USA*

100% of participants that reside in Philadelphia self-reported their ethnic group membership as follows:

<b>Ethnicity</b>	<b>Number of participants</b>
White/Caucasian	10
Asian	30
African American	3
Hispanic/Latino	4
Prefer not to say	43

***Singapore***

100% of participants that reside in Singapore self-reported their ethnic group membership as follows:

<b>Ethnicity</b>	<b>Number of participants</b>
Chinese	65
Indian	3
Singaporean	2
Indian/Chinese	2
Prefer not to say	17

***Bogota, Colombia***

100% of participants that reside in Singapore self-reported their ethnic group membership as follows:

<b>Ethnicity</b>	<b>Number of participants</b>
Black/ Afro-descendant/ Afro-Colombian	7
Mixed	104
Palenquero	1
Raizal	1
White/Caucasian	38
Prefer not to say	23

Participant surveys were separated into two: a pre-task questionnaire and a post-task questionnaire. The pre-task questionnaire aims to capture the participant's perceived skill level whereas the post-task questionnaire captures the participant's reflection on how well the task went. The list of questions for both questionnaires can be found in Table C1 with questions/answers designed for consistency and ease of filling in, as participants would be filling these out before/after each recording. This involved using multiple choice and Yes/No answers with open text fields being utilized sparingly.

	Question	Answer Type
Pre-Task	Recording Location	multiple choice
	How many times do you estimate you have done this task?	multiple choice
	How often do you carry out this task?	multiple choice
	How many years have you been doing this task?	multiple choice
	Have you taught this activity to others before?	Yes/No
	Have you recorded a video of yourself carrying out or explaining this task before?	Yes/No
	Have you watched videos of others doing this task before?	Yes/No
	Do you have any qualifications/professional training that are related to the task?	Yes/No
	How long does it typically take you to complete this task?*	text
	How long would you typically spend in one practice session of this task?†	text
Post-Task	Self Reported Quality	multiple choice
	Completed Route?‡	Yes/No
	What mistakes/errors did you make during this task?	text
	Any issues with the familiarity of the tools/location?	text
	Did it take longer/shorter than your initial expectation and why?	text
	How did you find wearing the camera?	multiple choice
	How easy was the setup for recording?	multiple choice
	Any other comments to take on board?	text

**Table C1:** Questions for the pre-task and post-task questionnaires. \*: Only applicable for non-dance/non-music scenarios. †: Only applicable for Dance/Music scenarios. ‡: Bouldering scenario only.

## D Language Descriptions

As introduced in the main paper, Ego-Exo4D provides three forms of parallel text corpora for the video: expert commentary, narrate and act, and atomic action descriptions. Table D2 shows examples from different scenarios highlighting their distinctions in style and point of view. Figure D7 shows word clouds per scenario and annotation type highlighting the differences in vocabulary and word frequency.

In Figure D8 we further emphasize the characteristics of each text corpus across three axes: total vocabulary size, average number of captions per video, and caption length. See caption for details.

### D.1 Expert Commentary Tool

To collect expert commentary, we developed a web-based tool, which is open sourced as part of the Ego-Exo4D dataset and benchmark suite. See Figure D9. Known as the Narrator, this application supports video playback for Ego-Exo4D skills demonstrations, records time-stamped verbal commentary, and allows exporting and viewing commented videos. As a web-based platform, the Narrator can be simply accessed through a browser, with minimal set-up and less restrictive system requirements compared to tools requiring local installation. These attributes made it efficient to onboard and manage our geographically distributed experts. We acknowledge the EPIC Narrator (Damen et al., 2018) as the open-sourced inspiration and source code for this initiative.

### D.2 Atomic Action Description Statistics

See Table D3 for atomic action descriptions summary statistics.

Domain	Atomic Action Description	Narrate and Act	Example commentary
Cooking	C turns on heat on the gas burner.	So I'm going to start out by boiling some water.	Here the preparer is checking the pasta for done-ness. It's important to do this and not rely on what a package says. Use a package that gives you cooking time as a guideline and start to check your pasta, you know, a few minutes before the maximum amount of time given for cooking that specific pasta.
Health	C inserts the nasal swab in the buffer test tube on the covid test kit pack with his right hand.	Open the newly picked up tube, place the swab in the tube, stirring the swab in the tube.	So this individual has done a great job of making sure that her nasal passageways have adequate time in contact with her nasal swab. Something that might make it a little bit easier for her is if she could tilt her head back just a bit so that she wouldn't have to strain quite so much to get that access. Additionally, she did a great job making sure that the nasal swab was about an inch into her nose.
Bike repair	C holds the bike wheel with her left hand.	And then I will locate the location of the valve cap and pull the tube out of the wheel.	It's a great method to always double check or do a pre-check before beginning work on a bicycle to make sure the issue that you are working to fix is the only issue that is occurring. If not, you could find a secondary issue or something else that may be greater than the one you are currently working on.
Music	C puts the bow on the violin with his right hand.	So regardless of how tricky left hand passage work is you want to always keep your bow completely independent.	This is a really great use of the bow and decision to play in this middle third of the bow. This is exactly where they should be playing. And we can hear that the note envelope is very consistent and that it's very controlled and that it also allows the rhythm to be stable...
Basketball	C runs towards the hoop with the basketball.	Now I'm going to do a reverse layup, stepping right, left, going up with the right hand.	As the ball goes through the basket, she catches the ball and does an excellent job of keeping the ball high, never allowing the ball to drop down to her waist area, but keeping the ball high in her upper chest, neck area throughout the drill...
Bouldering	C places both hands on a red hand hold.	So I know that a lot of these holds, I'm going to need my weight leaning to the left to utilize	Once the climber recovered from the foot cut, the climber pasted the right foot on this foot jib and then did a toe match. So brought this foot in and then dropped the right foot down and to the right to again counterbalance so that the climber can then move their left hand out left. But at this point the climber is just a little too gassed to be able to make this move, which is unfortunate.
Soccer	C kicks the ball to the right with his right foot.		Angle approach, start position is good, maybe slightly squarer than 45, but again because the intended outcome from previous actions is into the left, by being a little squarer is going to help him be able to rotate his hips to move to the left, but on a slight angle is good and help him with his technical action.
Dance	C moves her right leg forward while swinging both hands.	Ring and wing, one, two, one, two, three	She is doing these steps in place when she's traveling forward. At this point, she really could be further forward all the way, still on the screen, but towards the edge of the screen, if she was to take bigger steps. And she could take bigger steps if she bends her knees and lowers her center of gravity and then extends her leg outward...

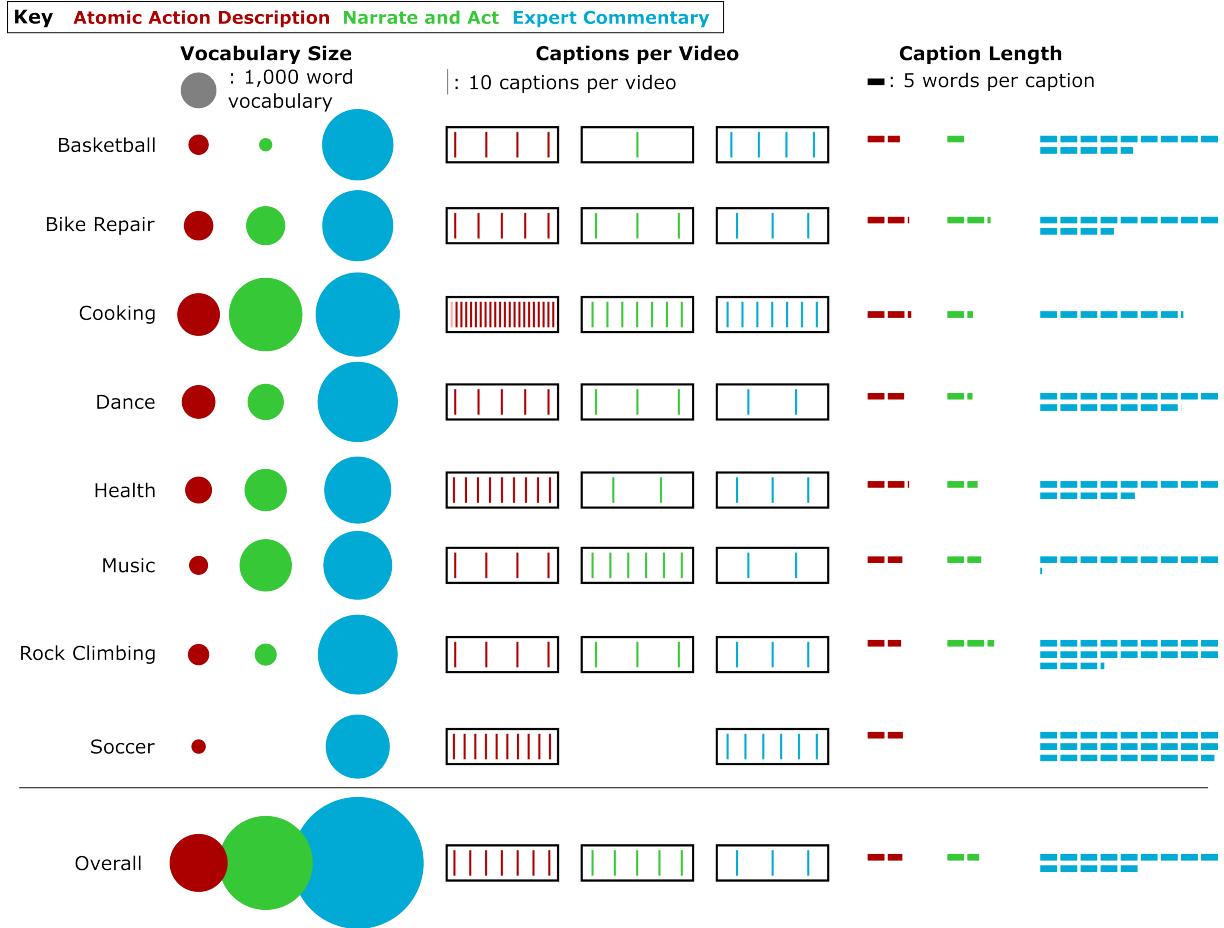
**Table D2:** Example excerpts from all three language types. Experts are charged with critiquing the performance of the participants, pointing out strengths and weaknesses and explaining how the participant's approach influences the quality of their skill demonstration. Narrate and act focuses more on what the camera wearer is doing and, sometimes, briefly why. Atomic action descriptions are about the specific actions seen.

Category	1x Coverage	2x Coverage	# of Descriptions	Descriptions Per Minute	Unique Nouns	Unique Verbs
Basketball	778	116	50299	53.330 (+- 26.049)	201	134
Bike Repair	202	160	31317	24.891 (+- 9.555)	642	393
Cooking	360	266	189225	27.745 (+- 12.843)	1744	823
Dance	307	417	43663	30.852 (+- 13.915)	504	468
Health	299	97	43769	24.304 (+- 11.234)	619	384
Music	85	75	10695	4.278 (+- 8.969)	255	163
Rock Climbing	1270	103	32246	32.350 (+- 11.974)	301	224
Soccer	225	53	31253	38.467 (+- 23.957)	229	125
All	3526	1287	432467	31.293 (+- 20.209)	2924	1481

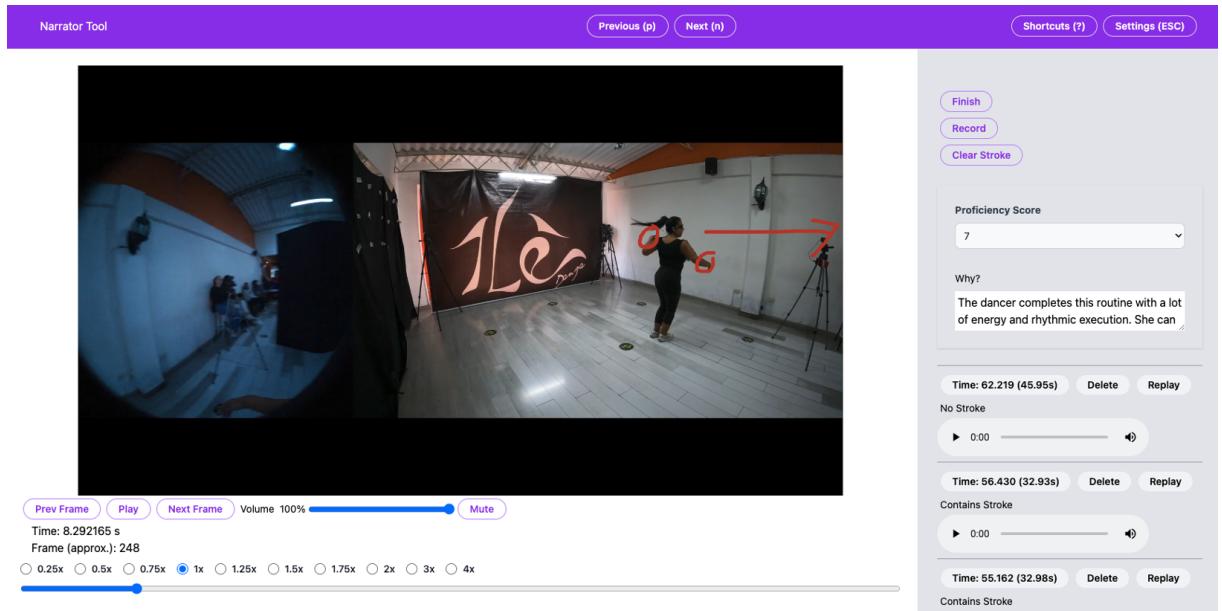
**Table D3:** Atomic action descriptions per-domain statistics.



**Fig. D7:** Word clouds for each scenario and annotation type. The vocabulary for atomic action descriptions typically focuses on the person’s hands and how they complete the actions (e.g. using left-/right/hand) whereas narrate and act describe the high level goals/objects. The expert commentary has the largest variety of words, including specialist words for each scenario such as swab/solution for health and axle/valve for bike repair.



**Fig. D8:** Comparisons between the vocabulary size (left) number of captions per video (center) and length of caption (right) for the atomic action description, narrate and act, and expert commentaries. Statistics are shown both per scenario and over the entire dataset. We see that the expert commentary tends to use a much larger vocabulary and more lengthy statements, since commentators are giving more elaborate statements of advice and explanation. The temporal density of the atomic action descriptions is greater than the other two forms, since the annotators are pausing to describe every single action of the camera wearer. Narrate-and-act comments use a vocabulary size in between the other two, reflecting the more free-form speech (compared to the written atomic actions) is used. Trends are mostly similar across scenarios, with the most noticeable differences being the temporal density; it is particularly high for both cooking and soccer. In the former, there are many procedural steps, whereas in the latter there are many instances of the drill being executed.



**Fig. D9:** Our expert commentary web tool called *Narrator* provides an easy-to-use platform for experts. Experts can stream video, record audio commentaries, and provide proficiency ratings and justifications. The tool also supports drawing on the video feed (see red arrow and circles on the right frame), allowing for manual spatial grounding during commentary.

## E Benchmarks: annotations and baselines

We employed a unified dataset split for EgoExo4D applicable for all benchmarks. Splits were defined at the “take” level, i.e. each take was allocated to either the training, validation, or testing set. Any derivatives from individual takes (e.g. segment clips) simply inherit the split assignment from the original take.

Additionally, the splits were stratified according to activity types and proficiency scores. This was done to ensure a balanced distribution of data across the training, validation, and testing sets, with proportional representation of activities and proficiency levels. To prevent data leakage, participants and all their associated artifacts were assigned exclusively to one of the splits. This constraint was particularly important because several participants contributed multiple takes, and it was crucial to avoid sharing their data between the splits. In the end we divided a total of 5045 takes into 3082 training, 842 validation and 1121 testing takes.

### E.1 Relation

#### *Annotations*

We used a multi-stage annotation process for annotating paired ego-exo videos:

- *Stage 0: Object Enumeration.* Annotator marks each object that is active at some point of the egocentric video with a bounding box in a frame where it is clearly visible and provides a free-form textual description.
- *Stage 1: Egocentric video annotation.* Annotator watches the egocentric video and is also shown (a) text and (b) a bounding box for one of the objects annotated in the previous stage. Annotator then marks a segmentation mask for that object in all the video frames where the object is visible. Segment Anything (Kirillov et al., 2023) is leveraged to generate segmentation masks efficiently using only point clicks.
- *Stage 2: Exocentric video annotation.* As shown in Figure E10, the annotator watches a temporally synchronized exocentric video and is also provided with the (a) text and (b) several ego segmentation masks of this object. Annotator

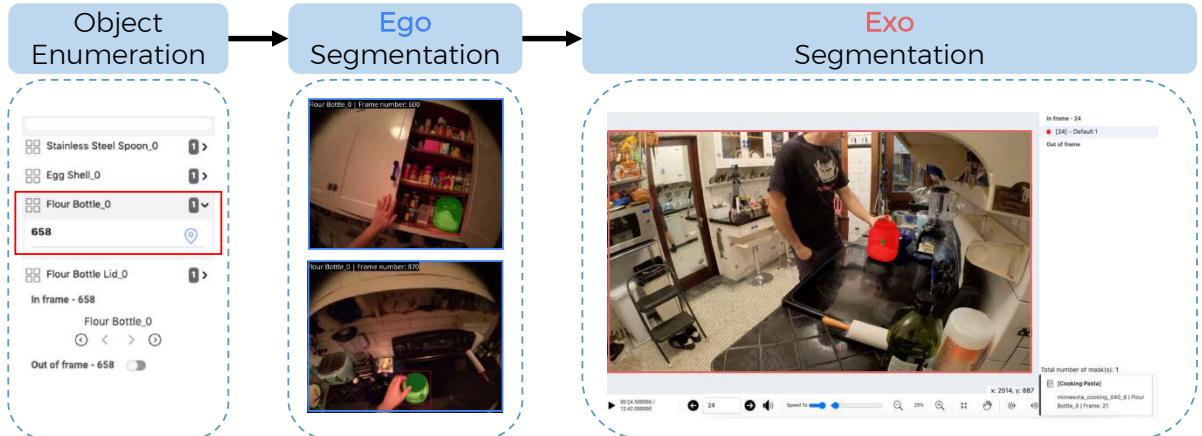
Scenario	# Takes	# Objects	# Ego Masks	# Exo Masks
basketball	394	602	21820	31165
bike	210	714	53886	71763
cooking	478	3481	549507	888384
health	127	570	77596	86585
music	112	153	33624	5599
soccer	12	22	2411	2475
Total	1335	5566	741965	1091135

**Table E4: Relation annotation statistics.** We show statistics for each scenario including the number of takes, total number of objects annotated and the number of egocentric and exocentric segmentation masks.

then marks a segmentation mask for this object in all the exo video frames, whenever the object is visible.

*What are the objects of interest?* We focus on objects that are *active* at some point during the execution of the activity. These objects are not only interesting because they are essential to the activity, but they are also challenging to track, since they are moving/changing state. In particular, our annotation guidelines requested annotators to list (a) objects that the camera-wearer interacts with through their body or tools; (b) other objects that are relevant to the activity (e.g., supporting surfaces like kitchen top); and (c) body parts (hands and legs). Note that every time an object changes visual state (adopting the Point-of-No-Return definition from (Grauman et al., 2022)), it is marked as a new object (e.g., annotators list *tomato* and *sliced tomato* as two distinct object instances).

*Which objects to annotate with masks?* For scenarios that involve few objects (Music, Basketball and Soccer), we annotated all object instances. Instead, for Cooking, Health and Bike Repair we sampled object instances based on their frequency of occurrence and their size, due to time and budget constraints. In particular, we binned each object annotated in the Object Enumeration stage into bins based on their frequency of occurrence across the dataset (high, low) and object size (small, large). We then uniformly sampled object instances from these bins while accounting for annotation time and budget and proceeded with segmentation mask annotations. We ignored all objects with area < 150 pixels. For Cooking,



**Fig. E10: Multi-stage annotation process for Ego-Exo Relation annotations.** After enumerating all active objects in the egocentric video, an object is selected and annotated with segmentation masks in all frames of the egocentric video. Then, annotators are given the exo video as well as the textual descriptions and sample egocentric segmentation masks for the object of interest, and mark segmentation masks for the specified object of interest in all the frames where it is visible.

specifically, we also filtered out a few objects such as spices, mixtures and liquids, as they tend to be too small to match in the exo view. Finally, we skipped exocentric mask annotations for objects that were visible in fewer than 10 frames of the egocentric video.

*What frame rate to annotate at?* We annotated segmentation masks at 1 frame per second, except for videos from the *Music scenario* which we annotated at 0.1 fps due to extremely long video durations.

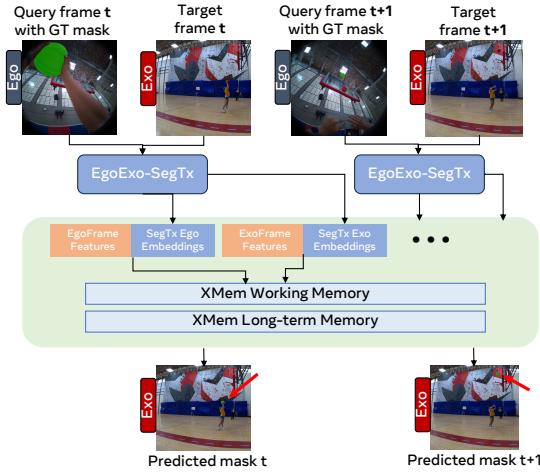
In total, our annotation process yielded segmentation masks for 5,566 objects in 1,335 ego-exo video-pairs. Approximately 4M million frames were annotated resulting in a total of 742K ego and 1.1M exo paired segmentation masks. Apart from this we also annotated 367K ego only segmentation masks. Collectively this results in a total of 2.2M segmentation masks. Table E4 shows a detailed breakdown per scenario for the paired masks.

### E.1.1 Ego-exo correspondence

#### Correspondence Baseline Implementation Details

*Spatial baseline model.* To adapt the architecture of SegSwap (Shen et al., 2022) for our correspondence problem, we additionally condition the model on the segmentation mask of the object of interest by feeding the query mask as a third input to the model. In particular, we first pass the egocentric frame, the exocentric frame and the query mask (as a binary mask) through a visual backbone network. We then flatten the resulting features into three sequences and pass them through the cross-image transformer with alternating self-attention and cross-attention layers. We first use the query mask features to attend to the features in the query view which are then used to cross-attend over features from the target view. This allows the model to reason over features from both views conditioned on the input mask. The resulting sequences for both views are “unflattened” and passed through a decoder to predict object segmentation masks in both views. We also pass the target view features through a classification head to classify if the query object is visible in the target view.

We train the model to perform mask prediction using a point-wise binary cross-entropy loss



**Fig. E11:** Overview of our spatio-temporal XView-XMem baseline model for the correspondence task.

and a dice loss over the predicted and ground truth masks. We use only pairs of frames where the object of interest is visible on both views and apply the losses on predicted masks in both the views. During inference, we only consider the mask predicted in the target view and discard the predicted mask in the query view. We train the head performing Visibility classification using a binary cross-entropy loss on all the frames of the sequence.

*Spatio-temporal baseline model.* To encourage the model to learn associations of the objects between egocentric and exocentric views, we train XView-XMem to track the object in a sequence of interleaved frames of egocentric and exocentric views, i.e., each egocentric frame is followed by an exocentric frame and vice versa, as shown in Figure E11.

To mitigate track drift (within and across views), we also explore feeding the XSegTx embeddings to the XMem working memory. Since these embeddings are trained to guide the mask decoder at each frame independently, they capture rich information about the object of interest. The extracted image features from the ResNet in XMem are fused with the encoded embeddings from multiple layers of SA (self-attention) and CA (cross-attention) layers of XSegTx. They are

then projected into keys and stored in memory for tracking.

For our spatial baseline model, we downsample the images to 480x480 resolution for all the views while using padding to keep the original aspect ratio of the images. For the image backbone we use the same ResNet50 (He et al., 2016) checkpoint as SegSwap and freeze its weights during training. Our cross-image transformer architecture also follows (Shen et al., 2022). We use a batch size of 32 and Adam (Kingma and Ba, 2014) as our optimizer with a learning rate of 0.0002 which decays to 0.0001 after 50,000 iterations. We run all our experiments on a single Nvidia RTX A6000 GPU for 200,000 iterations.

For our spatio-temporal baseline model, we use the same visual backbone (ResNet50 (He et al., 2016)) and architecture as XMem (Cheng and Schwing, 2022). Our only modification is in the information that gets inserted in the working memory at each frame. We first extract features from both ResNet and XSegTx for both both query and target frames. The corresponding features are then concatenated and projected to the original feature dimension through simple 2D convolution. We train on sequences of 8 interleaved ego and exo frames. The model is trained using AdamW as our optimizer with a learning rate of 0.00001 for 50,000 iterations and weight decay 0.05. The batch size is 8 clip pairs. We initialize our model with the original pretrained XMem, and keep both the ResNet backbone as well as our finetuned XSegTx models frozen. Note that we do not apply any data augmentations.

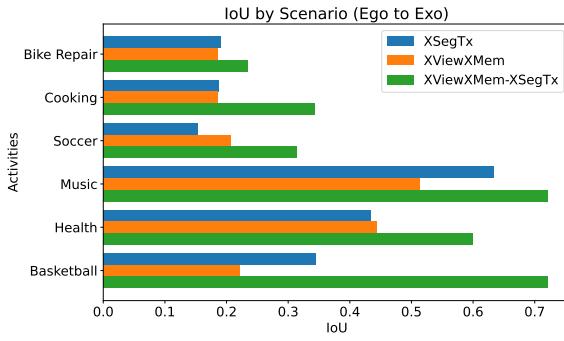
## Data

We use 1028 takes from the Ego-Exo dataset to train and evaluate models for this benchmark. In particular, we use the common split shared across benchmarks, with 838 takes for training, 201 takes for validation and 295 takes for testing. We extract pairs of images between egocentric and exocentric views which have corresponding object masks annotated for training. This gives us a total of about 193k pairs for training.

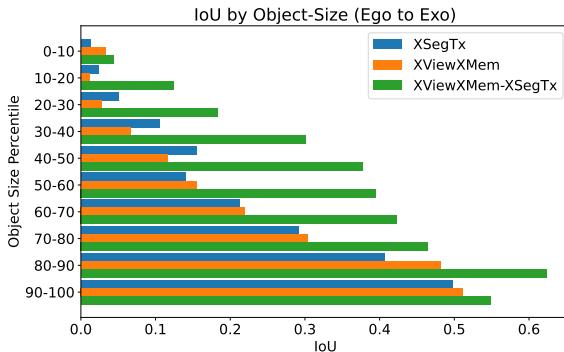
## Results

We break down our results across different activities in Fig. E12. We note that some activities are generally easier to model (e.g., basketball, soccer)

because of limited variation in object shape and appearance whereas some activities (e.g., cooking and bike repair) are much harder to model due to larger diversity in appearance, shape and size of the objects across views. We also explicitly evaluate our baselines on their ability to predict masks for very small objects. To do so, we split our validation set based on the predicted object size in proportion to pixels in the image. We see that, all our baselines struggle on very small objects and perform increasingly well on larger object sizes.



**Fig. E12:** Performance of both baselines per activity scenario.



**Fig. E13:** Correspondence evaluated across different object sizes in the target (exo) view. The object sizes range from  $7e^{-6}\%$  to 11% pixels in the target view.

### E.1.2 Ego-exo translation

In Table E5 we provide a break down of ego-exo translation results across different scenarios,

Scenario	IoU (%) ↑	LPIPS ↓
Basketball	14.5	0.41
Soccer	17.0	0.51
Music	4.5	0.30
Health	12.9	0.40
Bike	6.4	0.52
Cook	9.5	0.47

**Table E5:** Breakdown of ego-exo translation results per scenario for the subtasks of ego track prediction (IoU) and ego clip generation (LPIPS).

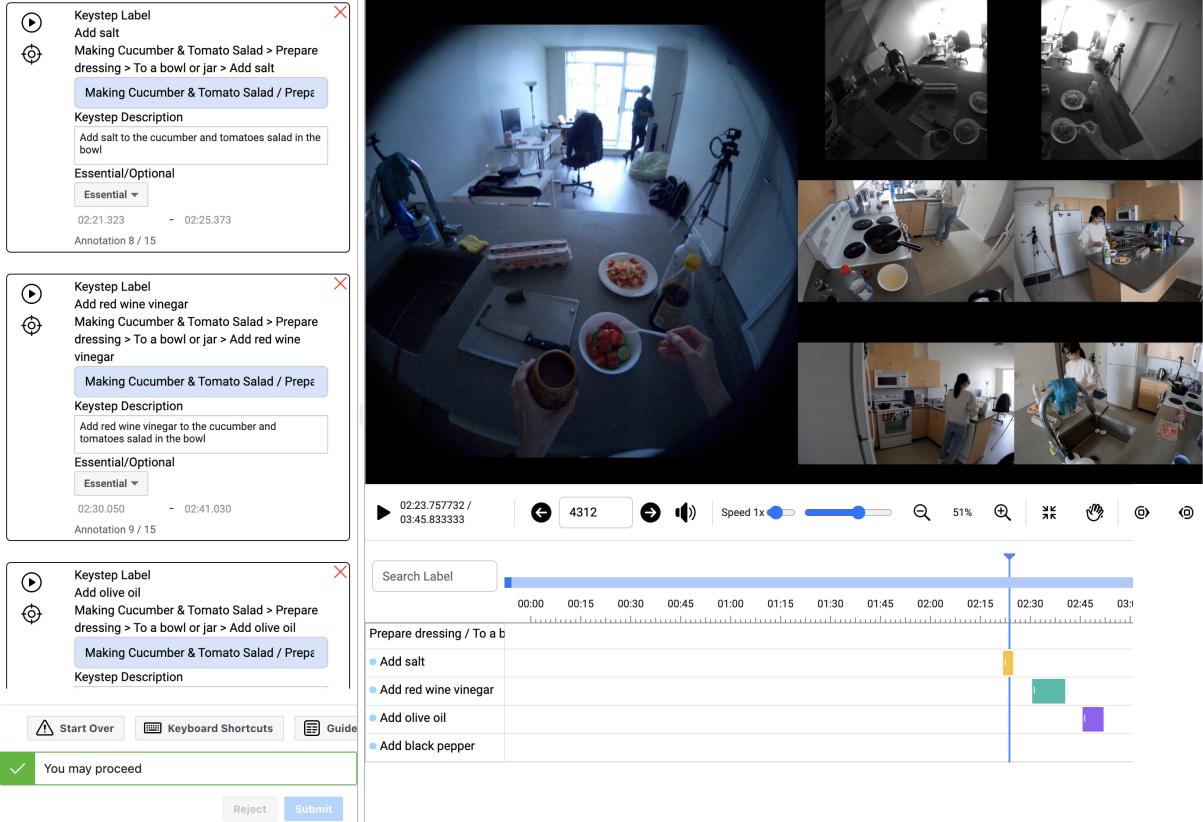
using GNT-mask for track prediction and DiT-pix for clip generation. We can observe similar trends for the two subtasks: the methods achieve better results in basketball and soccer scenarios than in bike and cook scenarios, which is reasonable as the objects in bike and cook scenarios are more complex and diverse.

## E.2 Fine-grained keystep recognition

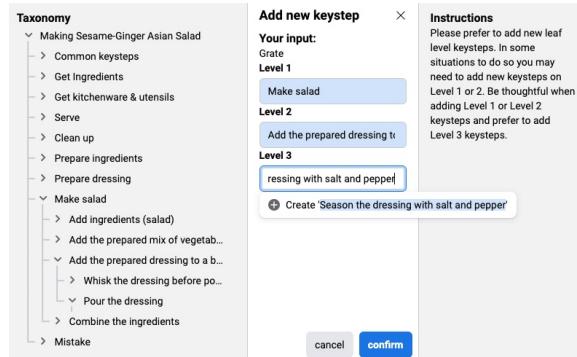
### Data annotation details

We collect manual annotations for keysteps (actions that contribute towards the completion of a procedural task) and build a keystep taxonomy in parallel. Figure E14 shows the annotation user interface. We provide annotators a composite view of time-synchronized ego and exo videos. Each keystep annotation contains the start and end timestamps, a category label, a natural language description, and a flag indicating whether the keystep is essential or optional for task completion. Annotators interact with a search widget which displays keystep labels with their complete path within a hierarchical tree, e.g., Making cucumber & tomato salad > Prepare dressing > To a bowl or jar > Add salt.

As the activities performed by the camera wearers are unscripted, it is not possible to establish a comprehensive keystep taxonomy prior to annotation. To address this challenge, we designed an iterative, data-driven process for taxonomy development. We first initialize the taxonomy using various resources including recipes and instruction articles from the Web. This initial taxonomy captures keysteps that are generally expected in the activities, but it is assumed to



**Fig. E14:** The keystep annotation tool shows a composite view of the time-synchronized ego-exo videos and the keystep time segment annotations. Each annotation consists of the start and end timestamps, a category label, a natural language description, and an essential/optional flag.



**Fig. E15:** Adding new keysteps to a taxonomy. Annotators utilize a specialized widget to introduce new keysteps at any level within the existing taxonomy hierarchy.

be incomplete for the specific variations the camera wearers performed in the recordings. Subsequently, in each iteration, annotators receive the

current taxonomy and are instructed to add new keysteps when they encounter actions not represented in it (see Figure E15). Any newly added keysteps are kept valid only for the duration of each annotation session and are not visible in other sessions. After a batch of videos have been annotated, we review the newly added keysteps to ensure their validity and update the taxonomy before repeating the process. We finalized the taxonomy after the third iteration, after which we re-annotated the entire set of videos with the final taxonomy for consistency.

#### Dataset splits

The keysteps in our dataset exhibit a very long-tailed distribution. To address this challenge, we set a cutoff threshold at 20 samples per keystep, limiting our analysis to 278 unique keysteps. For simplicity, we consider only the leaf node keysteps

Scenario	Takes			Ego Keystep Segments			Ego + Exo Keystep Segments			Taxonomy	
	Count	Duration (total / avg <sup>†</sup> )	Count (total / avg <sup>†</sup> )	Duration (total / avg <sup>‡</sup> )	Count (total / avg <sup>†</sup> )	Duration (total / avg <sup>‡</sup> )	Activity	Keystep			
Cooking	464	65.47h / 8.47m	19,034 / 41.02	58.08h / 10.99s	99,854 / 215.20	307.71h / 10.99s	11	527			
Bike repair	293	13.51h / 2.77m	2,573 / 8.78	11.82h / 16.54s	12,865 / 43.91	59.12h / 16.54s	4	82			
Health	331	18.72h / 3.39m	5,995 / 18.11	17.03h / 10.23s	30,723 / 92.82	86.99h / 10.23s	2	58			
Total	1,088	97.71h / 5.39m	27,602 / 25.37	86.94h / 11.34s	143,442 / 131.84	453.82h / 11.34s	17	664			

**Table E6: Keystep annotation statistics.** We report the statistics by grouping our 17 activities into three scenarios: cooking (11), bike repair (4), and health (2). Statistics are listed for takes<sup>†</sup> and keystep segments<sup>‡</sup>.

in the hierarchy. Exploring the hierarchical structure including parent nodes is a promising direction but we leave this as future work. In all, the dataset for keystep recognition comprises 130,979 segments, with an average duration of 11.34 seconds each. Specifically, the training set contains 74,342 segments, of which 14,326 are from the ego view and the rest from the exo view. The validation set consists of 23,636 segments, including 4,517 ego-view segments, and the test set has 33,001 segments with 6,373 in the ego view.

### Implementation details

We use clips of size  $8 \times 224 \times 224$ , with frames sampled at a rate of 1/32 for all baselines except for EgoVLPv2 (where we adhere to its pretraining scheme and sample 4 frames). The patch size is  $16 \times 16$ . For training, we resize the shorter side of the frame to a random value within the range of [256, 320], followed by randomly sampling a  $224 \times 224$  region from the resized video. For evaluation, we sample a single temporal clip in the middle of the video, scale down the shorter spatial side of the video to 224 pixels and select 3 spatial crops (top-left, center, bottom-right) from the temporal clip to cover a larger spatial extent within the clip. The final prediction is derived by averaging the scores obtained for these 3 crops. We train our model for a total of 100 epochs on 4 NVIDIA V100 GPUs with a batch size of 32. The model checkpoint yielding the best performance on the validation set is selected and evaluated on the test set.

### Results breakdown by keystep

We present a more detailed analysis of per-step performance in Figure E16, comparing training with ego-view videos and exo-view videos. We can observe that exo views show performance advantages over ego views in several steps, with the keystep ‘have a conversation asking different questions’ benefiting the most from exo. Conversely,

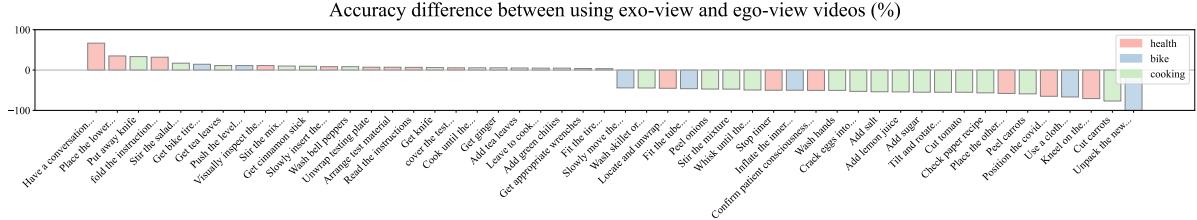
ego views are more effective in steps involving manipulation of small objects, like ‘cut carrots’ and ‘unpack the new tube’. This observation can be linked to the positioning of exo cameras, which are often placed further away from the subject, enabling them to capture a broader view, though possibly missing finer details of the activity. We hope these findings provide insight for future research on the effective use of exo-view videos during training.

## E.3 Energy-efficient multimodal keystep recognition

### Energy profiler

We adapt off-the-shelf profiler software built for PyTorch to compute the total multiply-accumulate operations (MACs) and memory transfer (MB) required to estimate total energy in Eqn. 1. The quantities are time-normalized — total energy consumption is expressed as power (mW). We describe each component of the profiler below.

- **Compute operations (MACs)** We use the native PyTorch FLOP counter to get the total FLOP count in the forward pass. We convert this to MACs (approximately 2 FLOPs = 1 MAC).
- **Memory transfer (bytes)** We consider GPUs as our processing device, and use the PyTorch memory profiler to get the list of all operations executed in the forward pass (`model.forward()` call) and their associated GPU memory usage. The total memory is the sum of the individual operation memory costs.
- **Sensor capture** For each modality, we measure the time for which it is active as the number of observations sampled containing the modality. We require that the sensors capture at least



**Fig. E16:** Keystep recognition evaluated per keystep label on a validation split, comparing training with only ego-view videos versus exo-view videos. The accuracy delta (exo-ego) is displayed, where a positive value indicates better exo and a negative value indicates better ego performance.

1 second worth of samples (roughly 100 samples) as energy consumption is ill-defined for an *instantaneous capture*.

### Energy tiers

As mentioned in the main paper, there is a natural trade-off between efficiency and better performance. Thus, we evaluate models in two tiers by setting a budget for the power consumption in each tier, namely 20 mW for the *high-efficiency* tier and 2.8W for the *high-performance* tier. We select the *high-efficiency* budget based on the energy consumption of current single-modality, efficient architectures (e.g., X3D-XS (Feichtenhofer, 2020)) with an eye to the future where multi-modal models operate within it. For the *high-performance* tier, we set the budget to a value that permits the use of powerful transformer-based action recognition models like LaViLa (Zhao et al., 2023). Once a model runs out of budget, in our setup it uses its latest prediction for all future steps.

### Complete baseline details

- **X3D-XS (Feichtenhofer, 2020).** This is a vision-only model comprising the X3D-XS feature encoder, which progressively expands the feature size and representational capacity of its layers, and later contracts them for achieving better performance-efficiency trade-off. This is the most lightweight model in our family of keystep predictors. The encoder has a depth factor of 2.2, and takes 4 RGB frames of size  $160 \times 160$  sampled at 15 fps, as inputs.
- **LaViLa (Zhao et al., 2023).** This is another vision-only model where the visual feature encoder is trained through CLIP-style video-language pre-training. To improve the feature

quality over vanilla CLIP-style pre-training, this method augments the number of video-text pairs by leveraging pre-trained large language models (LLM) to generate textual descriptions of un-annotated videos and rephrase existing narrations. In particular, we use the frozen TimeSformer (Bertasius et al., 2021)-Base (TSF-B) visual encoder pre-trained on the Ego4D dataset. To generate the feature for a target frame, the encoder samples 12 RGB frames of size  $224 \times 224$  at 30 fps from a time window centered around the target frame and pads the samples with the boundary frames on both ends to create a 16-frame clip.

- **Light-ASDNet (Liao et al., 2023).** This is an audio-only model that represents audio as spectrograms and efficiently encodes them by splitting 2D convolutions into 1D convolutions along the spectrogram temporal dimension (Liao et al., 2023). In our setup, the spectrograms are Kaldi (Povey et al., 2011)-compliant, and consist of 196 temporal windows and 160 Mel-frequency bins, respectively.
- **Audio-Visual Late Fusion (AV-LF).** This is an audio-visual model that does late fusion of visual features (encoded with X3D-XS or LaViLa) and audio features from Light-ASDNet by using linear layers.

### Experimental setup

We instantiate the task by considering keystep prediction episodes where the multimodal samples arrive in a streaming fashion. As mentioned above, we use vision and audio as our task modalities, where vision comprises RGB frames that are streaming at 30 frames per second (fps), and the audio modality is made up of time-aligned single-channel chunks that are 0.4 seconds long

and sampled at 16 kHz. However, the setup can be extended to include IMU, and potentially other sensors as well. We evaluate all models at the rate of 5 fps on a total of 211 test episodes of variable length, where the shortest episode is  $\sim$ 15 seconds, and the longest episode is  $\sim$ 34 minutes. We filter out episodes where all steps belong to the background class.

### Implementation details

We train all keystep prediction models for 150 epochs using the cross entropy loss. We use the AdamW (Loshchilov and Hutter, 2019) optimizer with an initial learning rate of  $10^{-4}$  and a weight decay of  $10^{-5}$ . We set the batch size to 512 for vision-only models, and 384 for audio-only and audio-visual models.

### Performance breakdown by modality

In Fig. E17, we present a detailed analysis of the keystep labels where the best vision-only model from the *high-efficiency* tier yields the maximum improvement or decline in performance compared to its audio-only counterpart. We observe that the vision-only model produces a large improvement over the audio-only model usually in steps where the activity does not produce distinctive sounds (e.g., *add green chillies*, *get celeries*, etc.). On the other hand, using audio alone helps the most when the activities involve sounds that are strongly indicative of the nature of the task (e.g., *stir fry egg mixture*, *cut butter*, etc.).

Finally, we envision that future work on this task will explore more sophisticated *learned* policies, potentially trained using reinforcement learning, in order to adaptively decide when to sample which modality instead of using fixed heuristics. Another promising direction is to investigate efficient transformer-based recognition backbones (Xu et al., 2021, Zhao and Krähenbühl, 2022) that can improve recognition performance without significantly affecting the model efficiency.

## E.4 Procedure Understanding

### Complete baseline details

*Graph-based baselines.* Graph-based baselines are composed of a keystep assignment and a procedural reasoning component (see Figure 23(a)).

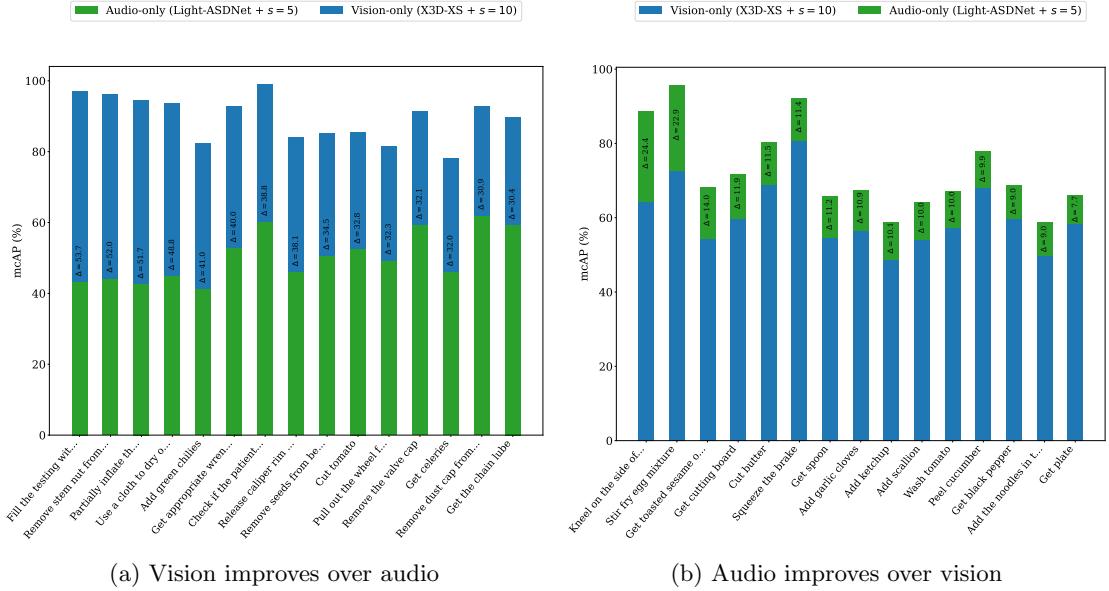
Keystep assignment (A1) is applied to obtain a pseudo-labeling of the provided video segments

when the supervision is at the procedure-level (i.e., segments are unlabeled and only keystep names are provided). This is achieved by means of an EgoVLPv2 model (Lin et al., 2022) pre-trained on ego-exo videos and narrations. Video segments and keystep names are projected to the shared video-language space using EgoVLPv2. We hence assigned each video segment to the closest keystep in the representation space according to the cosine distance. In the problem formulation with instance-level supervision (i.e., when keystep labels are available for all segments during both training and testing), we use ground truth labels instead of those obtained from keystep assignment. Additionally, we provide a baseline where the keystep assignment step is replaced with label predictions from the Keystep Recognition task (A2). Note that in the training set, segments with a confidence score below 20% have been discarded.

The procedural reasoning component (B) creates for each procedure a transition graph based on keystep co-occurrences. In the graph, each node represents a keystep category, while directed edges represent the probability of transitioning from one node to another one. An edge  $A \rightarrow B$  is assigned the following weight based on statistics collected from the training videos:

$$P(B|A) = \frac{\text{\# times keystep } B \text{ follows keystep } A}{\text{\# occurrences of keystep } A}$$

At test time, the graph is used to perform procedure understanding and answer the keystep-level questions. Specifically, given current segment  $s_i$ : 1) keystep  $y_{prev}$  is predicted as the previous keystep with confidence score equal to the transition probability  $P(y_i|y_{prev})$ , where  $y_i$  is the inferred or ground truth keystep label for segment  $s_i$ ; 2) segment  $s_i$  is predicted as optional based on the empirical probability  $\frac{\text{\# training videos containing } y_i}{\text{\# training videos}}$ ; 3) segment  $s_i$  is predicted as a procedural mistake with a score equal to the sum of the transition probabilities to  $y_i$  from keysteps  $y^{prev}$  that are missing from the keystep history, i.e.,  $\sum_{y^{prev}} [y^{prev} \notin S_{i-1}] \cdot P(y_i|y^{prev})$ , where  $[ \cdot ]$  is the indicator function; 4) keystep  $y$  is predicted as a possible missing keystep with probability  $[y_i \notin Y_{i-1}] \cdot P(y|y_i)$ ; 5) keystep  $y$  is predicted as a future keystep with probability  $P(y_i|y)$ .



(a) Vision improves over audio

(b) Audio improves over vision

**Fig. E17:** Improvement (left) or degradation (right) in keystep recognition performance per keystep label, when comparing the most efficient vision-only (X3D-XS ([Feichtenhofer, 2020](#)) +  $s = 10$ ) and audio-only (Light-ASDNet ([Liao et al., 2023](#)) +  $s = 5$ ) models from the *high-efficiency* tier. The plots show the 15 keysteps where improvement or degradation are largest.  $\Delta$  reports the amount of improvement/degradation.

*End-to-end baseline.* This baseline aims to provide an end-to-end approach to perform procedure understanding directly from the input clip. The baseline predicts previous keysteps, optional keysteps, and next keysteps by feeding video segment features extracted with EgoVLPv2 ([Pramanick et al., 2023](#)) to three dedicated MLPs. Figure 23(b) illustrates the architecture of the baseline. At training time, MLPs are supervised from the pseudo-labels obtained by graph-based baselines using Mean Squared Error (MSE) score to align the predicted probability distributions to the supervising ones. Missing keysteps and procedural mistakes are predicted from the outputs of the MLP components as in graph-based baselines.

## E.5 Proficiency estimation

### Annotations for demonstrator proficiency estimation

We derive annotations for this task from participant surveys (see Section C) and expert commentary (see Section D.1). Participant surveys contain responses to questions about prior experiences in the task such as “How many years have you been

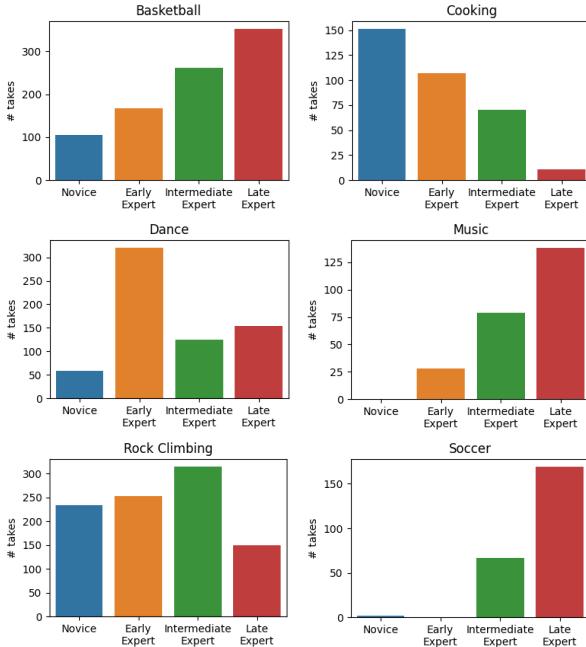
doing this task?”, and “Do you have any qualifications/professional training related to the task?” (see Table C1 for the complete list). On the other hand, expert commentary is performed by task-specific experts and includes 1 to 10 proficiency scores for each video from the participant (see Section D.1). After consulting with experts hired for each scenario, we designed scenario-specific conversion functions that use the surveys and expert commentaries to produce an estimate of a participant’s proficiency score (see Table E7). For example, in basketball and soccer, we use the years of experience to determine skill level since we found this to be an accurate indicator of skill based on analyzing the videos. On the other hand, to determine skill level in bouldering, we use the highest difficulty level of the route solved by the participant.

### Annotations for demonstration proficiency estimation

Table E8 shows examples of expert comments and corresponding annotation tags derived from them

Scenario	Novice	Early Expert	Intermediate Expert	Late Expert
Basketball	$X \in [0, 1)$	$X \in [1, 3)$	$X \in [3, 10)$	$X \geq 10$
Soccer	$X \in [0, 1)$	$X \in [1, 3)$	$X \in [3, 10)$	$X \geq 10$
Dancing	$X \in [0, 3)$	$X \in [3, 5)$	$(X \in [5, 10)) \vee ((X \geq 10) \wedge \neg P)$	$(X \geq 10) \wedge T$
Bouldering	$H \leq V3$	$H == V4$	$H == V5$	$H \geq V6$
Music (violin)	$(X \in [0, 3)) \vee (N \in [0, 500))$	$(X \in [3, 5)) \vee (N \in [500, 1000))$	$(X \in [5, 10)) \vee (N \in [1000, 10000))$	$(X \geq 10) \vee (N \geq 10000)$
Music (guitar)	$(X \in [0, 1)) \vee (N \in [0, 500))$	$(X \in [1, 3)) \vee (N \in [500, 1000))$	$(X \in [3, 10)) \vee (N \in [1000, 10000))$	$(X \geq 10) \vee (N \geq 10000)$
Music (piano)	$(X \in [0, 1)) \vee (N \in [0, 500))$	$(X \in [1, 5)) \vee (N \in [500, 1000))$	$(X \in [5, 10)) \vee (N \in [1000, 10000))$	$(X \geq 10) \vee (N \geq 10000)$
Cooking	$P < 3.5$	$P \in [3.5, 5)$	$P \in [5, 8)$	$P \geq 8$

**Table E7: Annotations for demonstrator proficiency estimation.** We designed scenario-specific conversion functions that take in participant surveys and expert commentary assessments to estimate proficiency of participants (i.e., novice, early expert, intermediate expert, and late expert). Legend:  $X$  = years of experience performing the task,  $T$  = professional training in the task,  $H$  = highest difficulty level solved by participant in bouldering,  $N$  = estimated number of times performing the task,  $P$  = average proficiency rating from expert commentary.



**Fig. E18:** Distribution of demonstrator proficiency scores per scenario.

indicating whether the comment suggests a good execution and/or tips for improvements.

### Baseline implementation details

**Demonstration proficiency estimation:** We use the TimeSFormer (Bertasius et al., 2021) architecture for the baseline. TimeSFormer is a video transformer designed for video action recognition/classification that introduces a novel decoupled spatiotemporal attention mechanism. We resize all videos to 448 pixels along the smallest dimension and use a clip size of 16 frames with a frame rate of 16 FPS. The models are trained to classify individual clips using the cross-entropy loss on 8 Quadro RTX 6000 GPUs for 15 epochs.

**Demonstration proficiency estimation:** We adopt ActionFormer (Zhang et al., 2022), a video action localization model for our experiments. Unlike traditional action localization that defines time windows as outputs, we instead perform timestamp regression since our annotations contain only a single point in time for each good execution or tip for improvement. We accordingly adapt ActionFormer’s post-processing strategy and evaluation metrics. In our task, the predicted timestamps correspond to frames retained after non-maximum suppression (NMS). We remove the regression head of ActionFormer and infer the predicted timestamps from the indices of frames retained after NMS. We also modify the NMS module of ActionFormer to rely on the  $\mathcal{L}_1$ -distance between predicted timestamps instead of the tIoU

Scenario	Expert comment	Good execution	Tips to improve
Basketball	Nice release. I like the follow through here. You'd like to see the guide hand maybe up a little bit higher on the release of that shot. Maybe to give it better ball control when you're letting go of the shot.	Yes	Yes
Basketball	Great footwork, left foot take off, lifting of the right knee and extending that body up. Love how he's looking up, checking out the backboard, shooting hand behind the basketball. Nice job.	Yes	No
Basketball	He's also really far away from his body and the more he can keep his arm up by his ear, it will give him the most opportunity to make the basket without the defense interrupting.	No	Yes
Bike repair	It's a great method to always double check or do a pre-check before beginning work on a bicycle to make sure the issue that you are working to fix is the only issue that is occurring. If not, you could find a secondary issue or something else that may be greater than the one you are currently working on.	Yes	No
Bike repair	As you can see she clearly slipped on loosening the nut which essentially creates damage to the surface of the nut itself and can round out the nut.	No	Yes
Bouldering	The climber was efficiently able to position herself with one hand on each hold at the start and had, once her hands were positioned, she matched her feet on the hold and efficiently moved to the next hold.	Yes	No
Bouldering	And since she popped out and is swinging out, she can't really keep the tension through her one arm because she's so locked off. So it caused her to kind of just fall off the wall and lose all tension throughout all of her body.	No	Yes
Cooking	You can see there, she's not able to stir properly. She has to push it around, which means that the lime is not gonna be very evenly distributed among the pieces of tomato and cucumber.	No	Yes
Cooking	Using a grinder for fresh pepper is an excellent way to get a lot of flavor. The fresh grind of pepper as opposed to buying already ground pepper really expels the oils and everything in those peppercorns and allows the flavor to be as big as it can possibly be.	Yes	No

**Table E8: Annotations for demonstration proficiency estimation.** We annotated expert comments about a participant’s task execution with tags indicating whether each comment describes a good execution or suggests tips for improving skills. Note that the same comment might describe one aspect of the task as being good while suggesting improvements in another aspect (e.g., see row 1).

between segments used in (Zhang et al., 2022). During training, we keep the classification loss from (Zhang et al., 2022) and replace the regression loss with the loss function defined in (Kwak et al., 2020). We train our models with Omnivore features (Girdhar et al., 2022) extracted with a clip size of 32 frames and a stride of 16 frames from all the videos.

#### *Scenario-specific results on demonstrator proficiency estimation*

We show scenario-specific results in Table E9. TimeSFormer achieves good performance with the

egocentric view in cooking since a close-up view of the objects of interest and hand poses is essential to assessing skill in these scenarios. On the other hand, the model performs better with the exocentric view in bouldering on the test data since the overall body pose is a useful indicator of proficiency. Unfortunately, it fails to improve over the majority-class baseline in most scenarios on the test splits except basketball, highlighting a distribution shift between the val and test splits.

Scenario	Majority-class	Val			Majority-class	Test		
		Ego	Exos	Ego + Exos		Ego	Exos	Ego + Exos
Basketball	35.66	<b>56.64</b>	<b>56.64</b>	51.75	46.71	<b>79.64</b>	71.86	76.65
Cooking	50.94	<b>56.60</b>	45.28	49.06	<b>39.54</b>	19.77	33.72	32.56
Dancing	<b>43.31</b>	42.52	31.50	32.28	<b>46.62</b>	43.24	40.55	42.57
Music	44.44	<b>69.44</b>	50.00	58.33	<b>73.24</b>	56.34	45.07	50.70
Bouldering	0.00	<b>27.67</b>	18.24	17.61	19.57	39.13	<b>43.04</b>	<b>43.04</b>
Soccer	<b>74.42</b>	65.11	37.21	41.86	<b>72.73</b>	65.15	31.82	30.30

**Table E9: Breakdown of results for demonstrator proficiency estimation across scenarios.**  
Top-1 accuracies per scenario for the TimeSFormer model.

	MPJPE	MPJVE
Location-based (baseline)	18.51	0.64
Levelwise att ViT	18.09	0.62
UCB Ego	17.19	0.57
Multi-Scale Model Fusion	15.32	0.55

**Table E10: Leaderboard of Ego-Exo4D Body Pose challenge at EgoVis CVPR 2024.**

## E.6 Ego Pose

### CVPR 2024 challenges

We organized two Ego Pose challenges as part of the EgoVis workshop at CVPR 2024, aiming to encourage researchers to explore and utilize the dataset. We had eight participants in total, divided evenly between the Body Pose challenge and the Hand Pose challenge.

#### Body Pose

There were a total of four submissions to the hand pose challenge, with three outperforming the best baseline method, as shown in Table E10. The best participant outperformed the baseline by 3.19 cm in MPJPE.

**First place: Multi-Scale model fusion** (*Baoqi Pei, Yifei Huang, Guo Chen, Jilan Xu, Yicheng Liu, Yuping He, Kanghua Pan, Tong Lu, Limin Wang, Yali Wang, Yu Qiao*)

Instead of implementing a new baseline model, this approach leverages the existing location-based baseline and focused on improving performance by fusing predictions from multiple models with varying numbers of transformer layers (ranging from 1 to 6). This approach allowed them to address variations in data distributions and stabilize predictions. By combining the outputs from the different versions of the baseline, they achieved an improved MPJPE of 15.32.

**Second place: UCB Ego** (*Brent Yi, Vickie Ye, Georgios Pavlakos, Lea Müller, Maya Zheng, Yi Ma, Jitendra Malik, Angjoo Kanazawa*)

The submission by UCB EGO to the Ego-Exo4D Body Pose Challenge centers around their development of a conditional human motion diffusion model that operates using 30Hz SLAM pose data to directly sample SMPL human body parameters. The model was trained exclusively on the AMASS dataset, which poses a limitation of domain transfer to the Ego-Exo4D dataset. To address this limitation, the team introduced an AdapterNet to estimate floor height and route procedural tasks to the baseline model while focusing the diffusion model on physical activities. This hybrid approach allowed the team to achieve a MPJPE of 17.19 cm.

**Third place: Levelwise attention ViT** (*Congsheng Xu, Jinfan Liu, Yifan Liu, Shuwen Wu*)

The authors proposed a model that combines two Vision Transformer (ViT) structures to leverage both coarse-grained and fine-grained information for 3D human pose estimation. They used a baseline ViT model with 8 attention heads and 3 layers for coarse-grained estimation, and a Huge-ViT model with 16 attention heads and 32 layers for fine-grained estimation. They combine the outputs from both ViTs through a weighting strategy. By assigning a higher weight to the fine-grained results, the method balances local detail and global information. This weighting approach resulted in a MPJPE of 18.09 cm.

#### Hand Pose

There are in total 4 submissions to the hand pose challenge, with 3 outperformed the baseline method POTTER (*Zheng et al., 2023*) as shown in Table E11. The best participant outperformed

	MPJPE	PA-MPJPE
POTTER (baseline)	28.94	11.07
PCIE EgoHandPose	25.51	8.49
Hand3D	30.52	9.30
POTTER-ensemble	28.68	10.24

**Table E11:** Leaderboard of Ego-Exo4D Hand Pose challenge at EgoVis CVPR 2024.

the baseline model by 11.85% in MPJPE, and by 23.31% in PA-MPJPE.

**First place: PCIE EgoHandPose** ([Chen et al., 2024](#))

The authors propose the Hand Pose Vision Transformer (HP-ViT). The HP-ViT comprises a ViT backbone and transformer head to estimate joint position in 3D, utilizing MPJPE and RLE loss function. To be more specific, the model employed regression loss based on RLE ([Li et al., 2021](#)) to minimize the gap between output and input distributions. The experiments show that the model with ViT-Huge ([Dosovitskiy, 2020](#)) as backbone achieves the best performance. The ensembling model with different setting also contributes to decrease the overall error.

**Second place: Hand3D** ([Pavlakos et al., 2024](#))

The authors apply recently introduced HaMeR out-of-the-box on the images of the EgoExo4D Challenge, and observe very strong performance. HaMeR is a feed-forward model that takes as input a single image of a hand and hand side (left or right), and estimates a 3D reconstruction of the hand in the form of the MANO parametric hand model ([Romero et al., 2017](#)). HaMeR adopts a fully transformerized architecture design using a ViT-H ([Dosovitskiy, 2020](#)) backbone, followed by a transformer head. The authors further show the limit of the method, particularly when the occlusions/truncations are very extreme (only a few visible fingers), the wrist location or hand orientation is ambiguous, and finally, if the original hand bounding box crop is not very accurate.

**Third place: POTTER-ensemble** (*Baoqi Pei, Yifei Huang, Guo Chen, Jilan Xu, Yicheng Liu, Yuping He, Kanghua Pan, Tong Lu, Limin Wang, Yali Wang, Yu Qiao*)

The authors handle the problems through a multi-scale model fusion method. To be specific, the authors enhance the baseline method POTTER ([Zheng et al., 2023](#)) by (1) including multiple upsampling dimensions, specifically 128, 256, and

512, (2) integrating a dynamic pooling operation before the 3D convolution operation of the model, and (3) combining/ensembling these models to improve generalization capabilities.

## F Detailed contribution statement

This project is the result of a large collaboration between many institutions over the last two years. Initial authors represent the leadership team of the project. Kristen Grauman initiated the project, served as the technical lead, initiated the recognition and proficiency benchmarks and expert commentary, and coordinated their working groups. Andrew Westbury served as the program manager and operations lead for all aspects of the project. Lorenzo Torresani led development of the capture domains, initiated the relation and ego-pose benchmarks, and coordinated their working groups. Kris Kitani led development of the multi-camera rig and supported the Ego-Exo4D engineering team on all aspects of the data annotation and organization. Jitendra Malik served as a scientific advisor. Authors with stars (\*) were key drivers of implementation, collection, and/or annotation development throughout the project. Authors with daggers (†) are faculty and senior researcher PIs for the project.

### *Camera rig*

Rawal Khirodkar proposed the hardware and software specifications for the ego-exo camera rig and helped design the capture protocol. Sean Crane investigated various hardware setups leading to the final rig configuration and helped draft the capture guidelines including recommended gear. Devansh Kukreja developed the sync and take separation algorithms, experimented with different equipment options (e.g., camera, timecode boxes, mount options), and designed the interface to transfer data; he also managed the ingestion pipeline, the collaboration with Aria, the integration of their code for EgoExo, and usage of the .vrs files.

### *Aria*

Jing Dong and Vijay Baiyya were responsible for obtaining camera poses, calibration, point-clouds and eye gaze using Aria MPS, created the 3D/4D visualizations for the paper and supplementary material, and acted as main contact points from the Aria team throughout the program; with Jing leading the algorithm development and verification, and Vijay leading the

Aria MPS workflow and infrastructure development. Jakob Engel acted as technical and scientific advisor, and led the team that built the Aria Localization and Point Cloud algorithms. Kiran Somasundram helped design the capture setup and time-synchronization. Xiaqing Pan helped to align the Aria engineering team to support the EgoExo4D project. Mingfei Yan, Prince Gupta, and Sach Lakhavani acted as product managers of Aria and organizational leads for the successful use of Aria in the program. Kelly Forbes helped setting up agreements and working through the legal requirements of using Aria devices for recording the EgoExo4D dataset across the globe. Richard Newcombe initiated the Aria/Ego4D collaboration and acted as a scientific advisor throughout the program. Furthermore, we want to acknowledge the contribution of the entire Project Aria team as listed in ([Engel et al., 2023](#)), including Carl Ren and Sean Diener leading the Aria software and hardware engineering organization, and Renzo De Nardi as technical lead for the Aria device.

### *Data collection*

**Los Andes University** Pablo Arbeláez - lead coordinator for data collection and collaborator on the overall project design; Maria Escobar - data collection for all phases, design of the collection setup and workflow, data inspection, ingestion, encoding, and metadata generation; Cristhian Forigua - data collection for all phases, participant recruitment, consent forms design, data inspection, communication with recording sites; Cristina González - data collection for phase 1, design of the collection setup and workflow, IRB management; Angela Castillo - data collection for phase 2, manual data inspection, and data analysis.

**Georgia Tech** James M. Rehg - lead coordinator for data collection and protocol design, and overall project manager; Bikram Boote - lead coordinator for data collection, including recruiting and ingestion; Fiona Ryan - contributed to data collection; Audrey Southerland - lead coordinator for IRB development, contributed to recruiting.

**National University Singapore** Mike Zheng Shou - lead coordinator for data collection and protocol design, and overall project manager;

Joya Chen - contributed to protocol design, camera setup design, data collection for all phases; Jia-Wei Liu - contributed to protocol design, camera setup design, data collection for all phases; Xinzhu Fu - contributed to data collection for all phases; Chenan Song - contributed to data collection for all phases.

**Meta** Andrew Westbury was the lead for data collection at our site, selecting scenarios, organizing capture sessions, recruiting participants, organizing and transferring data, and obtaining required approvals. In California, Hao Tang and Kevin Liang also supported all these functions, focused on bike repair. In New York, Devansh Kukreja and Alex Dinh lead collection for cooking scenarios. Miguel Martin also supported California-based collections and organized our local camera rig. Chefs Eton Chan and Dominic Ainza supported all culinary collections with technical guidance, recruitment, and coordination. Dimitri Elston coordinated and was the technical lead on bike collections. Adrian Salas supported pilot bouldering collections in California. Across all Ego-Exo4D collections, Devansh Kukreja continuously communicated and refined the recording procedure with universities, and problem-solved local recording issues.

**University of North Carolina at Chapel Hill** Gedas Bertasius - lead coordinator for data collection; Md Mohaiminul Islam - the main contributor to data collection and metadata processing across all scenarios; Wei Shan - contributed to data collection and metadata processing for the music and soccer scenarios; Jeff Zhuo - contributed to data collection and metadata processing for the soccer scenarios; Oluwatumininu Oguntola - contributed to participant recruiting and data collection for the music scenario.

**Carnegie Mellon University** Rawal Khirrodkar developed the automatic 3D body keypoints extraction pipeline and collected a subset of the soccer, bike mechanic, and cooking sequences for the CMU portion of the dataset. Sean Crane was in charge of the data collection, IRB documents, capturing data, working with participants and processing the data for CMU. Abrham Gebreselasie ran the actionformer-based baseline for demonstration proficiency benchmark. Eugene Byrne served as the engineering lead for the initial design and implementation of the dataset, camera rig, processing pipeline and keystep annotations

while at Meta. Subsequently at CMU, he assisted in the recognition benchmarks, implemented Ego-Exo transfer [Li et al. \(2021\)](#) (1 of the 3 baseline methods for keystep recognition) and the initial implementation of keystep action detection [Zhang et al. \(2022\)](#), and assisted in annotation/data quality generally.

**Simon Fraser University** Sanjay Haresh was the lead coordinator for data collection, including recruiting, data ingestion, and data analysis. Yongsen Mao also contributed to the data collection pipeline, recruiting, data ingestion, metadata annotations, and statistics computations. Manolis Savva advised on data collection, protocol design, and overall project management. We acknowledge the assistance of Hanxiao Jiang and Armin Kavian with recruitment and data collection.

**University of Pennsylvania** Edward Zhang led data collection efforts at UPenn and played a key role in subject recruitment, subject information collection, and on-site data recording. Jin Xu Zhang led data management, information logging, and data transfer. Shan Su is the overall project lead, focusing on determining good camera configurations based on 3D reconstruction feasibility and fixing issues in data post-processing of time synchronization and take separation.

**University of Tokyo** Yoichi Sato served as the primary coordinator for data collection, while Ryosuke Furuta was responsible for data collection across all three scenarios, participant recruitment, and IRB submission. Zecheng Yu and Masatoshi Tateno provided support for data collection and were responsible for managing and transferring data. Takuma Yagi helped with the IRB submission process.

**Indiana University** David Crandall oversaw the overall effort at Indiana University, including protocol design and data collection. Leslie Khoo led the IRB protocol design and compliance, arranged logistics such as ordering equipment, and designed and oversaw the cooking scenario data collection. Yuchen Wang and Ziwei Zhao co-led the participant recruitment and data collection for all three scenarios. Ziwei Zhao led data preparation and transfer. We also acknowledge Manasi Swaminathan who assisted with data collection and video synchronization.

**IIT-Hyderabad** Avijit Dasgupta was the lead on the ground for data collection in Hyderabad helping in organizing capture sessions, data collection, and managing and transferring data. Siddhant Bansal helped in the early stages with IRB application, consent forms, and pilot studies. C. V. Jawahar was the lead coordinator for data collection helped in selecting the scenarios, and recruiting the participants.

**University of Minnesota** Hyun Soo Park oversaw the overall effort at the University of Minnesota, Twin Cities, including protocol design and data collection. Zachary Chavis led the IRB protocol design and compliance, arranged logistics such as ordering equipment, participant recruitment, and designed and oversaw all scenarios of data collection. Anush Kumar assisted data collection for all scenarios.

### *Language annotations*

Kevin J Liang co-developed the expert commentary guidelines, interviewed and onboarded experts, helped test and suggest features for the narrator tool, and contributed to program management; he also developed the atomic action descriptions guidelines, helped coordinate annotations, and contributed to paper writing. Michael Wray contributed to the definition of the expert commentary guidelines; provided feedback for experts; co-developed the narrate-and-act guidelines and the pre-task/post-task questionnaires; and contributed to paper writing. Kristen Grauman proposed the expert commentary idea, co-developed the guidelines, interviewed and provided feedback to experts, and contributed to paper writing. Andrew Westbury implemented expert commentary, recruiting, mobilizing and managing our experts and workplan. Miguel Martin contributed to the atomic action descriptions annotation guidelines and produced the annotation files and associated tutorial code, and for expert commentary, he authored the initial version of the Narrator tool, transcribed the commentaries, and produced the annotation files and associated tutorial code.

Changan Chen contributed to the development of the narrator tool; provided feedback for experts; and contributed to paper writing. Siddhant Bansal contributed to the design of narrate-and-act, user

questionnaires, object dictionaries, expert commentary cooking. Dima Damen proposed narrate-and-act data collection and user questionnaires and contributed to their design, and also contributed to expert commentary for cooking scenarios. Tiffany Davis provided significant program management support throughout expert commentary. Devansh Kukreja built the render flow to generate video collages for annotations.

Domain resource people from our consortium were Dima Damen and Michael Wray [cooking], Kristen Grauman and Changan Chen [soccer], Gedas Bertasius [basketball], Kristen Grauman and Jianbo Shi [music], Andrew Westbury [health and bike repair], Kevin Liang [dance], Pablo Arbelaez and Maria Escobar [bouldering]

Ego-Exo4D's panel of expert commentators is: **Soccer**: John Bello, Phillip O'Kennedy, Lee Bakewell, Radcliffe McDougald, Thomas Harris **Music**: James Peterson, Trevor Minton, Andrea LaPlante, Ethan Fallis, Alex Rogers, Jacqueline Burd **Health**: Jasmine Higa, Angela Liszewski, Kristin Blanset, Melissa Robinson, Sonya Johnson **Dance**: Rolanda Williams, Deanna Martinez, Enya-Kalia Jordan, Rachel Repinz, Yauri Dalencour, Kathryn Hightower **Cooking**: Mark Manigault, Mary Drennen, Tiffany Davis, Reginald Howell, Rosanne Field, Donnie Murphy, Kiet Duong, Laura de Vera, Keegan Taylor **Bouldering**: Daniel Ramos, Mike Kimmel, Roy Quanstrom, Christopher Deal, Carmen Acuna, Kelsey Hanson **Bike repair**: Cesar Pineda, Walker Wilkson, Frank Trotter, Cordell Bushey, Dimitri Elston, Sam Arsenault, Aaron Hill **Basketball**: Elizabeth Blose, Raven Benton, Joseph McCarron, Cornelius Gilleyen, Cecil Brown. Aaron Jones

### *Benchmarks*

**Ego-exo correspondence** Manolis Savva co-led the correspondence benchmark and contributed to the task definition, the annotation guidelines, the baseline design, and paper writing. Effrosyni Mavroudi co-led the correspondence benchmark and contributed to the task definition, the annotation guidelines, the baseline design, and paper writing. Lorenzo Torresani contributed to the task formulation. Sanjay Haresh developed the spatial baselines and contributed to data analysis, experimental results, and paper writing.

Yongsen Mao developed the spatiotemporal baselines and contributed to data analysis, experimental results, and paper writing. Suyog Jain formulated the annotation pipeline, developed annotation tools and contributed to annotation guidelines and paper writing. Santhosh Ramakrishnan contributed to the annotation guidelines and the formulation of the annotation pipeline. Xitong Yang contributed to the annotation guidelines and the task definition. We would like to acknowledge Hanxiao Jiang for helpful discussions and preliminary ideas on baseline implementation. Devansh Kukreja built the render flow to generate frame-aligned videos of each camera for each take as model input.

**Ego-exo translation** Lorenzo Torresani co-led the translation benchmark, developed the task formulation, and advised the baseline development. Judy Hoffman co-led the translation benchmark and advised the baseline development. Feng Cheng led the baseline development and implemented the pix2pix and DiT models for track prediction and clip generation. Mi Luo implemented the GNT baseline and the evaluation pipeline for ego track prediction. Ziwei Zhao contributed the pix2pix baseline for multi-frame input and led the evaluation for ego clip generation. Huiyu Wang advised the baseline development and contributed to the task definition, the baseline design, and the metric selection and analysis.

**Fine-grained keystep recognition** Tushar Nagarajan co-led the keystep recognition benchmark, co-developed the task formulation, and advised the baseline development. David Crandall co-led the keystep recognition benchmark and contributed to the task formulation. Yale Song co-led the keystep recognition benchmark and led the keystep annotation effort, including design of annotation guidelines, taxonomy development and coordination of annotation workflows; he also contributed to the task formulation, advised baseline design, and facilitated the delivery of EgoVLPv2 pretrained backbone. Triantafyllos Afouras contributed to the taxonomy definition, managed the labeling effort, and developed software for post-processing the annotations. Zihui Xue led the baseline development effort, implemented the TimeSFormer, EgoVLP, VI Encoder and Viewpoint distillation baselines, and performed analysis of results. Eugene Byrne contributed to the

taxonomy development and to the Ego-Exo Transfer baseline implementation and analysis. Avijit Dasgupta contributed to the annotation and taxonomy development, and to the early stage of baseline design. Miguel Martin contributed to the annotation and the taxonomy development. Shraman Pramanick contributed the EgoVLPv2 pretrained backbone. Yifei Huang contributed to the early stages of task definition and baseline design. Devansh Kukreja built the render flow to generate frame-aligned videos of each camera for each take as model input, and to produce video collages for annotations. Kristen Grauman contributed to the task formulation.

**Energy-efficient multimodal keystep recognition** Tushar Nagarajan led the energy-efficient multimodal benchmark, co-developed the task formulation, and advised the baseline development. Sagnik Majumder led the baseline development effort and contributed all baseline implementations and analysis of results for the benchmark. Merey Ramazanova developed the energy profiler used to evaluate all baselines and contributed to the experimental analysis. Mitesh Kumar Singh helped design the energy formula. Miao Liu and Shengxin Cindy Zha initiated this benchmark and developed an early version of the task formulation.

**Procedure understanding** Antonino Furnari led the procedure understanding benchmark, and contributed to the task definition, the annotation guidelines, the baseline design, and paper writing. Giovanni Maria Farinella contributed to the task definition, the annotation guidelines, the baseline design, and paper writing. Luigi Seminara contributed to the annotation guidelines, the baseline design, and paper writing; he also developed tools for data annotation, and the baselines for the benchmark. Francesco Ragusa contributed to the annotation guidelines, the baseline design, and paper writing; he also developed tools for data annotation, and the baselines for the benchmark. Kumar Ashutosh contributed to the annotation guidelines, baseline design, and development of data annotation tools. Michael Wray contributed to the task definition, the annotation guidelines, the baseline design, and paper writing. Siddhant Bansal contributed to the task definition, the annotation guidelines, the baseline design, and paper writing. Gene Byrne

contributed to the task definition, the annotation guidelines and the baseline design. Tushar Nagarajan contributed to the task definition, and the baseline design.

**Ego-exo proficiency estimation** Santhosh Kumar Ramakrishnan co-led the proficiency estimation benchmark, co-developed the task formulation, and advised the baseline development. Gedas Bertasius co-led the proficiency estimation benchmark, co-developed the task formulation, and advised the baseline development. Arjun Somayazulu developed the demonstrator proficiency estimation baselines. Abrham Gebreselasie developed the demonstration proficiency estimation baselines. Maria Escobar contributed to the task definition and the baseline design. Eugene Byrne contributed to the task definition and the baseline design. Miguel Martin developed an interface for obtaining proficiency estimation scores from the recruited experts. Suyog Jain contributed to an annotation pipeline for demonstration proficiency estimation. Devansh Kukreja built the render flow to generate frame-aligned videos of each camera for each take as model input. Kristen Grauman contributed to the task formulation.

**Ego pose** Kris Kitani co-led the ego-pose benchmark, and provided directional guidance on the task definition, the annotation methodology, and the baseline development. Jianbo Shi co-led the ego-pose benchmark, and provided directional guidance on automatic 3D hand pose generation and development of ego hand pose baseline methods. Maria Escobar led the ego-pose body baseline development, implemented the IMU-based baseline and contributed to experiment analysis. Cristhian Forigua developed the static pose baseline and contributed to the implementation of the IMU-based baseline. Fu-Jen Chu developed the multi-view annotation UI, the hand pose annotation guidelines, and the data preprocessing code for ego-pose annotation; he also trained and evaluated the HandOccNet baseline. Rawal Khirodkar developed the multi-view triangulation and 3D body keypoint estimation pipeline. Zhengyi Luo contributed to the Kinpoly baseline and to the coordinate transform for Aria head poses. Shan Su led the ego-pose hands baseline development, and contributed to automatic 3D hand pose generation, task definition, and annotation development; she also evaluated the baseline using METRO. Suyog Jain developed the annotation pipeline to

scale the annotation collection, worked on training the annotators and managed the overall annotation process. Miguel Martin contributed to the automatic ground truth generation pipeline and provided high-level coordination of the body and hands automatic ground truth generation. Jinxu Zhang developed automatic ground truth generation for 3D hand pose; he also trained and evaluated baseline model POTTER. Yiming Huang trained and evaluated the baseline model THORnet. Zhifan Zhu developed the METRO hand pose baseline method. Jing Huang led the automatic ground truth generation effort, refined body pose annotation guidelines, coordinated ego-pose body baseline development, and ran the EgoEgo body pose baseline.

## References

- Flavell, J.H., Flavell, E.R., Green, F.L., Wilcox, S.A.: The development of three spatial perspective-taking rules. *Child Development* (1981)
- Newcombe, N.: The development of spatial perspective taking. *Advances in child development and behavior* (1989)
- Sigurdsson, G.A., Gupta, A., Schmid, C., Farhadi, A., Alahari, K.: Charades-ego: A large-scale dataset of paired third and first person videos. arXiv preprint arXiv:1804.09626 (2018)
- Sener, F., Chatterjee, D., Sheleporov, D., He, K., Singhania, D., Wang, R., Yao, A.: Assembly101: A large-scale multi-view video dataset for understanding procedural activities. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 21096–21106 (2022)
- Kwon, T., Tekin, B., Stühmer, J., Bogo, F., Pollefeyns, M.: H2o: Two hands manipulating objects for first person interaction recognition. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 10138–10148 (2021)
- Torre, F.D., Hodgins, J., Montano, J., Valcarcel, S., Forcada, R., Macey, J.: Guide to the carnegie mellon university multimodal activity (cmummact) database. In: Tech. Report CMU-RI-TR-08-22, Robotics Institute, Carnegie Mellon University (2009)
- Rai, N., Chen, H., Ji, J., Desai, R., Kozuka, K., Ishizaka, S., Adeli, E., Niebles, J.C.: Home action genome: Contrastive compositional action understanding. In: *CVPR* (2021)
- Damen, D., Doughty, H., Farinella, G.M., , Furnari, A., Ma, J., Kazakos, E., Moltisanti, D., Munro, J., Perrett, T., Price, W., Wray, M.: Rescaling egocentric vision. *IJCV* (2021)
- Grauman, K., Westbury, A., Byrne, E., Chavis, Z., Furnari, A., Girdhar, R., Hamburger, J., Jiang, H., Liu, M., Liu, X., Martin, M., Nagarajan, T., Radosavovic, I., Ramakrishnan, S.K., Ryan, F., Sharma, J., Wray, M., Xu, M., Xu, E.Z., Zhao, C., Bansal, S., Batra, D., Cartillier, V., Crane, S., Do, T., Doulaty, M., Erapalli, A., Feichtenhofer, C., Fragomeni, A., Fu, Q., Gebreselasie, A., González, C., Hillis, J., Huang, X., Huang, Y., Jia, W., Khoo, W., Kolář, J., Kotur, S., Kumar, A., Landini, F., Li, C., Li, Y., Li, Z., Mangalam, K., Modhugu, R., Munro, J., Murrell, T., Nishiyasu, T., Price, W., Ruiz, P., Ramazanova, M., Sari, L., Somasundaram, K., Southerland, A., Sugano, Y., Tao, R., Vo, M., Wang, Y., Wu, X., Yagi, T., Zhao, Z., Zhu, Y., Arbeláez, P., Crandall, D., Damen, D., Farinella, G.M., Fuegen, C., Ghanem, B., Ithapu, V.K., Jawahar, C.V., Joo, H., Kitani, K., Li, H., Newcombe, R., Oliva, A., Park, H.S., Rehg, J.M., Sato, Y., Shi, J., Shou, M.Z., Torralba, A., Torresani, L., Yan, M., Malik, J.: Ego4D: Around the world in 3,000 hours of egocentric video. In: *CVPR* (2022)
- Kay, W., Carreira, J., Simonyan, K., Zhang, B., Hillier, C., Vijayanarasimhan, S., Viola, F., Green, T., Back, T., Natsev, P., et al.: The kinetics human action video dataset. arXiv preprint arXiv:1705.06950 (2017)
- Gu, C., Sun, C., Ross, D.A., Vondrick, C., Pantofaru, C., Li, Y., Vijayanarasimhan, S., Toderici, G., Ricco, S., Sukthankar, R., Schmid, C., Malik, J.: Ava: A video dataset of spatio-temporally localized atomic visual actions. In: *CVPR* (2018)
- Monfort, M., Andonian, A., Zhou, B., Ramakrishnan, K., Bargal, S.A., Yan, T., Brown, L., Fan, Q., Gutfreund, D., Vondrick, C., Oliva, A.: Moments in time dataset: one million videos for event understanding. *PAMI* (2019)
- Soomro, K., Zamir, A.R., Shah, M.: Ucf101: A dataset of 101 human action classes from videos in the wild. In: *CRCV-TR-12-01* (2012)
- Miech, A., Zhukov, D., Alayrac, J.-B., Tapaswi, M., Laptev, I., Sivic, J.: HowTo100M: Learning a Text-Video Embedding by Watching Hundred Million Narrated Video Clips. In: *ICCV* (2019)

- Tang, Y., Lu, J., Zhou, J.: Comprehensive instructional video analysis: The coin dataset and performance evaluation. *IEEE transactions on pattern analysis and machine intelligence* (2020)
- Zhukov, D., Alayrac, J.-B., Cinbis, R.G., Fouhey, D., Laptev, I., Sivic, J.: Cross-task weakly supervised learning from instructional videos. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2019)
- Zhou, L., Louis, N., Corso, J.: Weakly-supervised video object grounding from text by loss weighting and object interaction. In: *BMVC* (2018)
- Engel, J., Somasundaram, K., Goesele, M., Sun, A., Gamino, A., Turner, A., Talatoff, A., Yuan, A., Souti, B., Meredith, B., Peng, C., Sweeney, C., Wilson, C., Barnes, D., DeTone, D., Caruso, D., Valleroy, D., Ginjupalli, D., Frost, D., Miller, E., Mueggler, E., Oleinik, E., Zhang, F., Somasundaram, G., Solaira, G., Lanaras, H., Howard-Jenkins, H., Tang, H., Kim, H.J., Rivera, J., Luo, J., Dong, J., Straub, J., Bailey, K., Eckenhoff, K., Ma, L., Pesqueira, L., Schwesinger, M., Monge, M., Yang, N., Charron, N., Raina, N., Parkhi, O., Borschowa, P., Moulou, P., Gupta, P., Mur-Artal, R., Pennington, R., Kulkarni, S., Miglani, S., Gondi, S., Solanki, S., Diener, S., Cheng, S., Green, S., Saarinen, S., Patra, S., Mourikis, T., Whelan, T., Singh, T., Balntas, V., Baiyya, V., Dreewes, W., Pan, X., Lou, Y., Zhao, Y., Mansour, Y., Zou, Y., Lv, Z., Wang, Z., Yan, M., Ren, C., Nardi, R.D., Newcombe, R.: Project Aria: A New Tool for Egocentric Multi-Modal AI Research (2023)
- Li, Y., Liu, M., Rehg, J.M.: In the eye of beholder: Joint learning of gaze and actions in first person video. In: *ECCV* (2018)
- Ragusa, F., Furnari, A., Livatino, S., Farinella, G.M.: The meccano dataset: Understanding human-object interactions from egocentric videos in an industrial-like domain. In: *WACV* (2021)
- Damen, D., Doughty, H., Farinella, G.M., Furnari, A., Kazakos, E., Ma, J., Moltisanti, D., Munro, J., Perrett, T., Price, W., *et al.*: Rescaling egocentric vision: collection, pipeline and challenges for epic-kitchens-100. *International Journal of Computer Vision* **130**(1), 33–55 (2022)
- Wang, X., Kwon, T., Rad, M., Pan, B., Chakraborty, I., Andrist, S., Bohus, D., Feniello, A., Tekin, B., Frujeri, F.V., *et al.*: Holoassist: an egocentric human interaction dataset for interactive ai assistants in the real world. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 20270–20281 (2023)
- Weinland, D., Ronfard, R., Boyer, E.: Free viewpoint action recognition using motion history volumes. *Computer Vision and Image Understanding (CVIU)* (2006)
- Corona, K., Osterdahl, K., Collins, R., Hoogs, A.: Meva: A large-scale multiview, multimodal video dataset for activity detection. In: *WACV* (2021)
- Sigurdsson, G.A., Gupta, A., Schmid, C., Farhadi, A., Alahari, K.: Actor and observer: Joint modeling of first and third-person videos. In: *CVPR* (2018)
- Jia, B., Chen, Y., Huang, S., Zhu, Y., Zhu, S.-c.: Lemma: A multi-view dataset for learning multi-agent multi-task activities. In: *European Conference on Computer Vision*, pp. 767–786 (2020). Springer
- Ohkawa, T., He, K., Sener, F., Hodan, T., Tran, L., Keskin, C.: Assemblyhands: Towards egocentric activity understanding via 3d hand pose estimation. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12999–13008 (2023)
- Huang, Y., Chen, G., Xu, J., Zhang, M., Yang, L., Pei, B., Zhang, H., Dong, L., Wang, Y., Wang, L., *et al.*: Egoexolearn: A dataset for bridging asynchronous ego-and exo-centric view of procedural activities in real world. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 22072–22086 (2024)
- Damen, D., Doughty, H., Farinella, G.M., Fidler, S., Furnari, A., Kazakos, E., Moltisanti, D.,

- Munro, J., Perrett, T., Price, W., Wray, M.: Scaling egocentric vision: The epic-kitchens dataset. In: European Conference on Computer Vision (ECCV) (2018)
- Tschernezki, V., Darkhalil, A., Zhu, Z., Fouhey, D., Larina, I., Larlus, D., Damen, D., Vedaldi, A.: EPIC Fields: Marrying 3D Geometry and Video Understanding. In: Proceedings of the Neural Information Processing Systems (NeurIPS) (2023)
- Lee, Y.J., Ghosh, J., Grauman, K.: Discovering important people and objects for egocentric video summarization. In: CVPR (2012)
- Pirsiavash, H., Ramanan, D.: Detecting activities of daily living in first-person camera views. In: CVPR (2012)
- Singh, K.K., Fatahalian, K., Efros, A.A.: Krishnacam: Using a longitudinal, single-person, egocentric dataset for scene understanding tasks. In: WACV (2016)
- Wong, B., Chen, J., Wu, Y., Lei, S.W., Mao, D., Gao, D., Shou, M.Z.: Assistq: Affordance-centric question-driven task completion for egocentric assistant. In: European Conference on Computer Vision (2022)
- Bansal, S., Arora, C., Jawahar, C.V.: My view is the best view: Procedure learning from egocentric videos. In: European Conference on Computer Vision (ECCV) (2022)
- Chang, A., Dai, A., Funkhouser, T., Nießner, M., Savva, M., Song, S., Zeng, A., Zhang, Y.: Matterport3d: Learning from rgb-d data in indoor environments. In: Proceedings of the International Conference on 3D Vision (3DV) (2017). MatterPort3D dataset license available at: [http://kaldir.vc.in.tum.de/matterport/MP\\_TOS.pdf](http://kaldir.vc.in.tum.de/matterport/MP_TOS.pdf).
- Xia, F., R. Zamir, A., He, Z.-Y., Sax, A., Malik, J., Savarese, S.: Gibson Env: real-world perception for embodied agents. In: CVPR (2018). IEEE. Gibson license is available at [http://svl.stanford.edu/gibson2/assets/GDS\\_agreement.pdf](http://svl.stanford.edu/gibson2/assets/GDS_agreement.pdf).
- Straub, J., Whelan, T., Ma, L., Chen, Y., Wijmans, E., Green, S., Engel, J.J., Mur-Artal, R., Ren, C., Verma, S., Clarkson, A., Yan, M., Budge, B., Yan, Y., Pan, X., Yon, J., Zou, Y., Leon, K., Carter, N., Briales, J., Gilligham, T., Mueggler, E., Pesqueira, L., Savva, M., Batra, D., Strasdat, H.M., Nardi, R.D., Goesele, M., Lovegrove, S., Newcombe, R.: The Replica dataset: A digital replica of indoor spaces. arXiv preprint arXiv:1906.05797 (2019)
- Ramakrishnan, S.K., Gokaslan, A., Wijmans, E., Maksymets, O., Clegg, A., Turner, J.M., Undersander, E., Galuba, W., Westbury, A., Chang, A.X., Savva, M., Zhao, Y., Batra, D.: Habitat-matterport 3d dataset (HM3d): 1000 large-scale 3d environments for embodied AI. In: Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (2021). <https://arxiv.org/abs/2109.08238>
- Xiao, J., Owens, A., Torralba, A.: Sun3d: A database of big spaces reconstructed using sfm and object labels. In: ICCV (2013)
- Reizenstein, J., Shapovalov, R., Henzler, P., Sborodone, L., Labatut, P., Novotny, D.: Common objects in 3d: Large-scale learning and evaluation of real-life 3d category reconstruction. In: ICCV (2021)
- Wu, Z., Song, S., Khosla, A., Yu, F., Zhang, L., Tang, X., Xiao, J.: 3d shapenets: A deep representation for volumetric shapes. In: Computer Vision and Pattern Recognition, IEEE Conference On (2015)
- Zhang, S., Ma, Q., Zhang, Y., Qian, Z., Kwon, T., Pollefeys, M., Bogo, F., Tang, S.: Egobody: Human body shape and motion of interacting people from head-mounted devices. In: ECCV (2022)
- Li, J., Liu, K., Wu, J.: Ego-body pose estimation via ego-head pose estimation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 17142–17151 (2023)
- Joo, H., Simon, T., Li, X., Liu, H., Tan, L., Gui, L., Banerjee, S., Godisart, T.S., Nabbe, B., Matthews, I., Kanade, T., Nobuhara, S.,

- Sheikh, Y.: Panoptic studio: A massively multi-view system for social interaction capture. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2017)
- Khirodkar, R., Bansal, A., Ma, L., Newcombe, R., Vo, M., Kitani, K.: Egohumans: An egocentric 3d multi-human benchmark. In: *ICCV* (2023)
- Guzov, V., Mir, A., Sattler, T., Pons-Moll, G.: Human poseitioning system (hps): 3d human pose estimation and self-localization in large scenes from body-mounted sensors. In: *CVPR* (2021)
- Pirsiavash, H., Vondrick, C., Torralba, A.: Assessing the quality of actions. In: *ECCV* (2014)
- Bertasius, G., Park, H.S., Yu, S., Shi, J.: Am i a baller? basketball performance assessment from first-person videos. In: *ICCV* (2017)
- Parmar, P., Morris, B.: Action quality assessment across multiple actions. In: *WACV* (2019)
- Doughty, H., Damen, D., Mayol-Cuevas, W.: Who's better? who's best? pairwise deep ranking for skill determination. In: *CVPR* (2018)
- Doughty, H., Mayol-Cuevas, W., Damen, D.: The Pros and Cons: Rank-aware Temporal Attention for Skill Determination in Long Videos (2019)
- Zhang, S., Dai, W., Wang, S., Shen, X., Lu, J., Zhou, J., Tang, Y.: Logo: A long-form video dataset for group action quality assessment. In: *CVPR* (2023)
- Ben-Shabat, Y., Yu, X., Saleh, F., Campbell, D., Rodriguez-Opazo, C., Li, H., Gould, S.: The ikea asm dataset: Understanding people assembling furniture through actions, objects and pose (2020)
- Elhamifar, E., Huynh, D.: Self-supervised multi-task procedure learning from instructional videos. In: *European Conference on Computer Vision*, pp. 557–573 (2020). Springer
- Xu, H., Ghosh, G., Huang, P.-Y., Okhonko, D., Aghajanyan, A., Metze, F., Zettlemoyer, L., Feichtenhofer, C.: Videoclip: Contrastive pre-training for zero-shot video-text understanding. *arXiv preprint arXiv:2109.14084* (2021)
- Miech, A., Alayrac, J.-B., Smaira, L., Laptev, I., Sivic, J., Zisserman, A.: End-to-end learning of visual representations from uncurated instructional videos. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9879–9889 (2020)
- Dvornik, N., Hadji, I., Pham, H., Bhatt, D., Martinez, B., Fazly, A., Jepson, A.D.: Flow graph to video grounding for weakly-supervised multi-step localization. In: *ECCV*, pp. 319–335 (2022). Springer
- Lin, X., Petroni, F., Bertasius, G., Rohrbach, M., Chang, S.-F., Torresani, L.: Learning to recognize procedural activities with distant supervision. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 13853–13863 (2022)
- Chang, C.-Y., Huang, D.-A., Xu, D., Adeli, E., Fei-Fei, L., Niebles, J.C.: Procedure planning in instructional videos. In: *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XI*, pp. 334–350 (2020). Springer
- Bi, J., Luo, J., Xu, C.: Procedure planning in instructional videos via contextual modeling and model-based policy learning. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 15611–15620 (2021)
- Zhao, H., Hadji, I., Dvornik, N., Derpanis, K.G., Wildes, R.P., Jepson, A.D.: P3iv: Probabilistic procedure planning from instructional videos with weak supervision. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2938–2948 (2022)
- Wang, H., Wu, Y., Guo, S., Wang, L.: Pdpp: Projected diffusion for procedure planning in instructional videos. *arXiv preprint arXiv:2303.14676* (2023)
- Zhong, Y., Yu, L., Bai, Y., Li, S., Yan, X., Li, Y.: Learning procedure-aware video representation from instructional videos and their narrations.

arXiv preprint arXiv:2303.17839 (2023)

Shvetsova, N., Chen, B., Rouditchenko, A., Thomas, S., Kingsbury, B., Feris, R.S., Harwath, D., Glass, J., Kuehne, H.: Everything at once-multi-modal fusion transformer for video retrieval. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 20020–20029 (2022)

Ko, D., Choi, J., Ko, J., Noh, S., On, K.-W., Kim, E.-S., Kim, H.J.: Video-text representation learning via differentiable weak temporal alignment. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 5016–5025 (2022)

Cao, M., Yang, T., Weng, J., Zhang, C., Wang, J., Zou, Y.: Locvtp: Video-text pre-training for temporal localization. In: European Conference on Computer Vision (2022)

Narasimhan, M., Yu, L., Bell, S., Zhang, N., Darrell, T.: Learning and verification of task structure in instructional videos. arXiv preprint arXiv:2303.13519 (2023)

Zhou, L., Xu, C., Corso, J.J.: Towards automatic learning of procedures from web instructional videos. In: AAAI (2018)

Alayrac, J.-B., Bojanowski, P., Agrawal, N., Sivic, J., Laptev, I., Lacoste-Julien, S.: Unsupervised learning from narrated instruction videos. In: CVPR (2016)

Ashutosh, K., Ramakrishnan, S.K., Afouras, T., Grauman, K.: Video-mined task graphs for keystep recognition in instructional videos. In: NeurIPS (2023)

Zhou, H., Martin-Martin, R., Kapadia, M., Savarese, S., Niebles, J.C.: Procedure-aware pretraining for instructional video understanding. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2023)

Ardeshir, S., Borji, A.: Ego2top: Matching viewers in egocentric and top-view videos. In: ECCV (2016)

Ardeshir, S., Borji, A.: Egocentric meets top-view. IEEE transactions on pattern analysis and machine intelligence **41**(6) (2018)

Fan, C., Lee, J., Xu, M., Kumar Singh, K., Jae Lee, Y., Crandall, D.J., Ryoo, M.S.: Identifying first-person camera wearers in third-person videos. In: CVPR (2017)

Xu, M., Fan, C., Wang, Y., Ryoo, M.S., Crandall, D.J.: Joint person segmentation and identification in synchronized first-and third-person videos. In: ECCV (2018)

Wen, Y., Singh, K.K., Anderson, M., Jan, W.-P., Lee, Y.J.: Seeing the unseen: Predicting the first-person camera wearer’s location and pose in third-person scenes. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) Workshops, pp. 3446–3455 (2021)

Ardeshir, S., Borji, A.: An exocentric look at egocentric actions and vice versa. Computer Vision and Image Understanding **171** (2018)

Sermanet, P., Lynch, C., Chebotar, Y., Hsu, J., Jang, E., Schaal, S., Levine, S.: Time-contrastive networks: Self-supervised learning from video. Proceedings of International Conference in Robotics and Automation (ICRA) (2018)

Yu, H., Cai, M., Liu, Y., Lu, F.: What i see is what you see: Joint attention learning for first and third person video co-analysis. In: ACM MM (2019)

Yu, H., Cai, M., Liu, Y., Lu, F.: First-and third-person video co-analysis by learning spatial-temporal joint attention. IEEE Transactions on Pattern Analysis and Machine Intelligence (2020)

Xue, Z., Grauman, K.: Learning fine-grained view-invariant representations from unpaired ego-exo videos via temporal alignment. In: NeurIPS (2023)

Li, Y., Nagarajan, T., Xiong, B., Grauman, K.: Ego-exo: Transferring visual representations from third-person to first-person videos. In:

- Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 6943–6953 (2021)
- Regmi, K., Borji, A.: Cross-view image synthesis using conditional gans. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2018)
- Regmi, K., Borji, A.: Cross-view image synthesis using geometry-guided conditional gans. Computer Vision and Image Understanding (2019) <https://doi.org/10.1016/j.cviu.2019.07.008>
- Tang, H., Xu, D., Sebe, N., Wang, Y., Corso, J.J., Yan, Y.: Multi-channel attention selection gan with cascaded semantic guidance for cross-view image translation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 2417–2426 (2019)
- Ren, B., Tang, H., Sebe, N.: Cascaded cross mlp-mixer gans for cross-view image translation. arXiv preprint arXiv:2110.10183 (2021)
- Luo, M., Xue, Z., Dimakis, A., Grauman, K.: Put myself in your shoes: Lifting the egocentric perspective from exocentric videos. In: ECCV (2024)
- Cheng, F., Luo, M., Wang, H., Dimakis, A., Torresani, L., Bertasius, G., Grauman, K.: 4DIFF: 3d-aware diffusion model for third-to-first viewpoint translation. In: ECCV (2024)
- Liu, A., Tucker, R., Jampani, V., Makadia, A., Snavely, N., Kanazawa, A.: Infinite nature: Perpetual view generation of natural scenes from a single image. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 14458–14467 (2021)
- Ren, X., Wang, X.: Look outside the room: Synthesizing a consistent long-term 3d scene video from a single image. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 3563–3573 (2022)
- Rombach, R., Esser, P., Ommer, B.: Geometry-free view synthesis: Transformers and no 3d priors. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 14356–14366 (2021)
- Wiles, O., Gkioxari, G., Szeliski, R., Johnson, J.: Synsin: End-to-end view synthesis from a single image. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 7467–7477 (2020)
- Watson, D., Chan, W., Martin-Brualla, R., Ho, J., Tagliasacchi, A., Norouzi, M.: Novel view synthesis with diffusion models. arXiv preprint arXiv:2210.04628 (2022)
- Tseng, H.-Y., Li, Q., Kim, C., Alsisan, S., Huang, J.-B., Kopf, J.: Consistent view synthesis with pose-guided diffusion models. arXiv preprint arXiv:2303.17598 (2023)
- Chan, E.R., Nagano, K., Chan, M.A., Bergman, A.W., Park, J.J., Levy, A., Aittala, M., De Mello, S., Karras, T., Wetzstein, G.: Generative novel view synthesis with 3d-aware diffusion models. arXiv preprint arXiv:2304.02602 (2023)
- Regmi, K., Shah, M.: Bridging the domain gap for ground-to-aerial image matching. In: ICCV (2019)
- Lin, T., Cui, Y., Belongie, S., Hays, J.: Learning deep representations for ground-to-aerial geolocation. In: CVPR (2015)
- Kukelova, Z., Heller, J., Fitzgibbon, A.: Efficient intersection of three quadrics and applications in computer vision. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 1799–1808 (2016)
- Lin, K.Q., Wang, A.J., Soldan, M., Wray, M., Yan, R., Xu, E.Z., Gao, D., Tu, R., Zhao, W., Kong, W., et al.: Egocentric video-language pretraining. NeurIPS (2022)
- Pramanick, S., Song, Y., Nag, S., Lin, K.Q., Shah, H., Shou, M.Z., Chellappa, R., Zhang, P.: Egovlpv2: Egocentric video-language pre-training with fusion in the backbone. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 5285–5297 (2023)

- Pan, B., Cai, H., Huang, D.-A., Lee, K.-H., Gaidon, A., Adeli, E., Niebles, J.C.: Spatio-temporal graph for video captioning with knowledge distillation. In: CVPR (2020)
- Iashin, V., Rahtu, E.: Multi-modal dense video captioning. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops, pp. 958–959 (2020)
- Zhao, Y., Misra, I., Krähenbühl, P., Girdhar, R.: Learning video representations from large language models. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 6586–6597 (2023)
- Radford, A., Kim, J.W., Xu, T., Brockman, G., McLeavey, C., Sutskever, I.: Robust speech recognition via large-scale weak supervision. In: International Conference on Machine Learning, pp. 28492–28518 (2023). PMLR
- Ashutosh, K., Girdhar, R., Torresani, L., Grauman, K.: Hiervl: Learning hierarchical video-language embeddings. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2023)
- Jiang, W., Trulls, E., Hosang, J., Tagliasacchi, A., Yi, K.M.: Cotr: Correspondence transformer for matching across images. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 6207–6217 (2021)
- Vicente, S., Rother, C., Kolmogorov, V.: Object cosegmentation. In: CVPR 2011, pp. 2217–2224 (2011). IEEE
- Tang, H., Liang, K., Grauman, K., Feiszli, M., Wang, W.: Egotracks: A long-term egocentric visual object tracking dataset. Advances in Neural Information Processing Systems (2023)
- Perazzi, F., Pont-Tuset, J., McWilliams, B., Van Gool, L., Gross, M., Sorkine-Hornung, A.: A benchmark dataset and evaluation methodology for video object segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 724–732 (2016)
- Brodersen, K.H., Ong, C.S., Stephan, K.E., Buhmann, J.M.: The balanced accuracy and its posterior distribution. In: 2010 20th International Conference on Pattern Recognition, pp. 3121–3124 (2010). IEEE
- Shen, X., Efros, A.A., Joulin, A., Aubry, M.: Learning co-segmentation by segment swapping for retrieval and discovery. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 5082–5092 (2022)
- Cheng, H.K., Schwing, A.G.: Xmem: Long-term video object segmentation with an atkinson-shiffrin memory model. In: European Conference on Computer Vision, pp. 640–658 (2022). Springer
- Lu, X., Li, Z., Cui, Z., Oswald, M.R., Pollefeys, M., Qin, R.: Geometry-aware satellite-to-ground image synthesis for urban areas. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2020)
- Liu, G., Tang, H., Latapie, H., Yan, Y.: Exocentric to egocentric image generation via parallel generative adversarial network. In: ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 1843–1847 (2020). IEEE
- Luo, M., Xue, Z., Dimakis, A., Grauman, K.: Put myself in your shoes: Lifting the egocentric perspective from exocentric videos. arXiv:2403.06351 (2024)
- Liu, G., Tang, H., Latapie, H.M., Corso, J.J., Yan, Y.: Cross-view exocentric to egocentric video synthesis. In: Proceedings of the 29th ACM International Conference on Multimedia, pp. 974–982 (2021)
- Hore, A., Ziou, D.: Image quality metrics: Psnr vs. ssim. In: 2010 20th International Conference on Pattern Recognition, pp. 2366–2369 (2010). IEEE
- Ding, K., Ma, K., Wang, S., Simoncelli, E.P.: Image quality assessment: Unifying structure and texture similarity. IEEE transactions on pattern analysis and machine intelligence **44**(5),

2567–2581 (2020)

Zhang, R., Isola, P., Efros, A.A., Shechtman, E., Wang, O.: The unreasonable effectiveness of deep features as a perceptual metric. In: CVPR (2018)

Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: International Conference on Machine Learning, pp. 8748–8763 (2021). PMLR

Isola, P., Zhu, J.-Y., Zhou, T., Efros, A.A.: Image-to-image translation with conditional adversarial networks. CVPR (2017)

Varma, M., Wang, P., Chen, X., Chen, T., Venugopalan, S., Wang, Z.: Is attention all that neRF needs? In: The Eleventh International Conference on Learning Representations (2023). <https://openreview.net/forum?id=xELtsE-xx>

Peebles, W., Xie, S.: Scalable diffusion models with transformers. arXiv preprint arXiv:2212.09748 (2022)

Reed, S., Lee, H., Anguelov, D., Szegedy, C., Erhan, D., Rabinovich, A.: Training deep neural networks on noisy labels with bootstrapping. arXiv preprint arXiv:1412.6596 (2014)

Sudre, C.H., Li, W., Vercauteren, T., Ourselin, S., Jorge Cardoso, M.: Generalised dice overlap as a deep learning loss function for highly unbalanced segmentations. In: Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support: Third International Workshop, DLMIA 2017, and 7th International Workshop, ML-CDS 2017, Held in Conjunction with MICCAI 2017, Québec City, QC, Canada, September 14, Proceedings 3, pp. 240–248 (2017). Springer

Perez, E., Strub, F., De Vries, H., Dumoulin, V., Courville, A.: Film: Visual reasoning with a general conditioning layer. In: Proceedings of the AAAI Conference on Artificial Intelligence, vol. 32 (2018)

Bertasius, G., Wang, H., Torresani, L.: Is space-time attention all you need for video understanding? In: ICML (2021)

Goyal, R., Ebrahimi Kahou, S., Michalski, V., Materzynska, J., Westphal, S., Kim, H., Haenel, V., Fruend, I., Yianilos, P., Mueller-Freitag, M., et al.: The “something something” video database for learning and evaluating visual common sense. In: Proceedings of the IEEE International Conference on Computer Vision, pp. 5842–5850 (2017)

Song, Y., Byrne, E., Nagarajan, T., Wang, H., Martin, M., Torresani, L.: Ego4d goal-step: Toward hierarchical understanding of procedural activities. In: NeurIPS (2023)

Mavroudi, E., Afouras, T., Torresani, L.: Learning to ground instructional articles in videos through narrations. (2022)

Oord, A.v.d., Li, Y., Vinyals, O.: Representation learning with contrastive predictive coding. arXiv preprint arXiv:1807.03748 (2018)

Sermanet, P., Lynch, C., Hsu, J., Levine, S.: Time-contrastive networks: Self-supervised learning from multi-view observation. In: 2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW), pp. 486–487 (2017). IEEE

Hinton, G., Vinyals, O., Dean, J.: Distilling the knowledge in a neural network. arXiv preprint arXiv:1503.02531 (2015)

Tong, Z., Song, Y., Wang, J., Wang, L.: Videomae: Masked autoencoders are data-efficient learners for self-supervised video pre-training. (2022). <https://arxiv.org/abs/2203.12602>

Feichtenhofer, C.: X3d: Expanding architectures for efficient video recognition. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 203–213 (2020)

Vasu, P.K.A., Gabriel, J., Zhu, J., Tuzel, O., Ranjan, A.: Mobileone: An improved one millisecond mobile backbone. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 7907–7917 (2023)

- Howard, A.G., Zhu, M., Chen, B., Kalenichenko, D., Wang, W., Weyand, T., Andreetto, M., Adam, H.: Mobilenets: Efficient convolutional neural networks for mobile vision applications. arXiv preprint arXiv:1704.04861 (2017)
- Zhang, X., Zhou, X., Lin, M., Sun, J.: Shufflenet: An extremely efficient convolutional neural network for mobile devices. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 6848–6856 (2018)
- Tan, M., Chen, B., Pang, R., Vasudevan, V., Sandler, M., Howard, A., Le, Q.V.: Mnasnet: Platform-aware neural architecture search for mobile. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 2820–2828 (2019)
- Mehta, S., Rastegari, M.: Mobilevit: light-weight, general-purpose, and mobile-friendly vision transformer. arXiv preprint arXiv:2110.02178 (2021)
- Gao, R., Oh, T.-H., Grauman, K., Torresani, L.: Listen to look: Action recognition by previewing audio. In: CVPR (2020)
- Korbar, B., Tran, D., Torresani, L.: Scsampl: Sampling salient clips from video for efficient action recognition. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (2019)
- Ghodrati, A., Bejnordi, B.E., Habibian, A.: Frameexit: Conditional early exiting for efficient video recognition. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2021)
- Meng, Y., Lin, C.-C., Panda, R., Sattigeri, P., Karlinsky, L., Oliva, A., Saenko, K., Feris, R.: Ar-net: Adaptive frame resolution for efficient action recognition. In: ECCV 2020 (2020)
- Tan, S., Nagarajan, T., Grauman, K.: Egodistill: Egocentric head motion distillation for efficient video understanding. NeurIPS (2023)
- Iandola, F.N., Han, S., Moskewicz, M.W., Ashraf, K., Dally, W.J., Keutzer, K.: Squeezezenet: Alexnet-level accuracy with 50x fewer parameters and 0.5 mb model size. arXiv preprint arXiv:1602.07360 (2016)
- Esser, S.K., McKinstry, J.L., Bablani, D., Appuswamy, R., Modha, D.S.: Learned step size quantization. arXiv preprint arXiv:1902.08153 (2019)
- Polino, A., Pascanu, R., Alistarh, D.: Model compression via distillation and quantization. arXiv preprint arXiv:1802.05668 (2018)
- Zhu, M., Gupta, S.: To prune, or not to prune: exploring the efficacy of pruning for model compression. arXiv preprint arXiv:1710.01878 (2017)
- Wu, Z., Nagarajan, T., Kumar, A., Rennie, S., Davis, L.S., Grauman, K., Feris, R.: Blockdrop: Dynamic inference paths in residual networks. In: CVPR (2018)
- Abrash, M.: Creating the future: Augmented reality, the next human-machine interface. In: 2021 IEEE International Electron Devices Meeting (IEDM) (2021)
- Chen, Y.-H., Yang, T.-J., Emer, J., Sze, V.: Eye-riss v2: A flexible accelerator for emerging deep neural networks on mobile devices. IEEE Journal on Emerging and Selected Topics in Circuits and Systems (2019)
- Yang, L., Radway, R.M., Chen, Y.-H., Wu, T.F., Liu, H., Ansari, E., Chandra, V., Mitra, S., Beigné, E.: Three-dimensional stacked neural network accelerator architectures for ar/vr applications. IEEE Micro (2022)
- Sze, V., Chen, Y.-H., Yang, T.-J., Emer, J.S.: How to evaluate deep neural network processors: Tops/w (alone) considered harmful. IEEE Solid-State Circuits Magazine (2020)
- Desislavov, R., Martínez-Plumed, F., Hernández-Orallo, J.: Trends in ai inference energy consumption: Beyond the performance-vs-parameter laws of deep learning. Sustainable Computing: Informatics and Systems **38**, 100857 (2023)

- Horowitz, M.: 1.1 computing's energy problem (and what we can do about it). In: 2014 IEEE International Solid-state Circuits Conference Digest of Technical Papers (ISSCC) (2014)
- Liu, C., Bainbridge, L., Berkovich, A., Chen, S., Gao, W., Tsai, T.-H., Mori, K., Ikeno, R., Uno, M., Isozaki, T., et al.: A  $4.6\ \mu\text{m}$ ,  $512 \times 512$ , ultra-low power stacked digital pixel sensor with triple quantization and 127db dynamic range. In: 2020 IEEE International Electron Devices Meeting (IEDM) (2020)
- De Geest, R., Gavves, E., Ghodrati, A., Li, Z., Snoek, C., Tuytelaars, T.: Online action detection. In: Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part V 14, pp. 269–284 (2016). Springer
- Liao, J., Duan, H., Feng, K., Zhao, W., Yang, Y., Chen, L.: A light weight model for active speaker detection. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 22932–22941 (2023)
- Chang, C.-Y., Huang, D.-A., Xu, D., Adeli, E., Fei-Fei, L., Niebles, J.C.: Procedure planning in instructional videos. In: European Conference on Computer Vision, pp. 334–350 (2020). Springer
- Bi, J., Luo, J., Xu, C.: Procedure planning in instructional videos via contextual modeling and model-based policy learning. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 15611–15620 (2021)
- Zhou, H., Martín-Martín, R., Kapadia, M., Savarese, S., Niebles, J.C.: Procedure-aware pretraining for instructional video understanding. In: CVPR, pp. 10727–10738 (2023)
- Seminara, L., Farinella, G.M., Furnari, A.: Differentiable Task Graph Learning: Procedural Activity Representation and Online Mistake Detection from Egocentric Videos (2024)
- Jang, Y., Sohn, S., Logeswaran, L., Luo, T., Lee, M., Lee, H.: Multimodal subtask graph generation from instructional videos. arXiv preprint arXiv:2302.08672 (2023)
- Xu, F.F., Ji, L., Shi, B., Du, J., Neubig, G., Bisk, Y., Duan, N.: A benchmark for structured procedural knowledge extraction from cooking videos. arXiv preprint arXiv:2005.00706 (2020)
- Soran, B., Farhadi, A., Shapiro, L.: Generating notifications for missing actions: Don't forget to turn the lights off! In: ICCV, pp. 4669–4677 (2015)
- Ding, G., Sener, F., Ma, S., Yao, A.: Every mistake counts in assembly. arXiv preprint arXiv:2307.16453 (2023)
- Parmar, P., Morris, B.T.: Learning To Score Olympic Events (2017)
- Parmar, P., Tran Morris, B.: What and how well you performed? a multitask learning approach to action quality assessment. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 304–313 (2019)
- Ismail Fawaz, H., Forestier, G., Weber, J., Idoumghar, L., Muller, P.-A.: Evaluating surgical skills from kinematic data using convolutional neural networks. In: Medical Image Computing and Computer Assisted Intervention – MICCAI 2018, pp. 214–221 (2018)
- Liu, D., Li, Q., Jiang, T., Wang, Y., Miao, R., Shan, F., Li, Z.: Towards unified surgical skill assessment. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 9522–9531 (2021)
- Zhang, Q., Li, B.: Relative hidden markov models for evaluating motion skill. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2013)
- Zia, A., Sharma, Y., Bettadapura, V., Sarin, E.L., Essa, I.A.: Video and accelerometer-based motion analysis for automated surgical skills assessment. CoRR [abs/1702.07772](https://arxiv.org/abs/1702.07772) (2017) [1702.07772](https://arxiv.org/abs/1702.07772)
- Yu, X., Rao, Y., Zhao, W., Lu, J., Zhou, J.: Group-aware contrastive regression for action quality assessment. In: Proceedings of the IEEE/CVF

- International Conference on Computer Vision (ICCV), pp. 7919–7928 (2021)
- Zhang, C.-L., Wu, J., Li, Y.: Actionformer: Localizing moments of actions with transformers. In: European Conference on Computer Vision. LNCS, vol. 13664, pp. 492–510 (2022)
- Girdhar, R., Singh, M., Ravi, N., Maaten, L., Joulin, A., Misra, I.: Omnivore: A single model for many visual modalities. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 16102–16112 (2022)
- Lin, T.-Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft COCO: Common objects in context. In: ECCV (2014)
- Jiang, H., Grauman, K.: Seeing invisible poses: Estimating 3d body pose from egocentric video. In: CVPR (2017)
- Yuan, Y., Kitani, K.: 3d ego-pose estimation via imitation learning. In: Proceedings of the European Conference on Computer Vision (ECCV) (2018)
- Yuan, Y., Kitani, K.: Ego-pose estimation and forecasting as real-time pd control. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) (2019)
- Luo, Z., Hachiuma, R., Yuan, Y., Kitani, K.: Dynamics-regulated kinematic policy for egocentric pose estimation. In: Advances in Neural Information Processing Systems (2021)
- Rhodin, H., Richardt, C., Casas, D., Insafutdinov, E., Shafiei, M., Seidel, H.-P., Schiele, B., Theobalt, C.: Egocap: egocentric marker-less motion capture with two fisheye cameras. ACM Transactions on Graphics (TOG) **35**(6), 1–11 (2016)
- Tome, D., Peluse, P., Agapito, L., Badino, H.: xr-egopose: Egocentric 3d human pose from an hmd camera. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) (2019)
- Xu, W., Chatterjee, A., Zollhoefer, M., Rhodin, H., Fua, P., Seidel, H.-P., Theobalt, C.: Mo 2 cap 2: Real-time mobile 3d motion capture with a cap-mounted fisheye camera. IEEE transactions on visualization and computer graphics **25**(5), 2093–2101 (2019)
- Ahuja, K., Harrison, C., Goel, M., Xiao, R.: Mecap: Whole-body digitization for low-cost vr/ar headsets. In: Proceedings of the 32nd Annual ACM Symposium on User Interface Software and Technology, pp. 453–462 (2019)
- Hwang, D.-H., Aso, K., Yuan, Y., Kitani, K., Koike, H.: Monoeye: Multimodal human motion capture system using a single ultra-wide fish-eye camera. In: Proceedings of the 33rd Annual ACM Symposium on User Interface Software and Technology, pp. 98–111 (2020)
- Simon, T., Joo, H., Matthews, I., Sheikh, Y.: Hand keypoint detection in single images using multiview bootstrapping. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 1145–1153 (2017)
- Moon, G., Yu, S.-I., Wen, H., Shiratori, T., Lee, K.M.: Interhand2. 6m: A dataset and baseline for 3d interacting hand pose estimation from a single rgb image. In: Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XX 16, pp. 548–564 (2020). Springer
- Hampali, S., Rad, M., Oberweger, M., Lepetit, V.: Honnote: A method for 3d annotation of hand and object poses. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 3196–3206 (2020)
- MMPoseContributors: OpenMMLab Pose Estimation Toolbox and Benchmark. <https://github.com/open-mmlab/mmpose> (2020)
- Teed, Z., Deng, J.: DROID-SLAM: Deep Visual SLAM for Monocular, Stereo, and RGB-D Cameras. Advances in neural information processing systems (2021)
- Ho, J., Jain, A., Abbeel, P.: Denoising diffusion probabilistic models. Advances in Neural Information Processing Systems **33**, 6840–6851

(2020)

- Mahmood, N., Ghorbani, N., F. Troje, N., Pons-Moll, G., Black, M.J.: Amass: Archive of motion capture as surface shapes. In: The IEEE International Conference on Computer Vision (ICCV) (2019). <https://amass.is.tue.mpg.de>
- Castillo, A., Escobar, M., Jeanneret, G., Pumarola, A., Arbeláez, P., Thabet, A., Sanakoyeu, A.: Bodiffusion: Diffusing sparse observations for full-body human motion synthesis. CV4Metaverse workshop, International Conference on Computer Vision (2023)
- Jiang, J., Strelci, P., Qiu, H., Fender, A., Laich, L., Snape, P., Holz, C.: Avatarposer: Articulated full-body pose tracking from sparse motion sensing. In: European Conference on Computer Vision, pp. 443–460 (2022). Springer
- Aboukhadra, A.T., Malik, J., Elhayek, A., Rober-tini, N., Stricker, D.: Thor-net: End-to-end graformer-based realistic two hands and object reconstruction with self-supervision. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision, pp. 1001–1010 (2023)
- Zhao, W., Wang, W., Tian, Y.: Graformer: Graph-oriented transformer for 3d pose estimation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 20438–20447 (2022)
- Park, J., Oh, Y., Moon, G., Choi, H., Lee, K.M.: Handoccnet: Occlusion-robust 3d hand mesh estimation network. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 1496–1505 (2022)
- He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: CVPR (2016)
- Lin, T.-Y., Dollár, P., Girshick, R., He, K., Hariharan, B., Belongie, S.: Feature pyramid networks for object detection. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 2117–2125 (2017)
- Romero, J., Tzionas, D., Black, M.J.: Embodied hands: Modeling and capturing hands and bodies together. ACM Transactions on Graphics, (Proc. SIGGRAPH Asia) **36**(6) (2017)
- Zheng, C., Liu, X., Qi, G.-J., Chen, C.: Potter: Pooling attention transformer for efficient human mesh recovery. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 1611–1620 (2023)
- Li, J., Xu, C., Chen, Z., Bian, S., Yang, L., Lu, C.: Hybrik: A hybrid analytical-neural inverse kinematics solution for 3d human pose and shape estimation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 3383–3393 (2021)
- Lin, K., Wang, L., Liu, Z.: End-to-end human pose and mesh reconstruction with transformers. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 1954–1963 (2021)
- Kolotouros, N., Pavlakos, G., Black, M.J., Daniilidis, K.: Learning to reconstruct 3d human pose and shape via model-fitting in the loop. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 2252–2261 (2019)
- Tendulkar, P., Surís, D., Vondrick, C.: Flex: Full-body grasping without full-body grasps. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2023)
- Taheri, O., Ghorbani, N., Black, M.J., Tzionas, D.: GRAB: A dataset of whole-body human grasping of objects. In: European Conference on Computer Vision (ECCV) (2020)
- Plizzari, C., Goletto, G., Furnari, A., Bansal, S., Ragusa, F., Farinella, G.M., Damen, D., Tommasi, T.: An outlook into the future of egocentric vision. International Journal of Computer Vision (2024)
- Ashutosh, K., Nagarajan, T., Pavlakos, G., Kitani, K., Grauman, K.: ExpertAF: Expert Actionable Feedback from Video (2024)

- Kirillov, A., Mintun, E., Ravi, N., Mao, H., Rolland, C., Gustafson, L., Xiao, T., Whitehead, S., Berg, A.C., Lo, W.-Y., Dollar, P., Girshick, R.: Segment anything. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV), pp. 4015–4026 (2023)
- Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014)
- Povey, D., Ghoshal, A., Boulianne, G., Burget, L., Glembek, O., Goel, N.K., Hannemann, M., Motlíček, P., Qian, Y., Schwarz, P., Silovsky, J., Stemmer, G., Veselý, K.: The kaldi speech recognition toolkit. (2011). <https://api.semanticscholar.org/CorpusID:1774023>
- Loshchilov, I., Hutter, F.: Decoupled weight decay regularization. In: International Conference on Learning Representations (2019). <https://openreview.net/forum?id=Bkg6RiCqY7>
- Xu, M., Xiong, Y., Chen, H., Li, X., Xia, W., Tu, Z., Soatto, S.: Long short-term transformer for online action detection. Advances in Neural Information Processing Systems **34**, 1086–1099 (2021)
- Zhao, Y., Krähenbühl, P.: Real-time online video detection with temporal smoothing transformers. In: European Conference on Computer Vision (2022)
- Kwak, I., Guo, J.-Z., Hantman, A., Kriegman, D., Branson, K.: Detecting the starting frame of actions in video. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision, pp. 489–497 (2020)
- Chen, F., Ding, L., Lertniphonphan, K., Li, J., Huang, K., Wang, Z.: Pcie\_egohandpose solution for egoexo4d hand pose challenge. arXiv preprint arXiv:2406.12219 (2024)
- Li, J., Bian, S., Zeng, A., Wang, C., Pang, B., Liu, W., Lu, C.: Human pose regression with residual log-likelihood estimation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 11025–11034 (2021)
- Dosovitskiy, A.: An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929 (2020)
- Pavlakos, G., Shan, D., Radosavovic, I., Kanazawa, A., Fouhey, D., Malik, J.: Reconstructing hands in 3d with transformers. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 9826–9836 (2024)
- Zhang, C., Wu, J., Li, Y.: Actionformer: Localizing moments of actions with transformers. arXiv preprint arXiv:2202.07925 (2022)