

# Ego4D: Around the World in 3,000 Hours of Egocentric Video

Kristen Grauman<sup>1,2</sup>, Andrew Westbury<sup>1</sup>, Eugene Byrne<sup>\*1</sup>, Zachary Chavis<sup>\*3</sup>, Antonino Furnari<sup>\*4</sup>, Rohit Girdhar<sup>\*1</sup>, Jackson Hamburger<sup>\*1</sup>, Hao Jiang<sup>\*5</sup>, Miao Liu<sup>\*6</sup>, Xingyu Liu<sup>\*7</sup>, Miguel Martin<sup>\*1</sup>, Tushar Nagarajan<sup>\*1,2</sup>, Ilija Radosavovic<sup>\*8</sup>, Santhosh Kumar Ramakrishnan<sup>\*1,2</sup>, Fiona Ryan<sup>\*6</sup>, Jayant Sharma<sup>\*3</sup>, Michael Wray<sup>\*9</sup>, Mengmeng Xu<sup>\*10</sup>, Eric Zhongcong Xu<sup>\*11</sup>, Chen Zhao<sup>\*10</sup>, Siddhant Bansal<sup>17</sup>, Dhruv Batra<sup>1</sup>, Vincent Cartillier<sup>1,6</sup>, Sean Crane<sup>7</sup>, Tien Do<sup>3</sup>, Morrie Doulaty<sup>13</sup>, Akshay Erapalli<sup>13</sup>, Christoph Feichtenhofer<sup>1</sup>, Adriano Fragnani<sup>9</sup>, Qichen Fu<sup>7</sup>, Christian Fuegen<sup>13</sup>, Abrham Gebreselasie<sup>12</sup>, Cristina González<sup>14</sup>, James Hillis<sup>5</sup>, Xuhua Huang<sup>7</sup>, Yifei Huang<sup>15</sup>, Wenqi Jia<sup>6</sup>, Weslie Khoo<sup>16</sup>, Jachym Kolar<sup>13</sup>, Satwik Kottur<sup>13</sup>, Anurag Kumar<sup>5</sup>, Federico Landini<sup>13</sup>, Chao Li<sup>5</sup>, Yanghao Li<sup>1</sup>, Zhenqiang Li<sup>15</sup>, Karttikeya Mangalam<sup>1,8</sup>, Raghava Modhug<sup>17</sup>, Jonathan Munro<sup>9</sup>, Tullie Murrell<sup>1</sup>, Takumi Nishiyasu<sup>15</sup>, Will Price<sup>9</sup>, Paola Ruiz Puentes<sup>14</sup>, Merey Ramazanova<sup>10</sup>, Leda Sari<sup>5</sup>, Kiran Somasundaram<sup>5</sup>, Audrey Southerland<sup>6</sup>, Yusuke Sugano<sup>15</sup>, Ruijie Tao<sup>11</sup>, Minh Vo<sup>5</sup>, Yuchen Wang<sup>16</sup>, Xindi Wu<sup>7</sup>, Takuma Yagi<sup>15</sup>, Yunyi Zhu<sup>11</sup>, Pablo Arbeláez<sup>†14</sup>, David Crandall<sup>†16</sup>, Dima Damen<sup>†9</sup>, Giovanni Maria Farinella<sup>†4</sup>, Bernard Ghanem<sup>†10</sup>, Vamsi Krishna Ithapu<sup>†5</sup>, C. V. Jawahar<sup>†17</sup>, Hanbyul Joo<sup>†1</sup>, Kris Kitani<sup>†7</sup>, Haizhou Li<sup>†11</sup>, Richard Newcombe<sup>†5</sup>, Aude Oliva<sup>†18</sup>, Hyun Soo Park<sup>†3</sup>, James M. Rehg<sup>†6</sup>, Yoichi Sato<sup>†15</sup>, Jianbo Shi<sup>†19</sup>, Mike Zheng Shou<sup>†11</sup>, Antonio Torralba<sup>†18</sup>, Lorenzo Torresani<sup>†1,20</sup>, Mingfei Yan<sup>†5</sup>, Jitendra Malik<sup>1,8</sup>

<sup>1</sup>Facebook AI Research (FAIR), <sup>2</sup>University of Texas at Austin, <sup>3</sup>University of Minnesota, <sup>4</sup>University of Catania,

<sup>5</sup>Facebook Reality Labs, <sup>6</sup>Georgia Tech, <sup>7</sup>Carnegie Mellon University, <sup>8</sup>UC Berkeley, <sup>9</sup>University of Bristol,

<sup>10</sup>King Abdullah University of Science and Technology, <sup>11</sup>National University of Singapore,

<sup>12</sup>Carnegie Mellon University Africa, <sup>13</sup>Facebook, <sup>14</sup>Universidad de los Andes, <sup>15</sup>University of Tokyo, <sup>16</sup>Indiana University,

<sup>17</sup>International Institute of Information Technology, Hyderabad, <sup>18</sup>MIT, <sup>19</sup>University of Pennsylvania, <sup>20</sup>Dartmouth

## Abstract

We introduce Ego4D, a massive-scale egocentric video dataset and benchmark suite. It offers 3,025 hours of daily-life activity video spanning hundreds of scenarios (household, outdoor, workplace, leisure, etc.) captured by 855 unique camera wearers from 74 worldwide locations and 9 different countries. The approach to collection is designed to uphold rigorous privacy and ethics standards with consenting participants and robust de-identification procedures where relevant. Ego4D dramatically expands the volume of diverse egocentric video footage publicly available to the research community. Portions of the video are accompanied by audio, 3D meshes of the environment, eye gaze, stereo, and/or synchronized videos from multiple egocentric cameras at the same event. Furthermore, we present a host of new benchmark challenges centered around understanding

the first-person visual experience in the past (querying an episodic memory), present (analyzing hand-object manipulation, audio-visual conversation, and social interactions), and future (forecasting activities). By publicly sharing this massive annotated dataset and benchmark suite, we aim to push the frontier of first-person perception. Project page: <https://ego4d-data.org/>

## 1. Introduction

Today’s computer vision systems excel at naming objects and activities in Internet photos or video clips. Their tremendous progress over the last decade has been fueled by major dataset and benchmark efforts, which provide the annotations needed to train and evaluate algorithms on well-defined tasks [47, 58, 59, 87, 102, 137].

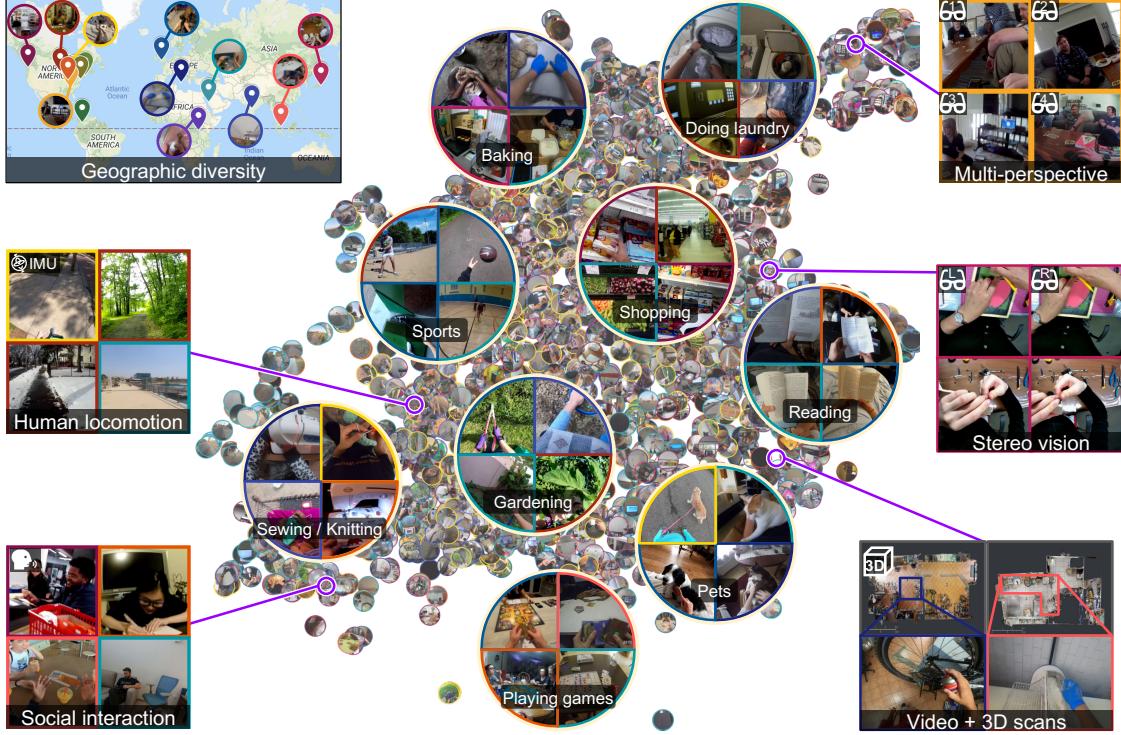


Figure 1. Ego4D is a massive-scale egocentric video dataset of daily life activity spanning 74 locations worldwide. Here we see a snapshot of the dataset (5% of the clips, randomly sampled) highlighting its diversity in geographic location, activities, and modalities. The data includes social videos where participants consented to remain unblurred. See <https://ego4d-data.org/fig1.html> for interactive figure.

While this progress is exciting, current datasets and models represent only a limited definition of visual perception. First, today’s influential Internet datasets capture brief, isolated moments in time from a third-person “spectator” view. However, in both robotics and augmented reality, the input is a long, fluid video stream from the *first-person* or “*egocentric*” point of view—where we see the world through the eyes of an agent actively engaged with its environment. Second, whereas Internet photos are intentionally captured by a human photographer, images from an always-on wearable egocentric camera lack this active curation. Finally, *first-person* perception requires a persistent 3D understanding of the camera wearer’s physical surroundings, and must interpret objects and actions in a human context—attentive to human-object interactions and high-level social behaviors.

Motivated by these critical contrasts, we present the Ego4D dataset and benchmark suite. Ego4D aims to catalyze the next era of research in first-person visual perception. *Ego* is for egocentric, and *4D* is for 3D spatial plus temporal information.

Our first contribution is the dataset: a massive ego-video collection of unprecedented scale and diversity that captures daily life activity around the world. See Figure 1. It consists of 3,025 hours of video collected by 855 unique participants from 74 worldwide locations in 9 different countries. The vast majority of the footage is unscripted and “in

the wild”, representing the natural interactions of the camera wearers as they go about daily activities in the home, workplace, leisure, social settings, and commuting. Based on self-identified characteristics, the camera wearers are of varying backgrounds, occupations, gender, and ages—not solely graduate students! The video’s rich geographic diversity supports the inclusion of objects, activities, and people frequently absent from existing datasets. Since each participant wore a camera for 1 to 10 hours at a time, the dataset offers long-form video content that displays the full arc of a person’s complex interactions with the environment, objects, and other people. In addition to RGB video, portions of the dataset also provide audio, 3D mesh scans, gaze, stereo, and/or synchronized multi-camera views that allow seeing one event from multiple perspectives. Our dataset draws inspiration from prior egocentric video data efforts [41, 42, 123, 125, 132, 173, 195, 198, 203], but makes significant advances in terms of scale, diversity, and realism.

Equally important to having the right data is to have the right research problems. Our second contribution is a suite of five benchmark tasks spanning the essential components of egocentric perception—indexing past experiences, analyzing present interactions, and anticipating future activity. To enable research on these fronts, we provide millions of rich annotations that resulted from over 250,000 hours of annotator effort and range from temporal, spatial, and seman-

tic labels, to dense textual narrations of activities, natural language queries, and speech transcriptions.

The Ego4D project is the culmination of an intensive two-year effort by Facebook and 13 universities around the world who came together for the common goal of spurring new research in egocentric perception. We will kickstart that work with a formal benchmark challenge to be held in June 2022. In the coming years, we believe our contribution can catalyze new research not only in vision, but also robotics, augmented reality, 3D sensing, multimodal learning, speech, and language. These directions will stem not only from the benchmark tasks we propose, but also alternative ones that the community will develop leveraging our massive, publicly available dataset.

## 2. Related Work

**Large-scale third-person datasets** In the last decade, annotated datasets have both presented new problems in computer vision and ensured their solid evaluation. Existing collections like Kinetics [102], AVA [87], UCF [200], ActivityNet [59], HowTo100M [151], ImageNet [47], and COCO [137] focus on third-person Web data, which have the benefit and bias of a human photographer. In contrast, Ego4D is first-person. Passively captured wearable camera video entails unusual viewpoints, motion blur, and lacks temporal curation. Notably, pre-training egocentric video models with third-person data [67, 214, 217, 232] suffers from the sizeable domain mismatch [133, 195].

**Egocentric video understanding** Egocentric video offers a host of interesting challenges, such as human-object interactions [25, 44, 157], activity recognition [104, 133, 236], anticipation [3, 71, 82, 138, 198], video summarization [46, 123, 125, 141, 142, 225], detecting hands [15, 128], parsing social interactions [63, 162, 224], and inferring the camera wearer’s body pose [101]. Our dataset can facilitate new work in all these areas and more, and our proposed benchmarks (and annotations thereof) widen the span of tasks researchers can consider moving forward. We defer discussion of how prior work relates to our benchmark tasks to Sec. 5.

**Egocentric video datasets** Multiple egocentric datasets have been developed over the last decade. Most relevant to our work are those containing unscripted daily life activity, which includes EPIC-Kitchens [41, 42], UT Ego [123, 203], Activities of Daily Living (ADL) [173], and the Disney dataset [63]. The practice of giving cameras to participants to take out of the lab, first explored in [63, 123, 173], inspires our approach. Others are (semi-)scripted, where camera wearers are instructed to perform a certain activity, as in Charades-Ego [195] and EGTEA [132]. Whereas today’s largest ego datasets focus solely on kitchens [42, 118, 132], Ego4D spans hundreds of environments

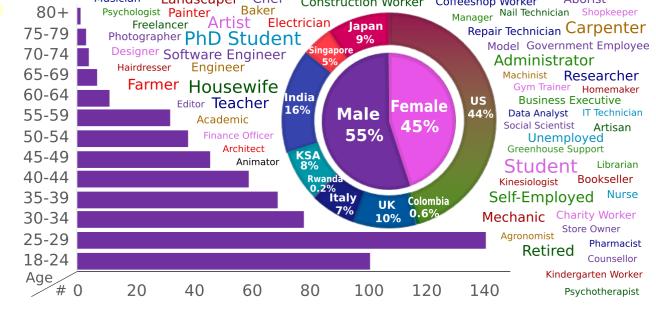


Figure 2. Ego4D camera wearer demographics—age, gender, countries of residence, and occupations (self-reported). Font size reflects relative frequency of the occupation.

both indoors and outdoors. Furthermore, while existing datasets rely largely on graduate students as camera wearers [41, 42, 63, 123, 123, 132, 162, 173, 188, 203], Ego4D camera wearers are of a much wider demographic, as detailed below. Aside from daily life activity, prior ego datasets focus on conversation [164], inter-person interactions [63, 162, 188, 224], place localization [177, 201], multimodal sensor data [118, 160, 197], human hands [15, 128] human-object interaction [178], and object tracking [54].

Ego4D is an order of magnitude larger than today’s largest egocentric datasets both in terms of hours of video (3,025 hours vs. 100 in [41]) and unique camera wearers (855 people vs. 71 in [195]); it spans hundreds of environments (rather than one or dozens, as in existing collections); and its video comes from 74 worldwide locations and 9 countries (vs. just one or a few cities). The Ego4D annotations are also of unprecedented scale and depth, with millions of annotations supporting multiple complex tasks. As such, Ego4D represents a step change in dataset scale and diversity. We believe both factors are paramount to pursue the next generation of perception for embodied AI.

## 3. Ego4D Dataset

Next we overview the dataset, which we are making publicly available under an Ego4D license.

### 3.1. Collection strategy and camera wearers

Not only do we wish to amass an ego-video collection that is substantial in scale, but we also want to ensure its diversity of people, places, objects, and activities. Furthermore, for realism, we are interested in unscripted footage captured by people wearing a camera for long periods of time.

To this end, we devised a distributed approach to data collection. The Ego4D project consists of 14 teams from universities and labs in 9 countries and 5 continents (see map in Figure 1). Each team recruited participants to wear a camera for 1 to 10 hours at a time, for a total of 855 unique camera wearers and 3,025 hours of video in this first dataset

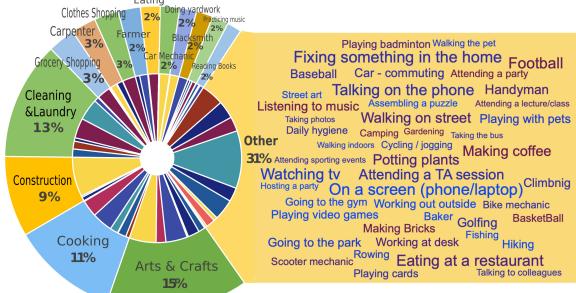


Figure 3. Scenarios in Ego4D. Outer circle shows the 14 most common scenarios (70% of the data). Wordle shows scenarios in the remaining 30%. Inner circle is color coded by the contributing partner (see map color legend in Fig 1).

release (Ego4D-3K). Participants in 74 total cities were recruited by word of mouth, ads, and postings on community bulletin boards. Some teams recruited participants with occupations that have interesting visual contexts, such as bakers, carpenters, landscapers, or mechanics.

Both the geographic spread of our team as well as our approach to recruiting participants were critical to arrive at a diverse demographic composition, as shown in Figure 2.<sup>1</sup> Participants cover a wide variety of occupations, span many age brackets, with 97 of them over 50 years old, and 45% are female. Two participants identified as non-binary, and three preferred not to say a gender.

### 3.2. Scenarios composing the dataset

What activities belong in an egocentric video dataset? Our research is motivated by problems in robotics and augmented reality, where vision systems will encounter *daily life scenarios*. Hence, we consulted a survey from the U.S. Bureau of Labor Statistics<sup>2</sup> that captures how people spend the bulk of their time in the home (e.g., cleaning, cooking, yardwork), leisure (e.g., crafting, games, attending a party), transportation (e.g., biking, car), errands (e.g., shopping, walking dog, getting car fixed), and in the workplace (e.g., talking with colleagues, making coffee).

To maximize coverage of such scenarios, our approach is a compromise between directing camera wearers and giving no guidance at all: 1) we recruited participants whose collective daily life activity would naturally encompass a spread of the scenarios (as selected freely by the participant), and 2) we asked participants to wear the camera at length (at least as long as the battery life of the device) so that the activity would unfold naturally in a longer context. A typical raw video clip in our dataset lasts 8 minutes—significantly longer

<sup>1</sup>for 69% of all participants; missing demographics are due to protocols or participants opting out of answering specific questions. Our dataset continues to grow in size, including additional incoming data not reflected in this chart from the partners in Colombia and Rwanda, who joined the consortium more recently and are actively collecting and processing data.

<sup>2</sup><https://www.bls.gov/news.release/atus.nr0.htm>

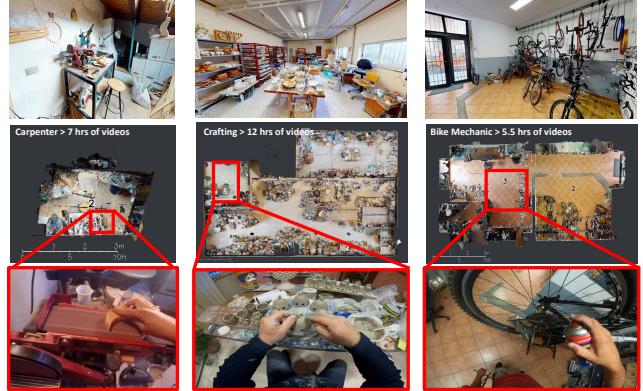


Figure 4. Some videos (bottom) have coupled 3D meshes (top) from Matterport3D scanners, allowing one to relate the dynamic video to the static 3D environment (middle).

than the 10 second clips often studied in third-person video understanding [102]. In this way, we capture unscripted activity while being mindful of the scenarios' coverage.

The exception is for certain multi-person scenarios, where, in order to ensure sufficient data for the audio-visual and social benchmarks, we asked participants at five sites who had consented to share their conversation audio and unblurred faces to take part in social activities, such as playing games. We leverage this portion of Ego4D for the audio-visual and social interaction benchmarks (Sec. 5.3 and 5.4).

Figure 3 shows the wide distribution of scenarios captured in our dataset. Note that within each given scenario there are typically dozens of actions taking place, e.g., the carpentry scenario includes hammering, drilling, moving wood, etc. Overall, the 855 camera wearers bestow our dataset with a glimpse of daily life activity around the world.

### 3.3. Cameras and modalities

To avoid models overfitting to a single capture device, seven different head-mounted cameras were deployed across the dataset: GoPro, Vuzix Blade, Pupil Labs, ZShades, OR-DRO EP6, iVue Rincon 1080, and Weevie. They offer tradeoffs in the modalities available (RGB, stereo, gaze), field of view, and battery life. The field of view and camera mounting are particularly influential: while a GoPro mounted on the head pointing down offers a high resolution view of the hands manipulating objects (Fig. 5, right), a heads-up camera like the Vuzix shares the vantage of a person's eyes, but will miss interactions close to the body (Fig. 5, left).

In addition to video, portions of Ego4D offer several other data modalities: 3D scans, audio, gaze<sup>3</sup>, stereo, multiple synchronized wearable cameras, and textual narrations. See Table 1. Each can support new research challenges. For example, having Matterport3D scans of the environment

<sup>3</sup>Eye trackers were deployed by Indiana University only.

Modality:	RGB video	Text narrations	Features	Audio	Faces	3D scans	Stereo	Gaze	IMU	Multi-cam
# hours:	3,025	3,025	3,025	2,207	612	491	80	70	836	224

Table 1. Modalities of data in Ego4D and their amounts. “Narrations” are dense, timestamped descriptions of camera wearer activity (cf. Sec. 4). “3D scans” are meshes from Matterport3D scanners for the full environment in which the video was captured. “Faces” refers to video where participants consented to remain unblurred. “Multi-cam” refers to synchronized video captured at the same event by multiple camera wearers. “Features” refers to precomputed SlowFast [67] video features. Gaze collected only by Indiana University.

coupled with ego-video clips (Figure 4) offers a unique opportunity for understanding dynamic activities in a persistent 3D context, as we exploit in the Episodic Memory benchmark (see Sec. 5.1). Multiple synchronized egocentric video streams allow accounting for the first and second-person view in social interactions. Audio allows analysis of conversation and acoustic scenes and events.

### 3.4. Privacy and ethics

From the onset, privacy and ethics standards were critical to this data collection effort. Each partner was responsible for developing a policy. While specifics vary per site, this generally entails:

- Comply with own institutional research policy, e.g., independent ethics committee review where relevant.
- Obtain informed consent of camera wearers, who can ask questions and withdraw at any time, and are free to review and redact their own video
- Respect rights of others in private spaces, and avoid capture of sensitive areas or activities
- Follow de-identification requirements for personally identifiable information (PII)

In short, these standards typically require that the video be captured in a controlled environment with informed consent by all participants, or else in public spaces where faces and other PII are blurred.

### 3.5. Possible sources of bias

While Ego4D pushes the envelope on massive everyday video from geographically and demographically diverse sources, we are aware of a few biases in our dataset. 74 locations is still a long way from complete coverage of the globe. In addition, the camera wearers are generally located in urban or college town areas. The COVID-19 pandemic led to ample footage in stay-at-home scenarios such as cooking, cleaning, crafts, etc. and more limited opportunities to collect video at major social public events. In addition, since battery life prohibits daylong filming, the videos—though unscripted—tend to contain more active portions of a participant’s day. Finally, Ego4D annotations are done by crowd-sourced workers in two sites in Africa. This means that there will be at least subtle ways in which the language-based narrations are biased towards their local word choices.



Figure 5. Example narrations. “C” refers to camera wearer.

## 4. Narrations of Camera Wearer Activity

Before any other annotation occurs, we pass all video through a *narration* procedure. Inspired by the *pause-and-talk narrator* [42], annotators are asked to watch a 5 minute clip of video, summarize it with a few sentences, and then re-watch, pausing repeatedly to write a sentence about each thing the camera wearer does. We record the timestamps and the associated free-form sentences. See Figure 5. Each video receives two independent narrations from different annotators. The narrations are temporally dense: on average we received 13.2 sentences per minute of video, for a total of 3.85M sentences. In total the narrations describe the Ego4D video using 1,772 unique verbs (activities) and 4,336 unique nouns (objects). See Appendix D for details.

The narrations allow us to 1) perform text mining for data-driven taxonomy construction for actions and objects, 2) sort the videos by their content to map them to relevant benchmarks, and 3) identify temporal windows where certain annotations should be seeded. Beyond these uses, the narrations are themselves a contribution of the dataset, potentially valuable for research on video with weakly aligned natural language. To our knowledge, ours is the largest repository of aligned language and video (e.g., HowTo100M [151], an existing Internet repository with narrations, contains noisy spoken narrations that only sometimes comment on the activities taking place).

## 5. Ego4D Benchmark Suite

First-person vision has the potential to transform many applications in augmented reality and robotics. However, compared to mainstream video understanding, egocentric perception requires new fundamental research to account for long-form video, attention cues, person-object interactions, multi-sensory data, and the lack of manual temporal curation inherent to a passively worn camera.



Figure 6. The Ego4D benchmark suite centers around the first-person visual experience—from remembering the past, to analyzing the present, to anticipating the future.

Inspired by all these factors, we propose a suite of challenging benchmark tasks. The five benchmarks tackle the *past*, *present*, and *future* of first-person video. See Figure 6. The following sections introduce each task and its annotations. The first dataset release has annotations for 50-800 hours of data per benchmark, on top of the 3,025 hours of data that is narrated. The Appendix describes how we sampled videos per benchmark to maximize relevance to the task while maintaining geographic diversity.

We developed baseline models drawing on state-of-the-art components from the literature in order to test drive all Ego4D benchmarks. **The Appendix presents the baseline models and quantitative results.** We will run a formal Ego4D competition in June 2022 inviting the research community to improve on these baselines.

## 5.1. Episodic Memory

**Motivation** Egocentric video from a wearable camera records the who/what/when/where of an individual’s daily life experience. This makes it ideal for what Tulving called *episodic memory* [206]: specific first-person experiences (“what did I eat and who did I sit by on my first flight to France?”), to be distinguished from *semantic memory* (“what’s the capital of France?”). An augmented reality assistant that processes the egocentric video stream could give us super-human memory if it could appropriately index our visual experience and answer queries.

**Task definition** Given an egocentric video and a query, the Ego4D Episodic Memory task requires localizing where the answer can be seen within the user’s past video. We consider three query types. (1) *Natural language queries* (NLQ), in which the query is expressed in text (e.g., “What did I put in the drawer?”), and the output response is the temporal window where the answer is visible or deducible. (2) *Visual queries* (VQ), in which the query is a static image of an object, and the output response localizes the object the last time it was seen in the video, both temporally and spatially. The spatial response is a 2D bounding box on the object, and optionally a 3D displacement vector from the current camera position to the object’s 3D bounding box. VQ captures how a user might teach the system an object with an image example, then later ask for its location (“Where

in this new video  
or just more  
annotation?  
why release  
separately?

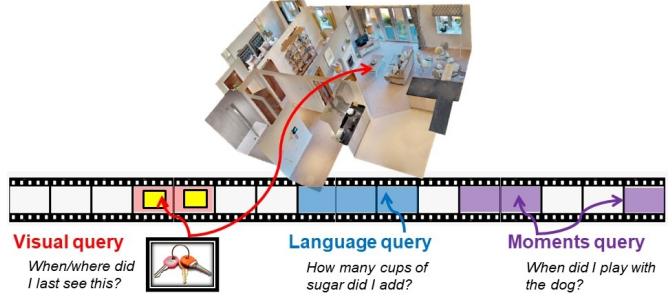


Figure 7. Episodic Memory’s three query types

is this [picture of my keys]?”). (3) *Moments queries* (MQ), in which the query is the name of a high-level activity or “moment”, and the response consists of all temporal windows where the activity occurs (e.g., “When did I read to my children?”). See Figure 7.

**Annotations** For language queries, we devised a set of 13 template questions meant to span things a user might ask to augment their memory, such as “*what is the state of object X?*”, e.g., “did I leave the window open?”. Annotators express the queries in free-form natural language, and also provide the slot filling (e.g., X = window). For moments, we established a taxonomy of 110 activities in a data-driven, semi-automatic manner by mining the narration summaries. Moments capture high-level activities in the camera wearer’s day, e.g., *setting the table* is a moment, whereas *pick up* is an action in our Forecasting benchmark (Sec. 5.5).

For NLQ and VQ, we ask annotators to generate language/visual queries and couple them with the “response track” in the video. For MQ, we provide the taxonomy of labels and ask annotators to label clips with each and every temporal segment containing a moment instance. In total, we have ~74K total queries spanning 800 hours of video.

**Evaluation metrics and baselines** For NLQ, we use top-k recall at a certain temporal intersection over union (tIoU) threshold. MQ adopts a popular metric used in temporal action detection: *mAP at multiple tIoU thresholds*, as well as *top-kx recall at multiple tIoU thresholds*. VQ adopts temporal and spatio-temporal localization metrics as well as timeliness metrics that encourage speedy searches. Ap-

pendix E presents the baseline models we developed and reports results.

**Relation to existing tasks** Episodic Memory has some foundations in existing vision problems, but also adds new challenges. All three queries call for spatial reasoning in a static environment coupled with dynamic video of a person who moves and changes things; current work largely treats these two elements separately. The timeliness metrics encourage work on intelligent contextual search. While current literature on language+vision focuses on captioning and question answering for isolated instances of Internet data [11, 33, 113, 221], NLQ is motivated by queries about the camera wearer’s own visual experience and operates over long-term observations. VQ upgrades object instance recognition [22, 81, 120, 149] to deal with video (frequent FoV changes, objects entering/exiting the view) and to reason about objects in the context of a 3D environment. Finally, MQ can be seen as activity detection [135, 222, 230] but for the activities of the camera wearer.

## 5.2. Hands and Objects

**Motivation** While Episodic Memory aims to make *past* video queryable, our next benchmark aim to understand the camera wearer’s *present* activity—in terms of interactions with objects and other people. Specifically, the **Hands and Objects benchmark** captures how the camera wearer changes the state of an object by using or manipulating it—which we call an *object state change*. Though cutting a piece of lumber in half can be achieved through many methods (*e.g.*, various tools, force, speed, grasps, end-effectors), all should be recognized as the same state change. This generalization ability will enable us to understand human actions better, as well as to train robots to learn from human demonstrations in video.

**Task definitions** We interpret an object state change to include various **physical changes**, including changes in **size**, **shape**, **composition**, and **texture**. Object state changes can be viewed along **temporal**, **spatial** and **semantic axes**, leading to these three tasks: (1) *Point-of-no-return temporal localization*: given a short video clip of a state change, the goal is to estimate the keyframe that contains the point-of-no-return (PNR) (the time at which a state change begins); (2) *State change object detection*: given three temporal frames (pre, post, PNR), the goal is to regress the bounding box of the object undergoing a state change; (3) *Object state change classification*: given a short video clip, the goal is to classify whether an object state change has taken place or not.

**Annotations** We select the data to annotate based on activities that are likely to involve hand-object interactions (*e.g.*, knitting, carpentry, baking, *etc.*). We start by labeling each narrated hand-object interaction. For each, we label three moments in time (pre, PNR, post) and the bounding boxes



Figure 8. Hand and Objects: Example object state changes defined by pre-condition, PNR, and post-condition frames.

for the hands, tools, and objects in each of the three frames. We also annotate the **state change types** (remove, burn, *etc.*, see Fig. 8), **action verbs**, and **nouns** for the objects.

**Evaluation metrics and baselines** Object state change temporal localization is evaluated using absolute temporal error measured in seconds. Object state change classification is evaluated by classification accuracy. State change object detection is evaluated by average precision (AP). Appendix F details the annotations and presents baseline model results for the three Hands and Objects tasks.

**Relation to existing tasks** Limited prior work considers object state change in photos [97, 158] or video [7, 65, 235]; Ego4D is the first video benchmark dedicated to the task of understanding object state changes. The task is similar to action recognition (*e.g.*, [95, 104, 133, 214, 236]) because in some cases a specific action can correspond to a specific state change. However, a single state change (*e.g.*, cutting) can also be observed in many forms (various object-tool-action combinations). It is our hope that the proposed benchmarks will lead to the development of more explicit models of object state change, while avoiding approaches that simply overfit to action or object observations.

## 5.3. Audio-Visual Diarization

**Motivation** Our next two tasks aim to understand the camera wearer’s present interactions with *people*. People communicate using **spoken language**, making the capture of conversational content in business meetings and social settings a problem of great scientific and practical interest. While diarization has been a standard problem in the speech recognition community, Ego4D brings in two new aspects (1) simultaneous capture of video and audio (2) the egocentric perspective of a participant in the conversation.

**Task definition and annotations** The Audio-Visual Diarization (AVD) benchmark is composed of four tasks (see Figure 9):



Figure 9. Audio-Visual and Social benchmark annotations

- **Localization and tracking** of the participants (i.e., candidate speakers) in the visual field of view (FoV). A bounding box is annotated around each participant’s face.
- **Active speaker detection** where each tracked speaker is assigned an anonymous label, including the camera wearer who never appears in the visual FoV.
- **Diarization** of each speaker’s speech activity, where we provide the time segments corresponding to each speaker’s voice activity in the clip.
- **Transcription** of each speaker’s speech content (only English speakers are considered for this version).

**Evaluation metrics and baselines** We use standardized object tracking (**MOT**) metrics [17, 18] to evaluate speaker localization and tracking in the visual FoV. Speaker detection with anonymous labels is evaluated using the speaker error rate, which measures the proportion of wrongly assigned labels. We adopt the well studied diarization error rate (DER) [10] and word error rate (WER) [108] for diarization and transcription, respectively. We present AVD baseline models and results in Appendix G.

**Relation to existing tasks** The past few years have seen audio studied in computer vision tasks [238] for action classification [104, 219], object categorization [119, 227], source localization and tracking [13, 191, 205] and embodied navigation [31]. Meanwhile, visual information is increasingly used in historically audio-only tasks like speech transcription, voice recognition, audio spatialization [4, 76, 99, 155], speaker diarization [9, 79], and source separation [55, 74, 78]. Datasets like VoxCeleb [37], AVA Speech [30], AVA active speaker [186], AVDIAR [79], and EasyCom [51] support this research. However, these datasets are mainly non-egocentric. Unlike Ego4D, they do not capture natural conversational characteristics involving a variety of noisy backgrounds, overlapping, interrupting and unintelligible speech, environment variation, moving camera wearers, and speakers facing away from the camera wearer.

## 5.4. Social Interactions

**Motivation** An egocentric video provides a unique lens for studying social interactions because it captures utterances and nonverbal cues [109] from each participant’s unique view and enables embodied approaches to social understanding. Progress in egocentric social understanding could lead to more capable virtual assistants and social robots. Computational models of social interactions can also provide new tools for diagnosing and treating disorders of socialization and communication such as autism [182], and could support novel prosthetic technologies for the hearing-impaired.

**Task definition** While the Ego4D dataset can support such a long-term research agenda, our initial Social benchmark focuses on multimodal understanding of conversational interactions via attention and speech. Specifically, we focus on identifying communicative acts that are directed towards the camera-wearer, as distinguished from those directed to other social partners: (1) **Looking at me (LAM)**: given a video in which the faces of social partners have been localized and identified, classify whether each visible face is looking at the camera wearer; and (2) **Talking to me (TTM)**: given a video and audio segment with the same tracked faces and an additional label that identifies speaker status, classify whether each visible face is talking to the camera wearer.

**Annotations** Social annotations build on those from AV diarization (Sec. 5.3). Given (1) face bounding boxes labeled with participant IDs and tracked across frames, and (2) associated active speaker annotations that identify in each frame whether the social partners whose faces are visible are speaking, annotators provide the ground truth labels for LAM and TTM as a binary label for each face in each frame. For LAM, annotators label the time segment (start and end time) of a visible person when the individual is looking at the camera wearer. For TTM, we use the vocal activity annotation from AVD, then identify the time segment when the speech is directed at the camera wearer. See Figure 9.

**Evaluation metrics and baselines** We use mean average precision (mAP) and Top-1 accuracy to quantify the classification performance for both tasks. Unlike AVD, we measure precision at every frame. Appendix H provides details and presents Social baseline models and results.

**Relation to existing tasks** Compared to [64], Ego4D contains substantially more participants, hours of recording, and variety of sensors and social contexts. The LAM task is most closely related to prior work on eye contact detection in ego-video [34, 153], but addresses more diverse and challenging scenarios. Mutual gaze estimation [52, 144–146, 166, 170] and gaze following [35, 62, 105, 180] are also relevant. The TTM task is related to audio-visual speaker detection [6, 187] and meeting understanding [20, 126, 148].

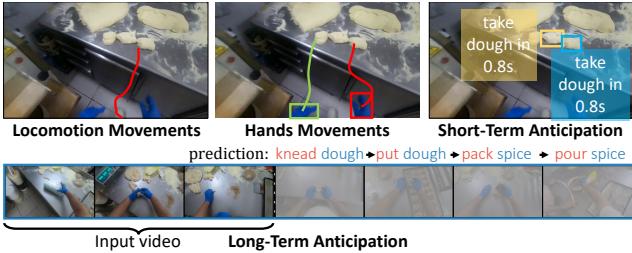


Figure 10. The Forecasting benchmark aims to predict future locomotion, movement of hands, next object interactions, and sequences of future actions.

## 5.5. Forecasting

**Motivation** Having addressed the past and present of the camera wearer’s visual experience, our last benchmark moves on to anticipating the future. Forecasting movements and interactions requires comprehending the camera wearer’s *intention*. It has immediate applications in AR and human-robot interaction, such as anticipatively turning on appliances or moving objects for the human’s convenience. The scientific motivation can be seen by analogy with language models such as GPT-3 [23], which implicitly capture knowledge needed by many other tasks. Rather than predict the next word, visual forecasting models the dynamics of an agent acting in the physical world.

**Task definition** The Forecasting benchmark includes four tasks (Fig. 10): (1) *Locomotion prediction*: predict a set of possible future ground plane trajectories of the camera wearer. (2) *Hand movement prediction*: predict the hand positions of the camera wearer in future frames. (3) *Short-term object interaction anticipation*: detect a set of possible future interacted objects in the most recent frame of the clip. To each object, assign a verb indicating the possible future interaction and a “time to contact” estimate of when the interaction is going to begin. (4) *Long-term action anticipation*: predict the camera wearer’s future sequence of actions.

**Annotations** Using the narrations, we identify the occurrence of each object interaction, assigning a verb and a target object class. The verb and noun taxonomies are seeded from the narrations and then hand-refined. For each action, we identify a contact frame and a pre-condition frame in which we annotate bounding boxes around active objects. The same objects as well as hands are annotated in three frames preceding the pre-condition frame by 0.5s, 1s and 1.5s. We obtain ground truth ego-trajectories of the camera wearer using structure from motion.

**Evaluation metrics and baselines** We evaluate future locomotion movement and hand movement prediction using L2 distance. Short-term object interaction anticipation is evaluated using a Top-5 mean Average Precision metric which “discounts” the Top-4 false negative predictions. Long-term

action anticipation is evaluated using edit distance. Appendix I details the Forecasting tasks, annotations, baseline models, and results.

**Relation to existing tasks** Predicting future events from egocentric vision has increasing interest [185]. Previous work considers future localization [107, 114, 168, 223], action anticipation [72, 73, 82, 112, 121, 212], next active object prediction [19, 70], future event prediction [143, 161], and future frame prediction [139, 140, 147, 208, 211, 220]. Whereas past work relies on different benchmarks and task definitions, we propose a unified benchmark to assess progress in the field.

## 6. Conclusion

Ego4D is a first-of-its-kind dataset and benchmark suite aimed at advancing multimodal perception of egocentric video. Compared to existing work, our dataset is orders of magnitude larger in scale and diversity. The data will allow AI to learn from daily life experiences around the world—seeing what we see and hearing what we hear—while our benchmark suite provides solid footing for innovations in video understanding that are critical for augmented reality, robotics, and many other domains. We look forward to the research that will build on Ego4D in the years ahead.

## Contribution statement

Project led and initiated by Kristen Grauman. Program management and operations led by Andrew Westbury. Scientific advising by Jitendra Malik. Authors with stars (\*) were key drivers of implementation, collection, and/or annotation development throughout the project. Authors with daggers (†) are faculty PIs and working group leads in the project. The benchmarks brought together many researchers from all institutions including cross-institution baseline evaluations. Appendix E through H detail the contributions of individual authors for the various benchmarks. The video collected by Facebook Reality Labs used Vuzix Blade® Smart Glasses and was done in a closed environment in Facebook’s buildings by paid participants who signed consents to share their data. All other video collection and participant recruitment was managed by the university partners. Appendix A provides details about the data collection done per site and acknowledges the primary contributors. The annotation effort was led by Facebook AI.

## Acknowledgements

We gratefully acknowledge the following colleagues for valuable discussions and support of our project: Aaron Adcock, Andrew Allen, Behrouz Behmardi, Serge Belongie, Mark Broyles, Xiao Chu, Samuel Clapp, Irene D’Ambra, Peter Dodds, Jacob Donley, Ruohan Gao, Tal Hassner, Ethan

Henderson, Jiabo Hu, Guillaume Jeanneret, Sanjana Krishnan, Tsung-Yi Lin, Bobby Otillar, Manohar Paluri, Maja Pantic, Lucas Pinto, Vivek Roy, Jerome Pesenti, Joelle Pineau, Luca Sbordone, Rajan Subramanian, Helen Sun, Mary Williamson, and Bill Wu. We also acknowledge Jacob Chalk for setting up the Ego4D AWS backend and Prasanna Sridhar for developing the Ego4D website. Thank you to the Common Visual Data Foundation (CVDF) for hosting the Ego4D dataset.

The universities acknowledge the usage of commercial software for de-identification of video. brighter.ai was used for redacting videos by some of the universities. Personal data from the University of Bristol was protected by Primloc's Secure Redact software suite.

UNICT is supported by MIUR AIM - Attrazione e MobilitàInternazionale Linea 1 - AIM1893589 - CUP E64118002540007. Bristol is supported by UKRIEngineering and Physical Sciences Research Council (EPSRC) Doctoral Training Program (DTP), EPSRC Fellowship UMPIRE (EP/T004991/1). KAUST is supported by the KAUST Office of Sponsored Research through the Visual Computing Center (VCC) funding. National University of Singapore is supported by Mike Shou's Start-Up Grant. Georgia Tech is supported in part by NSF award 2033413 and NIH award R01MH114999.

## References

- [1] Github repository of the ESPNet model zoo. [https://github.com/espnet/espnet\\_model\\_zoo](https://github.com/espnet/espnet_model_zoo). We used the `ShinjiWatanabe/gigaspeech_asr_train_asr_raw_en_bpe5000_valid.acc.ave` model. 51, 56, 57
- [2] NIST SRE 2000 Evaluation Plan. [https://www.nist.gov/sites/default/files/documents/2017/09/26/spk-2000-plan-v1.0.htm\\_.pdf](https://www.nist.gov/sites/default/files/documents/2017/09/26/spk-2000-plan-v1.0.htm_.pdf). 51
- [3] Yazan Abu Farha, Alexander Richard, and Juergen Gall. When will you do what?-anticipating temporal occurrences of activities. In *Computer Vision and Pattern Recognition*, pages 5343–5352, 2018. 3
- [4] Triantafyllos Afouras, Joon Son Chung, Andrew Senior, Oriol Vinyals, and Andrew Zisserman. Deep audio-visual speech recognition. *IEEE transactions on pattern analysis and machine intelligence*, 2018. 8, 47
- [5] Triantafyllos Afouras, Joon Son Chung, and Andrew Zisserman. The conversation: Deep audio-visual speech enhancement. In *Interspeech*, 2018. 47
- [6] Triantafyllos Afouras, Andrew Owens, Joon Son Chung, and Andrew Zisserman. Self-supervised Learning of Audio-Visual Objects from Video. In *Proceedings of the European Conference on Computer Vision (ECCV 20)*, volume 12363 LNCS, pages 208–224, 2020. 8
- [7] Jean-Baptiste Alayrac, Josef Sivic, Ivan Laptev, and Simon Lacoste-Julien. Joint discovery of object states and manipulation actions. *ICCV*, 2017. 7, 41, 42
- [8] Humam Alwassel, Fabian Caba Heilbron, Victor Escorcia, and Bernard Ghanem. Diagnosing error in temporal action detectors. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018. 38, 39
- [9] Xavier Anguera, Simon Bozonnet, Nicholas Evans, Corinne Fredouille, Gerald Friedland, and Oriol Vinyals. Speaker diarization: A review of recent research. *IEEE Transactions on audio, speech, and language processing*, 20(2):356–370, 2012. 8, 49, 51
- [10] Xavier Anguera Miró. *Robust speaker diarization for meetings*. Universitat Politècnica de Catalunya, 2006. 8, 50
- [11] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C. Lawrence Zitnick, and Devi Parikh. VQA: Visual Question Answering. In *International Conference on Computer Vision (ICCV)*, 2015. 7
- [12] Mehmet Ali Arabacı, Fatih Özkan, Elif Surer, Peter Jančovič, and Alptekin Temizel. Multi-modal egocentric activity recognition using audio-visual features. *arXiv preprint arXiv:1807.00612*, 2018. 47
- [13] Relja Arandjelović and Andrew Zisserman. Objects that sound. In *ECCV*, 2018. 8, 47
- [14] Alexei Baevski, Henry Zhou, Abdelrahman Mohamed, and Michael Auli. wav2vec 2.0: A framework for self-supervised learning of speech representations. *arXiv preprint arXiv:2006.11477*, 2020. 55
- [15] Sven Bambach, Stefan Lee, David J. Crandall, and Chen Yu. Lending a hand: Detecting hands and recognizing activities in complex egocentric interactions. In *The IEEE International Conference on Computer Vision (ICCV)*, December 2015. 3
- [16] Mark A Bee and Christophe Micheyl. The cocktail party problem: what is it? how can it be solved? and why should animal behaviorists study it? *Journal of comparative psychology*, 122(3):235, 2008. 48
- [17] Keni Bernardin, Alexander Elbs, and Rainer Stiefelhagen. Multiple object tracking performance metrics and evaluation in a smart room environment. In *Sixth IEEE International Workshop on Visual Surveillance, in conjunction with ECCV*, volume 90. Citeseer, 2006. 8, 49
- [18] Keni Bernardin and Rainer Stiefelhagen. Evaluating multiple object tracking performance: the clear mot metrics. *EURASIP Journal on Image and Video Processing*, 2008:1–10, 2008. 8, 49
- [19] Gedas Bertasius, Hyun Soo Park, Stella X. Yu, and Jianbo Shi. First-person action-object detection with egonet. In *Proceedings of Robotics: Science and Systems*, July 2017. 9
- [20] Cigdem Beyan, Francesca Capozzi, Cristina Becchio, and Vittorio Murino. Prediction of the leadership style of an emergent leader using audio and visual nonverbal features. *IEEE Transactions on Multimedia*, 20(2):441–456, 2018. 8
- [21] Goutam Bhat, Martin Danelljan, Luc Van Gool, and Radu Timofte. Know Your Surroundings: Exploiting Scene Information for Object Tracking. *arXiv:2003.11014 [cs]*, May 2020. 32, 33
- [22] Eric Brachmann, Alexander Krull, Frank Michel, Stefan Gumhold, Jamie Shotton, and Carsten Rother. Learning 6d object pose estimation using 3d object coordinates. In *European conference on computer vision*, pages 536–551. Springer, 2014. 7
- [23] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020. 9
- [24] Ian M Bullock, Thomas Feix, and Aaron M Dollar.

- The yale human grasping dataset: Grasp, object, and task data in household and machine shop environments. *IJRR*, 2015. 42
- [25] Minjie Cai, Kris M Kitani, and Yoichi Sato. Understanding hand-object manipulation with grasp types and object attributes. In *RSS*, 2016. 3
- [26] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European Conference on Computer Vision*, pages 213–229. Springer, 2020. 45, 46
- [27] Jean Carletta, Simone Ashby, Sébastien Bourban, Mike Flynn, Mael Guillemot, Thomas Hain, Jaroslav Kadlec, Vasilis Karaikos, Wessel Kraaij, Melissa Kronenthal, et al. The AMI meeting corpus: A pre-announcement. In *International workshop on machine learning for multimodal interaction*, pages 28–39. Springer, 2006. 51
- [28] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6299–6308, 2017. 44, 45
- [29] Chien-Yi Chang, De-An Huang, Danfei Xu, Ehsan Adeli, Li Fei-Fei, and Juan Carlos Niebles. Procedure planning in instructional videos. *arXiv preprint arXiv:1907.01172*, 2019. 42
- [30] Sourish Chaudhuri, Joseph Roth, Daniel PW Ellis, Andrew Gallagher, Liat Kaver, Radhika Marvin, Caroline Pantofaru, Nathan Reale, Loretta Guarino Reid, Kevin Wilson, et al. Ava-speech: A densely labeled dataset of speech activity in movies. *arXiv preprint arXiv:1808.00606*, 2018. 8, 48
- [31] Changan Chen, Unnat Jain, Carl Schissler, Sebastia Vicenc Amengual Gari, Ziad Al-Halah, Vamsi Krishna Ithapu, Philip Robinson, and Kristen Grauman. Audio-visual embodied navigation. *environment*, 97:103, 2019. 8
- [32] Guoguo Chen, Shuzhou Chai, Guanbo Wang, Jiayu Du, Wei-Qiang Zhang, Chao Weng, Dan Su, Daniel Povey, Jan Trmal, Junbo Zhang, et al. Gigaspeech: An evolving, multi-domain asr corpus with 10,000 hours of transcribed audio. *arXiv preprint arXiv:2106.06909*, 2021. 56
- [33] Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco captions: Data collection and evaluation server. *arXiv preprint arXiv:1504.00325*, 2015. 7
- [34] Eunji Chong, Elysha Clark-Whitney, Audrey Southerland, Elizabeth Stubbs, Chanel Miller, Eliana L Ajo-dan, Melanie R Silverman, Catherine Lord, Agata Rozga, Rebecca M Jones, and James M Rehg. Detection of eye contact with deep neural networks is as accurate as human experts. *Nature Communications*, 11(1):6386, dec 2020. 8
- [35] Eunji Chong, Yongxin Wang, Nataniel Ruiz, and James M. Rehg. Detecting Attended Visual Targets in Video. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR 20)*, pages 5395–5405, Seattle, WA, 2020. 8, 60
- [36] Joon Son Chung, Jaesung Huh, Seongkyu Mun, Minjae Lee, Hee Soo Heo, Soyeon Choe, Chiheon Ham, Sunghwan Jung, Bong-Jin Lee, and Icksang Han. In defence of metric learning for speaker recognition. In *Interspeech*, 2020. 59
- [37] Joon Son Chung, Jaesung Huh, Arsha Nagrani, Triantafyllos Afouras, and Andrew Zisserman. Spot the conversation: speaker diarisation in the wild. *arXiv preprint arXiv:2007.01216*, 2020. 8, 48
- [38] J. S. Chung, A. Nagrani, and A. Zisserman. Vox-Celeb2: Deep Speaker Recognition. In *INTERSPEECH*, 2018. 48
- [39] Kenneth Church and Patrick Hanks. Word association norms, mutual information, and lexicography. *Computational linguistics*, 16(1):22–29, 1990. 66
- [40] Dima Damen, Hazel Doughty, Giovanni Farinella, Sanja Fidler, Antonino Furnari, Evangelos Kazakos, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, et al. The epic-kitchens dataset: Collection, challenges and baselines. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, (01):1–1, 2020. 48
- [41] Dima Damen, Hazel Doughty, Giovanni Maria Farinella, , Antonino Furnari, Jian Ma, Evangelos Kazakos, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, and Michael Wray. Rescaling egocentric vision. *IJCV*, 2021. 2, 3, 42
- [42] Dima Damen, Hazel Doughty, Giovanni Maria Farinella, Sanja Fidler, Antonino Furnari, Evangelos Kazakos, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, and Michael Wray. Scaling egocentric vision: The epic-kitchens dataset. In *European Conference on Computer Vision (ECCV)*, 2018. 2, 3, 5, 19, 48
- [43] Dima Damen, Teesid Leelasawassuk, Osian Haines, Andrew Calway, and Walterio Mayol-Cuevas. You-Do, I-Learn: Discovering task relevant objects and their modes of interaction from multi-user egocentric video. In *BMVC*, 2014. 41, 42
- [44] Dima Damen, Teesid Leelasawassuk, and Walterio Mayol-Cuevas. You-do, i-learn: Egocentric unsupervised discovery of objects and their modes of interaction towards video-based guidance. *CVIU*, 2016. 3
- [45] Fred J Damerau. A technique for computer detection and correction of spelling errors. *Communications of the ACM*, 1964. 68

- [46] Ana Garcia Del Molino, Cheston Tan, Joo-Hwee Lim, and Ah-Hwee Tan. Summarization of egocentric videos: A comprehensive survey. *IEEE Transactions on Human-Machine Systems*, 47(1), 2016. 3
- [47] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *CVPR*, 2009. 1, 3
- [48] Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Superpoint: Self-supervised interest point detection and description. In *CVPR Workshop*, 2018. 35
- [49] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv:1810.04805*, 2018. 26
- [50] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. 37
- [51] Jacob Donley, Vladimir Tourbabin, Jung-Suk Lee, Mark Broyles, Hao Jiang, Jie Shen, Maja Pantic, Vamsi Krishna Ithapu, and Ravish Mehra. Easycom: An augmented reality dataset to support algorithms for easy communication in noisy environments. *arXiv preprint arXiv:2107.04174*, 2021. 8, 48
- [52] Bardia Doosti, Ching-Hui Chen, Raviteja Vemulapalli, Xuhui Jia, Yukun Zhu, and Bradley Green. Boosting image-based mutual gaze detection using pseudo 3d gaze. In *Thirty-Fifth AAAI Conference on Artificial Intelligence*, pages 1273–1281, 2021. 8
- [53] Hazel Doughty, Ivan Laptev, Walterio Mayol-Cuevas, and Dima Damen. Action modifiers: Learning from adverbs in instructional videos. *arXiv preprint arXiv:1912.06617*, 2019. 42
- [54] Matteo Dunnhofer, Antonino Furnari, Giovanni Maria Farinella, and Christian Micheloni. Is first person vision challenging for object tracking? In *IEEE/CVF International Conference on Computer Vision Workshops (ICCVW) - Visual Object Tracking Challenge*, 2021. 3
- [55] Ariel Ephrat, Inbar Mosseri, Oran Lang, Tali Dekel, Kevin Wilson, Avinatan Hassidim, William T Freeman, and Michael Rubinstein. Looking to listen at the cocktail party: A speaker-independent audio-visual model for speech separation. In *SIGGRAPH*, 2018. 8, 47
- [56] Dave Epstein, Boyuan Chen, and Carl Vondrick. Oops! predicting unintentional action in video. In *Arxiv*, 2019. 41
- [57] N. Ryant et. al. The Second DIHARD Diarization Challenge: Dataset, task, and baselines. In *Proceedings of Interspeech*, 2019. 51
- [58] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International journal of computer vision*, 88(2):303–338, 2010. 1, 67
- [59] Bernard Ghanem Fabian Caba Heilbron, Victor Escorcia and Juan Carlos Niebles. Activitynet: A large-scale video benchmark for human activity understanding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 961–970, 2015. 1, 3, 30
- [60] Heng Fan, Haibin Ling, Liting Lin, Fan Yang, Peng Chu, Ge Deng, Sijia Yu, Hexin Bai, Yong Xu, and Chunyuan Liao. LaSOT: A High-Quality Benchmark for Large-Scale Single Object Tracking. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5369–5378, Long Beach, CA, USA, June 2019. IEEE. 33
- [61] Haoqi Fan, Bo Xiong, Karttikeya Mangalam, Yang-hao Li, Zhicheng Yan, Jitendra Malik, and Christoph Feichtenhofer. Multiscale vision transformers. *arXiv preprint arXiv:2104.11227*, 2021. 71
- [62] Yi Fang, Jiapeng Tang, Wang Shen, Wei Shen, Xiao Gu, Li Song, and Guangtao Zhai. Dual Attention Guided Gaze Target Detection in the Wild. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR 21)*, 2021. 8
- [63] Alireza Fathi, Jessica K. Hodgins, and James M. Rehg. Social interactions: A first-person perspective. In *CVPR*, 2012. 3
- [64] A. Fathi, J. K. Hodgins, and J. M. Rehg. Social interactions: A first-person perspective. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR 12)*, pages 1226–1233. IEEE, jun 2012. 8
- [65] A. Fathi and J. Rehg. Modeling actions through state changes. In *CVPR*, 2013. 7
- [66] Alireza Fathi and James M Rehg. Modeling actions through state changes. In *CVPR*, 2013. 41, 42
- [67] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for video recognition. In *ICCV*, 2019. 3, 5, 37, 45
- [68] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for video recognition. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6202–6211, 2019. 37, 70, 71
- [69] Jianglin Fu, Ivan V Bajić, and Rodney G Vaughan. Datasets for face and object detection in fisheye images. *Data in brief*, 27:104752, 2019. 53
- [70] Antonino Furnari, Sebastiano Battiatto, Kristen Grau-

- man, and Giovanni Maria Farinella. Next-active-object prediction from egocentric videos. *Journal of Visual Communication and Image Representation*, 49:401–411, 2017. 9
- [71] Antonino Furnari and Giovanni Farinella. Rolling-unrolling lstms for action anticipation from first-person video. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020. 3
- [72] Antonino Furnari and Giovanni Maria Farinella. What would you expect? anticipating egocentric actions with rolling-unrolling lstms and modality attention. In *International Conference on Computer Vision*, 2019. 9
- [73] Jiyang Gao, Zhenheng Yang, and Ram Nevatia. Red: Reinforced encoder-decoder networks for action anticipation. *BMVC*, 2017. 9
- [74] R. Gao, R. Feris, and K. Grauman. Learning to separate object sounds by watching unlabeled video. In *ECCV*, 2018. 8
- [75] Ruohan Gao, Rogerio Feris, and Kristen Grauman. Learning to separate object sounds by watching unlabeled video. In *ECCV*, 2018. 47
- [76] Ruohan Gao and Kristen Grauman. 2.5d visual sound. In *CVPR*, 2019. 8, 47
- [77] Ruohan Gao and Kristen Grauman. Co-separating sounds of visual objects. In *ICCV*, 2019. 47
- [78] R. Gao and K. Grauman. VisualVoice: Audio-visual speech separation with cross-modal consistency. In *CVPR*, 2021. 8, 47
- [79] I. Gebru, S. Ba, X. Li, and R. Horaud. Audio-visual speaker diarization based on spatiotemporal bayesian fusion. *PAMI*, 2018. 8, 47
- [80] Israel D. Gebru, Silèye Ba, Xiaofei Li, and Radu Horaud. Audio-visual speaker diarization based on spatiotemporal bayesian fusion. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39, 2017. 48
- [81] Georgios Georgakis, Md Alimoor Reza, Arsalan Mousavian, Phi-Hung Le, and Jana Košecká. Multi-view rgb-d dataset for object instance detection. In *2016 Fourth International Conference on 3D Vision (3DV)*, pages 426–434. IEEE, 2016. 7
- [82] Rohit Girdhar and Kristen Grauman. Anticipative video transformer. In *ICCV*, 2021. 3, 9
- [83] Ross Girshick. Fast r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 1440–1448, 2015. 70
- [84] Georgia Gkioxari and Jitendra Malik. Finding action tubes. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 759–768, Boston, MA, USA, June 2015. IEEE. 30
- [85] Raghav Goyal, Samira Ebrahimi Kahou, Vincent Michalski, Joanna Materzynska, Susanne Westphal, Heuna Kim, Valentin Haenel, Ingo Fruend, Peter Yianilos, Moritz Mueller-Freitag, et al. The “something something” video database for learning and evaluating visual common sense. In *ICCV*, 2017. 42
- [86] Alex Graves, Santiago Fernández, and Jürgen Schmidhuber. Bidirectional lstm networks for improved phoneme classification and recognition. In *International conference on artificial neural networks*, pages 799–804. Springer, 2005. 45
- [87] Chunhui Gu, Chen Sun, David A Ross, Carl Vondrick, Caroline Pantofaru, Yeqing Li, Sudheendra Vijayanarasimhan, George Toderici, Susanna Ricco, Rahul Sukthankar, et al. Ava: A video dataset of spatio-temporally localized atomic visual actions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6047–6056, 2018. 1, 3, 66
- [88] Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, Yu Zhang, Jiahui Yu, Wei Han, Shibo Wang, Zhengdong Zhang, Yonghui Wu, et al. Conformer: Convolution-augmented transformer for speech recognition. *arXiv preprint arXiv:2005.08100*, 2020. 57
- [89] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask R-CNN. *arXiv:1703.06870 [cs]*, Jan. 2018. 31
- [90] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 44, 45, 46
- [91] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 52
- [92] Farnoosh Heidarivincheh, Majid Mirmehdi, and Dima Damen. Detecting the moment of completion: Temporal models for localising action completion. In *BMVC*, 2018. 41
- [93] Matthew Honnibal, Ines Montani, Sofie Van Landeghem, and Adriane Boyd. spaCy: Industrial-strength Natural Language Processing in Python, 2020. 19
- [94] Lianghua Huang, Xin Zhao, and Kaiqi Huang. GOT-10k: A Large High-Diversity Benchmark for Generic Object Tracking in the Wild. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43(5):1562–1577, May 2021. 33
- [95] Noureddien Hussein, Efstratios Gavves, and Arnold WM Smeulders. Timeception for complex action recognition. In *CVPR*, 2019. 7
- [96] Go Irie, Mirela Ostrek, Haochen Wang, Hirokazu Kameoka, Akisato Kimura, Takahito Kawanishi, and Kunio Kashino. Seeing through sounds: Predicting visual semantic segmentation results from multichannel audio signals. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 3961–3964. IEEE, 2019. 47

- [97] Phillip Isola, Joseph J. Lim, and Edward H. Adelson. Discovering states and transformations in image collections. In *CVPR*, 2015. 7
- [98] Phillip Isola, Joseph J Lim, and Edward H Adelson. Discovering states and transformations in image collections. In *CVPR*, 2015. 42
- [99] Koji Iwano, Tomoaki Yoshinaga, Satoshi Tamura, and Sadaoki Furui. Audio-visual speech recognition using lip information extracted from side-face images. *EURASIP Journal on Audio, Speech, and Music Processing*, 2007:1–9, 2007. 8, 47
- [100] Andrew Jaegle, Felix Gimeno, Andrew Brock, Andrew Zisserman, Oriol Vinyals, and Joao Carreira. Perceiver: General perception with iterative attention. *arXiv preprint arXiv:2103.03206*, 2021. 45
- [101] Hao Jiang and Kristen Grauman. Seeing invisible poses: Estimating 3d body pose from egocentric video. In *CVPR*, 2017. 3
- [102] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, et al. The kinetics human action video dataset. *arXiv preprint arXiv:1705.06950*, 2017. 1, 3, 4, 37
- [103] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, et al. The kinetics human action video dataset. *arXiv preprint arXiv:1705.06950*, 2017. 71
- [104] Evangelos Kazakos, Arsha Nagrani, Andrew Zisserman, and Dima Damen. Epic-fusion: Audio-visual temporal binding for egocentric action recognition. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5492–5501, 2019. 3, 7, 8, 47
- [105] Petr Kellnhofer, Simon Stent, Wojciech Matusik, and Antonio Torralba. Gaze360: Physically Unconstrained Gaze Estimation in the Wild. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV 19)*, 2019. 8, 58, 59, 61
- [106] Suyoun Kim, Takaaki Hori, and Shinji Watanabe. Joint ctc-attention based end-to-end speech recognition using multi-task learning. In *2017 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pages 4835–4839. IEEE, 2017. 57
- [107] Kris M. Kitani, Brian Ziebart, James D. Bagnell, and Martial Hebert. Activity forecasting. In *ECCV*, 2012. 9
- [108] Dietrich Klakow and Jochen Peters. Testing the correlation of word error rate and perplexity. *Speech Communication*, 38(1-2):19–28, 2002. 8, 50
- [109] Mark L. Knapp, Judith A. Hall, and Terrence G. Horgan. *Nonverbal Communication in Human Interaction*. Wadsworth Cengage Learning, 8th edition, 2014. 8
- [110] Ross A Knepper, Todd Layton, John Romanishin, and Daniela Rus. Ikeabot: An autonomous multi-robot coordinated furniture assembly system. In *2013 IEEE International conference on robotics and automation*, pages 855–862. IEEE, 2013. 41
- [111] Andrew J Kolarik, Brian CJ Moore, Pavel Zahorik, Silvia Cirstea, and Shahina Pardhan. Auditory distance perception in humans: a review of cues, development, neuronal bases, and effects of sensory loss. *Attention, Perception, & Psychophysics*, 78(2):373–395, 2016. 47
- [112] Hema S. Koppula and Ashutosh Saxena. Anticipating human activities using object affordances for reactive robotic response. *Pattern Analysis and Machine Intelligence*, 38(1):14–29, 2016. 9
- [113] Ranjay Krishna, Kenji Hata, Frederic Ren, Li Fei-Fei, and Juan Carlos Niebles. Dense-captioning events in videos. In *International Conference on Computer Vision (ICCV)*, 2017. 7
- [114] Alexei A. Efros Krishna Kumar Singh, Kayvon Fatahalian. Krishnacam: Using a longitudinal, single-person, egocentric dataset for scene understanding tasks. In *IEEE Winter Conference on Applications of Computer Vision (WACV)*, 2016. 9
- [115] Matej Kristan, Ales Leonardis, Jiri Matas, Michael Felsberg, Roman Pflugfelder, Joni-Kristian Kamarainen, Luka Čehovin Zajc, Martin Danelljan, Alan Lukezic, Ondrej Drbohlav, Linbo He, Yushan Zhang, Song Yan, Jinyu Yang, Gustavo Fernandez, and et al. The eighth visual object tracking VOT2020 challenge results, 2020. 30
- [116] Taku Kudo and John Richardson. Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. *arXiv preprint arXiv:1808.06226*, 2018. 57
- [117] Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Mallocci, Alexander Kolesnikov, et al. The open images dataset v4. *International Journal of Computer Vision*, 128(7):1956–1981, 2020. 53
- [118] F. De la Torre, J. Hodgins, J. Montano, S. Valcarcel, R. Forcada, and J. Macey. Guide to the carnegie mellon university multimodal activity (cmu-mmact) database. In *Tech. report CMU-RI-TR-08-22, Robotics Institute, Carnegie Mellon University*, 2009. 3
- [119] Loic Lachèze, Yan Guo, Ryad Benosman, Bruno Gas, and Charlie Couverture. Audio/video fusion for objects recognition. In *2009 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 652–657. IEEE, 2009. 8
- [120] Kevin Lai, Liefeng Bo, and Dieter Fox. Unsupervised feature learning for 3d scene labeling. In *2014 IEEE International Conference on Robotics and Automation*

- (ICRA), pages 3050–3057. IEEE, 2014. 7
- [121] Tian Lan, Tsung-Chuan Chen, and Silvio Savarese. A hierarchical representation for future action prediction. In *ECCV*, 2014. 9
- [122] Federico Landini, Ján Profant, Mireia Diez, and Lukáš Burget. Bayesian hmm clustering of x-vector sequences (vbx) in speaker diarization: theory, implementation and analysis on standard tasks. *Computer Speech & Language*, 71:101254, 2022. 51
- [123] Y. J. Lee, J. Ghosh, and K. Grauman. Discovering important people and objects for egocentric video summarization. In *CVPR*, 2012. 2, 3
- [124] Y. J. Lee, J. Ghosh, and K. Grauman. Discovering important people and objects for egocentric video summarization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012. 42
- [125] Yong Jae Lee and Kristen Grauman. Predicting important objects for egocentric video summarization. *IJCV*, 2015. 2, 3
- [126] Bruno Lepri, Ramanathan Subramanian, Kyriaki Kalimeri, Jacopo Staiano, Fabio Pianesi, and Nicu Sebe. Connecting meeting behavior with extraversion—a systematic study. *IEEE Transactions on Affective Computing*, 3(4):443–455, 2012. 8
- [127] Vladimir I Levenshtein et al. Binary codes capable of correcting deletions, insertions, and reversals. In *Soviet physics doklady*, 1966. 68
- [128] Cheng Li and Kris Kitani. Model recommendation with virtual probes for ego-centric hand detection. In *ICCV*, 2013. 3
- [129] Yin Li, Alireza Fathi, and James M. Rehg. Learning to predict gaze in egocentric video. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3216–3223, 2013. 60
- [130] Y. Li, M. Liu, and J. Rehg. In the eye of beholder: Joint learning of gaze and actions in first person video. In *ECCV*, 2018. 42
- [131] Yin Li, Miao Liu, and Jame Rehg. In the Eye of the Beholder: Gaze and Actions in First Person Video. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021. 60
- [132] Yin Li, Miao Liu, and James M Rehg. In the eye of beholder: Joint learning of gaze and actions in first person video. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 619–635, 2018. 2, 3
- [133] Yanghao Li, Tushar Nagarajan, Bo Xiong, and Kristen Grauman. Ego-exo: Transferring visual representations from third-person to first-person videos. In *CVPR*, 2021. 3, 7
- [134] Tianwei Lin, Xiao Liu, Xin Li, Errui Ding, and Shilei Wen. Bmn: Boundary-matching network for temporal action proposal generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3889–3898, 2019. 44, 45
- [135] Tianwei Lin, Xu Zhao, Haisheng Su, Chongjing Wang, and Ming Yang. Bsn: Boundary sensitive network for temporal action proposal generation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 3–19, 2018. 7, 23, 30, 37
- [136] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature Pyramid Networks for Object Detection. *arXiv:1612.03144 [cs]*, Apr. 2017. 31
- [137] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft COCO: Common objects in context. In *ECCV*, 2014. 1, 3, 43
- [138] Miao Liu, Siyu Tang, Yin Li, and James M Rehg. Forecasting human-object interaction: joint prediction of motor attention and actions in first person video. In *ECCV*, 2020. 3
- [139] Wen Liu, Weixin Luo, Dongze Lian, and Shenghua Gao. Future frame prediction for anomaly detection—a new baseline. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6536–6545, 2018. 9
- [140] William Lotter, Gabriel Kreiman, and David Cox. Deep predictive coding networks for video prediction and unsupervised learning. *arXiv preprint arXiv:1605.08104*, 2016. 9
- [141] Cewu Lu, Renjie Liao, and Jiaya Jia. Personal object discovery in first-person videos. *TIP*, 2015. 3
- [142] Zheng Lu and Kristen Grauman. Story-driven summarization for egocentric video. In *CVPR*, 2013. 3
- [143] Tahmida Mahmud, Mahmudul Hasan, and Amit K Roy-Chowdhury. Joint prediction of activity labels and starting times in untrimmed videos. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5773–5782, 2017. 9
- [144] Manuel J Marin-Jimenez, Vicky Kalogeiton, Pablo Medina-Suarez, and Andrew Zisserman. Laeo-net: revisiting people looking at each other in videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3477–3485, 2019. 8
- [145] Manuel Jesús Marín-Jiménez, Andrew Zisserman, Marcin Eichner, and Vittorio Ferrari. Detecting people looking at each other in videos. *International Journal of Computer Vision*, 106(3):282–296, 2014.
- [146] Manuel J Marín-Jiménez, Andrew Zisserman, and Vittorio Ferrari. Here's looking at you, kid. *Detecting people looking at each other in videos*. In *BMVC*, 5, 2011. 8
- [147] Michael Mathieu, Camille Couprie, and Yann LeCun. Deep multi-scale video prediction beyond mean

- square error. *arXiv preprint arXiv:1511.05440*, 2015. 9
- [148] Iain McCowan, Jean Carletta, Wessel Kraaij, Simone Ashby, Sebastien Bourban, Mike Flynn, Mael Guillemot, Thomas Hain, Jaroslav Kadlec, Vasilis Karaikos, Melissa Kronenthal, Guillaume Lathoud, Mike Lincoln, Agnes Lisowska, Wilfried Post, Dennis Reidsma, and Pierre Wellner. The AMI meeting corpus. In *Proceedings of Measuring Behavior 2005, the 5th International Conference on Methods and Techniques in Behavioral Research*, pages 137–140, 2005. 8
- [149] Jean-Philippe Mercier, Mathieu Garon, Philippe Giguere, and Jean-Francois Lalonde. Deep template-based object instance detection. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 1507–1516, January 2021. 7
- [150] Christophe Micheyl, Christian Kaernbach, and Laurent Demany. An evaluation of psychophysical models of auditory change perception. *Psychological review*, 115(4):1069, 2008. 47
- [151] Antoine Miech, Dimitri Zhukov, Jean-Baptiste Alayrac, Makarand Tapaswi, Ivan Laptev, and Josef Sivic. HowTo100M: Learning a Text-Video Embedding by Watching Hundred Million Narrated Video Clips. In *ICCV*, 2019. 3, 5
- [152] Ishan Misra, Abhinav Gupta, and Martial Hebert. From red wine to red tomato: Composition with context. In *CVPR*, 2017. 42
- [153] Yu Mitsuzumi, Atsushi Nakazawa, and Toyoaki Nishida. Deep eye contact detector: Robust eye contact bid detection using convolutional neural network. In *BMVC*, 2017. 8
- [154] Davide Moltisanti, Michael Wray, Walterio Mayol-Cuevas, and Dima Damen. Trespassing the boundaries: Labelling temporal bounds for object interactions in egocentric video. In *ICCV*, 2017. 41
- [155] Pedro Morgado, Nono Vasconcelos, Timothy Langlois, and Oliver Wang. Self-supervised generation of spatial audio for 360° video. In *NeurIPS*, 2018. 8, 47
- [156] Matthias Müller, Adel Bibi, Silvio Giancola, Salman Alsubaihi, and Bernard Ghanem. TrackingNet: A Large-Scale Dataset and Benchmark for Object Tracking in the Wild. In Vittorio Ferrari, Martial Hebert, Cristian Sminchisescu, and Yair Weiss, editors, *Computer Vision – ECCV 2018*, volume 11205, pages 310–327. Springer International Publishing, Cham, 2018. 33
- [157] Tushar Nagarajan, Christoph Feichtenhofer, and Kristen Grauman. Grounded human-object interaction hotspots from video. *ICCV*, 2019. 3
- [158] Tushar Nagarajan and Kristen Grauman. Attributes as operators: factorizing unseen attribute-object compositions. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 169–185, 2018. 7, 42
- [159] A. Nagrani, J. S. Chung, and A. Zisserman. Vox-Celeb: a large-scale speaker identification dataset. In *INTERSPEECH*, 2017. 48
- [160] Katsuyuki Nakamura, Serena Yeung, Alexandre Alahi, and Li Fei-Fei. Jointly learning energy expenditures and activities using egocentric multimodal signals. In *CVPR*, 2017. 3
- [161] Lukas Neumann, Andrew Zisserman, and Andrea Vedaldi. Future event prediction: If and when. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 0–0, 2019. 9
- [162] Evonne Ng, Donglai Xiang, Hanbyul Joo, and Kristen Grauman. You2me: Inferring body pose in egocentric video via first and second person interactions. In *CVPR*, 2020. 3
- [163] Joonas Nikunen and Tuomas Virtanen. Direction of arrival based spatial covariance model for blind sound source separation. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 22(3):727–739, 2014. 47
- [164] C. Northcutt, S. Zha, S. Lovegrove, and R. Newcombe. Egocom: A multi-person multi-modal egocentric communications dataset. *PAMI*, 2020. 3
- [165] Andrew Owens and Alexei A Efros. Audio-visual scene analysis with self-supervised multisensory features. In *ECCV*, 2018. 47
- [166] Cristina Palmero, Elsbeth A van Dam, Sergio Escalera, Mike Kelia, Guido F Lichtert, Lucas PJJ Noldus, Andrew J Spink, and Astrid van Wieringen. Automatic mutual gaze detection in face-to-face dyadic interaction videos. *Measuring Behavior 2018*, 2018. 8
- [167] Daniel S Park, William Chan, Yu Zhang, Chung-Cheng Chiu, Barret Zoph, Ekin D Cubuk, and Quoc V Le. SpecAugment: A simple data augmentation method for automatic speech recognition. *arXiv preprint arXiv:1904.08779*, 2019. 56
- [168] H. S. Park, J.-J. Hwang, Y. Niu, and J. Shi. Egocentric future localization. In *CVPR*, 2016. 9
- [169] H. S. Park, J.-J. Hwang, Y. Niu, and J. Shi. Egocentric future localization. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 68
- [170] Hyun Soo Park, Eakta Jain, and Yaser Sheikh. 3D social saliency from head-mounted cameras. In *Advances in Neural Information Processing Systems*, volume 1, pages 422–430, 2012. 8
- [171] Tae Jin Park, Naoyuki Kanda, Dimitrios Dimitriadis, Kyu J Han, Shinji Watanabe, and Shrikanth Narayanan. A review of speaker diarization: Re-

- cent advances with deep learning. *arXiv preprint arXiv:2101.09624*, 2021. 49, 51
- [172] David R Perrott and Kourosh Saberi. Minimum audible angle thresholds for sources varying in both elevation and azimuth. *The Journal of the Acoustical Society of America*, 87(4):1728–1731, 1990. 47
- [173] Hamed Pirsiavash and Deva Ramanan. Detecting activities of daily living in first-person camera views. In *2012 IEEE conference on computer vision and pattern recognition*, pages 2847–2854. IEEE, 2012. 2, 3
- [174] H. Pirsiavash and D. Ramanan. Detecting activities of daily living in first-person camera views. In *Computer Vision and Pattern Recognition (CVPR)*, 2012. 42
- [175] Daniel Povey, Arnab Ghoshal, Gilles Boulianne, Lukas Burget, Ondrej Glembek, Nagendra Goel, Mirko Hannemann, Petr Motlicek, Yanmin Qian, Petr Schwarz, Jan Silovsky, Georg Stemmer, and Karel Vesely. The Kaldi speech recognition toolkit. In *IEEE 2011 Workshop on Automatic Speech Recognition and Understanding*, 2011. 51
- [176] Senthil Purushwalkam, Maximilian Nickel, Abhinav Gupta, and Marc’ Aurelio Ranzato. Task-driven modular networks for zero-shot compositional learning. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3593–3602, 2019. 42
- [177] F. Ragusa, A. Furnari, S. Battiatto, G. Signorello, and G. M. Farinella. Egocentric visitors localization in cultural sites. *Journal on Computing and Cultural Heritage (JOCCH)*, 2019. 3
- [178] Francesco Ragusa, Antonino Furnari, Salvatore Livotino, and Giovanni Maria Farinella. The meccano dataset: Understanding human-object interactions from egocentric videos in an industrial-like domain. In *IEEE Winter Conference on Application of Computer Vision (WACV)*, 2021. 3
- [179] René Ranftl, Alexey Bochkovskiy, and Vladlen Koltun. Vision transformers for dense prediction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12179–12188, 2021. 35, 36
- [180] Adria Recasens, Aditya Khosla, Carl Vondrick, and Antonio Torralba. Where are they looking? In *Advances in Neural Information Processing Systems*, pages 199–207, 2015. 8
- [181] Joseph Redmon and Ali Farhadi. Yolov3: An incremental improvement. *arXiv preprint arXiv:1804.02767*, 2018. 53
- [182] James M. Rehg, Agata Rozga, Gregory D. Abowd, and Matthew S. Goodwin. Behavioral Imaging and Autism. *IEEE Pervasive Computing*, 13(2):84–87, 2014. 8
- [183] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *NeurIPS*, 2015. 31
- [184] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in neural information processing systems*, 28:91–99, 2015. 45, 46
- [185] Ivan Rodin, Antonino Furnari, Dimitrios Mavrodis, and Giovanni Maria Farinella. Predicting the future from first person (egocentric) vision: A survey. *Computer Vision and Image Understanding*, 2021. 9
- [186] Joseph Roth, Sourish Chaudhuri, Ondrej Klejch, Radhika Marvin, Andrew Gallagher, Liat Kaver, Sharadh Ramaswamy, Arkadiusz Stopczynski, Cordelia Schmid, Zhonghua Xi, et al. Avactivespeaker: An audio-visual dataset for active speaker detection. *arXiv preprint arXiv:1901.01342*, 2019. 8, 48, 49
- [187] Joseph Roth, Sourish Chaudhuri, Ondrej Klejch, Radhika Marvin, Andrew Gallagher, Liat Kaver, Sharadh Ramaswamy, Arkadiusz Stopczynski, Cordelia Schmid, Zhonghua Xi, and Caroline Pantofaru. Ava Active Speaker: An Audio-Visual Dataset for Active Speaker Detection. In *ICASSP, IEEE International Conference on Acoustics, Speech and Signal Processing - Proceedings*, volume 2020-May, pages 4492–4496, 2020. 8
- [188] M. S. Ryoo and L. Matthies. First-person activity recognition: What are they doing to me? In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2013. 3
- [189] Paul-Edouard Sarlin, Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Superglue: Learning feature matching with graph neural networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4938–4947, 2020. 34, 35
- [190] Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4104–4113, 2016. 34
- [191] A. Senocak, T.-H. Oh, J. Kim, M. Yang, and I. S. Kweon. Learning to localize sound sources in visual scenes: Analysis and applications. *TPAMI*, 2019. 8, 47
- [192] Dandan Shan, Jiaqi Geng, Michelle Shu, and David Fouhey. Understanding human hands in contact at internet scale. In *CVPR*, 2020. 42, 43
- [193] Dandan Shan, Jiaqi Geng, Michelle Shu, and David Fouhey. Understanding human hands in contact at internet scale. In *CVPR*, 2020. 45, 46
- [194] Mohit Sharma, Kevin Zhang, and Oliver Kroemer. Learning semantic embedding spaces for slicing vegetables. *arXiv preprint arXiv:1904.00303*, 2019. 41
- [195] Gunnar A Sigurdsson, Abhinav Gupta, Cordelia

- Schmid, Ali Farhadi, and Karteek Alahari. Charadesego: A large-scale dataset of paired third and first person videos. *arXiv preprint arXiv:1804.09626*, 2018. 2, 3, 42
- [196] Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob Fergus. Indoor segmentation and support inference from rgbd images. In *European conference on computer vision*, pages 746–760. Springer, 2012. 35
- [197] Michel Silva, Washington Ramos, João Ferreira, Felipe Chamone, Mario Campos, and Erickson R. Nascimento. A weighted sparse sampling and smoothing frame transition approach for semantic fast-forward first-person videos. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 3
- [198] Krishna Kumar Singh, Kayvon Fatahalian, and Alexei A Efros. Krishnacam: Using a longitudinal, single-person, egocentric dataset for scene understanding tasks. In *WACV*, 2016. 2, 3
- [199] David Snyder, Daniel Garcia-Romero, Gregory Sell, Daniel Povey, and Sanjeev Khudanpur. X-vectors: Robust DNN embeddings for speaker recognition. In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2018. 51
- [200] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. Ucf101: A dataset of 101 human action classes from videos in the wild. In *CRCV-TR-12-01*, 2012. 3
- [201] Emiliano Spera, Antonino Furnari, Sebastiano Battatito, and Giovanni Maria Farinella. Egocentric shopping cart localization. In *International Conference on Pattern Recognition (ICPR)*, 2018. 3
- [202] Julian Straub, Thomas Whelan, Lingni Ma, Yufan Chen, Erik Wijmans, Simon Green, Jakob J Engel, Raul Mur-Artal, Carl Ren, Shobhit Verma, et al. The replica dataset: A digital replica of indoor spaces. *arXiv preprint arXiv:1906.05797*, 2019. 14
- [203] Yu-Chuan Su and Kristen Grauman. Detecting engagement in egocentric video. In *ECCV*, 2016. 2, 3, 42
- [204] Ruijie Tao, Zexu Pan, Rohan Kumar Das, Xinyuan Qian, Mike Zheng Shou, and Haizhou Li. Is someone speaking? exploring long-term temporal features for audio-visual active speaker detection. *arXiv preprint arXiv:2107.06592*, 2021. 49, 54, 55, 56
- [205] Y. Tian, J. Shi, B. Li, Z. Duan, and C. Xu. Audio-visual event localization in unconstrained videos. In *ECCV*, 2018. 8, 47
- [206] E. Tulving. Episodic and semantic memory. In E. Tulving and W. Donaldson, editors, *Organization of memory*. Academic Press, 1972. 6
- [207] TwentyBN. The 20BN-jester Dataset V1. <https://20bn.com/datasets/jester>. 42
- [208] Joost Van Amersfoort, Anitha Kannan, Marc’Aurelio Ranzato, Arthur Szlam, Du Tran, and Soumith Chintala. Transformation-based models of video sequences. *arXiv preprint arXiv:1701.08435*, 2017. 9, 37
- [209] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017. 46
- [210] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017. 57
- [211] Ruben Villegas, Jimei Yang, Seunghoon Hong, Xunyu Lin, and Honglak Lee. Decomposing motion and content for natural video sequence prediction. *arXiv preprint arXiv:1706.08033*, 2017. 9
- [212] Carl Vondrick, Hamed Pirsiavash, and Antonio Torralba. Anticipating visual representations from unlabeled video. In *CVPR*, 2016. 9
- [213] He Wang, Sören Pirk, Ersin Yumer, Vladimir G Kim, Ozan Sener, Srinath Sridhar, and Leonidas J Guibas. Learning a generative model for multi-step human-object interactions from videos. In *Eurographics*, 2019. 42
- [214] Limin Wang, Yuanjun Xiong, Zhe Wang, Yu Qiao, Dahua Lin, Xiaoou Tang, and Luc Van Gool. Temporal segment networks: Towards good practices for deep action recognition. In *ECCV*, 2016. 3, 7
- [215] Qiang Wang, Li Zhang, Luca Bertinetto, Weiming Hu, and Philip H. S. Torr. Fast online object tracking and segmentation: A unifying approach, 2019. 16
- [216] Xiaolong Wang, Ali Farhadi, and Abhinav Gupta. Actions~ transformations. In *CVPR*, 2016. 42
- [217] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *CVPR*, 2018. 3
- [218] Yuxin Wu, Alexander Kirillov, Francisco Massa, Wan-Yen Lo, and Ross Girshick. Detectron2. 33
- [219] Fanyi Xiao, Yong Jae Lee, Kristen Grauman, Jitendra Malik, and Christoph Feichtenhofer. Audiovisual slowfast networks for video recognition. 2020. 8, 47
- [220] SHI Xingjian, Zhourong Chen, Hao Wang, Dit-Yan Yeung, Wai-Kin Wong, and Wang-chun Woo. Convolutional lstm network: A machine learning approach for precipitation nowcasting. In *Advances in neural information processing systems*, pages 802–810, 2015. 9
- [221] Jun Xu, Tao Mei, Ting Yao, and Yong Rui. Msrvtt: A large video description dataset for bridging video and language. IEEE International Conference on Computer Vision and Pattern Recognition (CVPR), June 2016. 7

- [222] Mengmeng Xu, Chen Zhao, David S Rojas, Ali Thabet, and Bernard Ghanem. G-tad: Sub-graph localization for temporal action detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10156–10165, 2020. [7](#), [23](#), [30](#), [37](#)
- [223] Takuma Yagi, Karttikeya Mangalam, Ryo Yonetani, and Yoichi Sato. Future person localization in first-person videos. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. [9](#)
- [224] Ryo Yonetani, Kris M. Kitani, and Yoichi Sato. Recognizing micro-actions and reactions from paired ego-centric videos. In *CVPR*, 2016. [3](#)
- [225] Ryo Yonetani, Kris M Kitani, and Yoichi Sato. Visual motif discovery via first-person vision. In *ECCV*, 2016. [3](#)
- [226] Fisher Yu, Dequan Wang, Evan Shelhamer, and Trevor Darrell. Deep layer aggregation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2403–2412, 2018. [46](#)
- [227] Hua Zhang, Xiaochun Cao, and Rui Wang. Audio visual attribute discovery for fine-grained object recognition. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018. [8](#)
- [228] Hao Zhang, Aixin Sun, Wei Jing, and Joey Tianyi Zhou. Span-based localizing network for natural language video localization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6543–6554, Online, July 2020. Association for Computational Linguistics. [37](#)
- [229] Songyang Zhang, Houwen Peng, Jianlong Fu, and Jiebo Luo. Learning 2d temporal adjacent networks for moment localization with natural language. In *AAAI*, 2020. [24](#), [30](#), [36](#)
- [230] Chen Zhao, Ali Thabet, and Bernard Ghanem. Video self-stitching graph network for temporal action localization. *arXiv preprint arXiv:2011.14598*, 2020. [7](#), [23](#), [30](#), [37](#), [39](#)
- [231] Hang Zhao, Chuang Gan, Andrew Rouditchenko, Carl Vondrick, Josh McDermott, and Antonio Torralba. The sound of pixels. In *ECCV*, 2018. [47](#)
- [232] Bolei Zhou, Alex Andonian, Aude Oliva, and Antonio Torralba. Temporal relational reasoning in videos. In *ECCV*, 2018. [3](#)
- [233] Hao Zhou, Chongyang Zhang, Yan Luo, Yanjun Chen, and Chuiping Hu. Embracing uncertainty: Decoupling and de-bias for robust temporal grounding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8445–8454, 2021. [31](#)
- [234] Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl. Objects as points. *arXiv preprint arXiv:1904.07850*, 2019. [45](#), [46](#)
- [235] Y. Zhou and T. Berg. Learning temporal transforma-  
tions from time-lapse videos. In *ECCV*, 2016. [7](#)
- [236] Yipin Zhou and Tamara L Berg. Temporal perception and prediction in ego-centric video. In *ICCV*, 2015. [3](#), [7](#)
- [237] Yipin Zhou and Tamara L Berg. Learning temporal transformations from time-lapse videos. In *ECCV*, 2016. [42](#)
- [238] Hao Zhu, Man-Di Luo, Rui Wang, Ai-Hua Zheng, and Ran He. Deep audio-visual learning: A survey. *International Journal of Automation and Computing*, pages 1–26, 2021. [8](#)