

Detecting Engagement in Egocentric Video

Yu-Chuan Su and Kristen Grauman

The University of Texas at Austin

Abstract. In a wearable camera video, we see what the camera wearer sees. While this makes it easy to know roughly *what he chose to look at*, it does not immediately reveal *when he was engaged with the environment*. Specifically, at what moments did his focus linger, as he paused to gather more information about something he saw? Knowing this answer would benefit various applications in video summarization and augmented reality, yet prior work focuses solely on the “what” question (estimating saliency, gaze) without considering the “when” (engagement). We propose a learning-based approach that uses long-term egomotion cues to detect engagement, specifically in browsing scenarios where one frequently takes in new visual information (e.g., shopping, touring). We introduce a large, richly annotated dataset for ego-engagement that is the first of its kind. Our approach outperforms a wide array of existing methods. We show engagement can be detected well independent of both scene appearance and the camera wearer’s identity.

1 Introduction

Imagine you are walking through a grocery store. You may be mindlessly plowing through the aisles grabbing your usual food staples, when a new product display—or an interesting fellow shopper—captures your interest for a few moments. Similarly, in the museum, as you wander the exhibits, occasionally your attention is heightened and you draw near to examine something more closely.

These examples illustrate the notion of *engagement* in ego-centric activity, where one pauses to inspect something more closely. While engagement happens throughout daily life activity, it occurs frequently and markedly during “*browsing*” scenarios in which one traverses an area with the intent of taking in new information and/or locating certain objects—for example, in a shop, museum, library, city sightseeing, or touring a campus or historic site.

Problem definition We explore engagement from the first-person vision perspective. In particular, we ask: Given a video stream captured from a head-mounted camera during a browsing scenario, can we automatically detect those time intervals where the recorder experienced a heightened level of engagement? What cues are indicative of first-person engagement, and how do they differ from traditional saliency metrics? To what extent are engagement cues independent of the particular person wearing the camera (the “recorder”), or the particular environment they are navigating? See Fig. 1.

While engagement is interesting in a variety of daily life settings, for now we restrict our focus to browsing scenarios. This allows us to concentrate on cases



Fig. 1: The goal is to identify time intervals when the camera wearer’s engagement is heightened, meaning he interrupts his ongoing activity to gather more information about some object in the environment. Note that this is different than detecting what the camera wearer sees or gazes upon, which comes for “free” with a head-mounted camera and/or eye tracking devices.

where 1) engagement naturally ebbs and flows repeatedly, 2) the environment offers discrete entities (products in the shop, museum paintings, etc.) that may be attention-provoking, which aids objectivity in evaluation, and 3) there is high potential impact for emerging applications.

Applications A system that can successfully address the above questions would open up several applications. For example, it could facilitate camera control, allowing the user’s attention to trigger automatic recording/zooming. Similarly, it would help construct video summaries. Knowing when a user’s engagement is waning would let a system display info on a heads-up display when it is least intrusive. Beyond such “user-centric” applications, third parties would relish the chance to gather data about user attention at scale—for instance, a vendor would like to know when shoppers linger by its new display. Such applications are gaining urgency as wearable cameras become increasingly attractive tools in the law enforcement, healthcare, education, and consumer domains.

Novelty of the problem The rich literature on visual saliency—including video saliency [1–8]—does not address this problem. First and foremost, as discussed above, detecting moments of engagement is different than estimating saliency. While a person always sees something, he does not pay attention to everything he sees; knowing *what* a person is looking at does not reveal *when* the person is engaging with the environment. Nearly all prior work studies visual saliency from the *third person* perspective and equates saliency with gaze: salient points are those upon which a viewer would fixate his gaze, when observing a previously recorded image/video on a static screen. In contrast, our problem entails detecting *temporal intervals of engagement as perceived by the person capturing the video as he moves about his environment*. Thus, *recorder engagement* is distinct from *viewer attention*. To predict it from video requires identifying time intervals of engagement as opposed to spatial regions that are salient (gaze worthy) per frame. As such, estimating egocentric gaze [9–11] is also insufficient to predict first-person engagement.

Challenges Predicting first-person engagement presents a number of challenges. First of all, the motion cues that are significant in third-person video taken with an actively controlled camera (e.g., zoom [4, 12–14]) are absent in

passive wearable camera data. Instead, first-person data contains both scene motion and unstable body motions, which are difficult to stabilize with traditional methods [15]. Secondly, whereas third-person data is inherently already focused on moments of interest that led the recorder to turn the camera on, a first-person camera is “always on”. Thirdly, whereas traditional visual attention metrics operate with instantaneous motion cues [1, 2, 16, 17] and fixed sliding temporal window search strategies, detecting engagement *intervals* requires long-term descriptors and handling intervals of variable length. Finally, it is unclear whether there are sufficient visual cues that transcend user- or scene-specific properties, or if engagement is strongly linked to the specific content a user observes (in which case, an exorbitant amount of data might be necessary to learn a general-purpose detector).

Our approach We propose a learning approach to detect time intervals where first-person engagement occurs. In an effort to maintain independence of the camera wearer as well as the details of his environment, we employ motion-based features that span long temporal neighborhoods and integrate out local head motion effects. We develop a search strategy that integrates instantaneous frame-level estimates with temporal interval hypotheses to detect intervals of varying lengths, thereby avoiding a naive sliding window search. To train and evaluate our model, we undertake a large-scale data collection effort.

Contributions Our main contributions are as follows. First, we precisely define egocentric engagement and systematically evaluate under that definition. Second, we collect a large annotated dataset spanning 14 hours of activity explicitly designed for ego-engagement in browsing situations. Third, we propose a learned motion-based model for detecting first-person engagement. Our model shows better accuracy than an array of existing methods. It also generalizes to unseen browsing scenarios, suggesting that some properties of ego-engagement are independent of appearance content.

2 Related Work

Third-person image and video saliency Researchers often equate human gaze fixations as the gold standard with which a *saliency* metric ought to correlate [18, 19]. There is increasing interest in estimating saliency from video. Initial efforts examine simple motion cues, such as frame-based motion and flicker [8, 18, 20]. One common approach to extend spatial (image) saliency to the video domain is to sum image saliency scores within a temporal segment, e.g., [21]. Most methods are unsupervised and entail no learning [4–8, 18, 20]. However, some recent work develops learned measures, using ground truth gaze data as the target output [1–3, 16, 22].

Our problem setting is quite different than saliency. Saliency aims to *predict viewer attention* in terms of where in the frame a third party is likely to fixate his gaze; it is an image property analyzed independent of the behavior of the person recording the image. In contrast, we aim to *detect recorder engagement* in terms of when (which time intervals) the recorder has paused to examine

something in his environment.¹ Accounting for this distinction is crucial, as we will see in results. Furthermore, prior work in video saliency is evaluated on short video clips (e.g., on the order of 10 seconds [23]), which is sufficient to study gaze movements. In contrast, we evaluate on long sequences—30 minutes on average per clip, and a total of 14 hours—in order to capture the broad context of ego-behavior that affects engagement in browsing scenarios.

Third-person video summarization In video summarization, the goal is to form a concise representation for a long input video. Motion cues can help detect “important” moments in third-person video [12–14, 17, 21], including temporal differences [17] and cues from active camera control [12–14]. Whereas prior methods try to extract what will be interesting to a third-party viewer, we aim to capture *recorder* engagement.

First-person video saliency and gaze Researchers have long expected that ego-attention detection requires methods distinct from bottom-up saliency [24]. In fact, traditional motion saliency can actually *degrade* gaze prediction for first-person video [11]. Instead, it is valuable to separate out camera motion [10] or use head motion and hand locations to predict gaze [9]. Whereas these methods aim to predict spatial coordinates of a recorder’s gaze at every frame, we aim to predict time intervals where his engagement is heightened. Furthermore, whereas they study short sequences in a lab [10] or kitchen [9], we analyze long data in natural environments with substantial scene changes per sequence.

We agree that first-person attention, construed in the most general sense, will inevitably require first-person “user-in-the-loop” feedback to detect [24]; accordingly, our work does not aim to detect arbitrary subjective attention events, but instead to detect moments of engagement to examine an object more closely.

Outside of gaze, there is limited work on attention in terms of head fixation detection [15] and “physical analytics” [25]. In [15], a novel “cumulative displacement curve” motion cue is used to categorize the recorder’s activity (walking, sitting, on bus, etc.) and is also shown to reveal periods with fixed head position. They use a limited definition of attention: a period of more than 5 seconds where the head is still but the recorder is walking. In [25], inertial sensors are used in concert with optical flow magnitude to decide when the recorder is examining a product in a store. Compared to both [15, 25], engagement has a broader definition, and we discover its scope from data from the crowd (vs. hand-crafting a definition on visual features). Crucially, the true positives reflect that a person can have heightened engagement yet still be in motion.

First-person activity and summarization Early methods for egocentric video summarization extract the camera motion and define rules for important moments (e.g., intervals when camera rotation is below a threshold) [26, 27], and test qualitatively on short videos. Rather than inject hand-crafted rules, we propose to *learn* what constitutes an engagement interval. Recent methods

¹ Throughout, we will use the term “recorder” to refer to the photographer or the first-person camera-wearer; we use the term “viewer” to refer to a third party who is observing the data captured by some other recorder.

explore ways to predict the “importance” of spatial regions (objects, people) using cues like hand detection and frame centrality [28, 29], detect novelty [30], and infer “social saliency” when multiple cameras capture the same event [31–33]. We tackle engagement, not summarization, though likely our predictions could be another useful input to a summarization system.

In a sense, detecting engagement could be seen as detecting a particular ego-activity. An array of methods for classifying activity in egocentric video exist, e.g., [34–41]. However, they do not address our scenario: 1) they learn models specific to the objects [34, 36–38, 40, 41] or scenes [39] with which the activity takes place (e.g., making tea, snowboarding), whereas engagement is by definition object- and scene-independent, since arbitrary things may capture one’s interest; and 2) they typically focus on recognition of trimmed video clips, versus temporal detection in ongoing video.

3 First-Person Engagement: Definition and Data

Next we define first-person engagement. Then we describe our data collection procedure, and quantitatively analyze the consistency of the resulting annotations. We introduce our approach for predicting engagement intervals in Sec. 4.

3.1 Definition of first-person engagement

This research direction depends crucially on having 1) a precise definition of engagement, 2) realistic video data captured in natural environments, and 3) a systematic way to annotate the data for both learning and evaluation.

Accordingly, we first formalize our meaning of first-person engagement. There are two major requirements. First, the engagement must be related to external factors, either induced by or causing the change in visual signals the recorder perceives. This ensures predictability from video, excluding high-attention events that are imperceptible (by humans) from visual cues. Second, an engagement interval must reflect the *recorder’s* intention, as opposed to the reaction of a third-person viewer of the same video.

Based on these requirements, we **define heightened ego-engagement in a browsing scenario** as follows. A time interval is considered to have a high engagement level if *the recorder is attracted by some object(s), and he interrupts his ongoing flow of activity to purposefully gather more information about the object(s)*. We stress that this definition is scoped specifically for *browsing* scenarios; while the particular objects attracting the recorder will vary widely, we assume the person is traversing some area with the intent of taking in new information and/or locating certain objects.

The definition captures situations where the recorder reaches out to touch or grasp an object of interest (e.g., when closely inspecting a product at the store), as well as scenarios where he examines something from afar (e.g., when he reads a sign beside a painting at the museum). Having an explicit definition allows annotators to consistently identify video clips with high engagement, and it lets us directly evaluate the prediction result of different models.

We stress that ego-engagement differs from gaze and traditional saliency. While a recorder always has a gaze point per frame (and it is correlated with the

	Mall	Market	Museum	All
Attention Ratio	0.305	0.451	0.580	0.438
#intervals (per min.)	1.19	1.22	1.50	1.30
Length Median (sec)	7.5	12.1	13.3	11.3
Length IQR (sec)	11.6	18.2	20.1	17.6

Table 1: Basic statistics for ground truth intervals.

frame center), periods of engagement are more sparsely distributed across time, occupy variable-length intervals, and are a function of his activity and changing environment. Furthermore, as we will see below, moments where a person is approximately still are *not* equivalent to moments of engagement, making observer motion magnitude [25] an inadequate signal.

3.2 Data collection

To collect a dataset, we ask multiple recorders to take videos during browsing behavior under a set of *scenarios*, or scene and event types. We aim to gather scenarios with clear distinctions between high and low engagement intervals that will be apparent to a third-party annotator. Based on that criterion, we collect videos under three scenarios: 1) shopping in a market, 2) window shopping in shopping mall, and 3) touring in a museum. All three entail spontaneous stimuli, which ensures that variable levels of engagement will naturally occur.

The videos are recorded using Looxcie LX2 with 640×480 resolution and 15 fps frame rate, which we chose for its long battery life and low profile. We recruited 9 recorders—5 females and 4 males—all students between 20-30 years old. Other than asking them to capture instances of the scenarios above, we did not otherwise instruct the recorders to behave in any way. Among the 9 recorders, 5 of them record videos in all 3 scenarios. The other 4 record videos in 2 scenarios. Altogether, we obtained 27 videos, each averaging 31 minutes, for a total dataset of 14 hours. To keep the recorder behavior as natural as possible, we asked the recorders to capture the video when they planned to go to such scenarios anyway; as such, it took about 1.5 months to collect the video.

After collecting the videos, we crowdsource the ground truth annotations on Amazon Mechanical Turk. Importantly, we ask annotators to put themselves in the camera-wearer’s shoes. They must precisely mark the start and end points of each engagement interval from the recorder’s perspective, and record their confidence.² We break the source videos into 3 minutes overlapping chunks to make each annotation task manageable yet still reveal temporal context for the clip. We estimate the annotations took about 450 worker-hours and cost \$3,000. Our collection strategy is congruous with the goals stated above in Sec. 3.1, in that annotators are shown only the visual signal (without audio) and are asked to consider engagement from the point of view of the recorder. See Supp. file.

Despite our care in the instructions, there remains room for annotator subjectivity, and the exact interval boundaries can be ambiguous. Thus, we ask 10

² For a portion of the video, we also ask the original recorders to label all frames for their own video; this requires substantial tedious effort, hence to get the full labeled set in a scalable manner we apply crowdsourcing.

Turkers to annotate each video. Positive intervals are those where a majority agree engagement is heightened. To avoid over-segmentation, we ignore intervals shorter than 1 second. For each positive interval, we select the tightest annotation that covers more than half of the interval as the final ground truth.

The resulting dataset contains examples that are diverse in content and duration. The recorders are attracted by a variety of objects: groceries, household items, clothes, paintings, sculptures, other people. In some cases, the attended object is out of the field of view, e.g., a recorder grabs an item without directly looking at it, in which case Turkers infer the engagement from context.

Table 1 summarizes some statistics of the labeled data. On average, the recorder is engaged about 44% of the time (see “Attention Ratio”), and it increases once to twice per minute. This density reflects the browsing scenarios on which we focus the data. The length of a positive interval varies substantially: the interquartile range (IQR) is 17.6 seconds, about 50% longer than the median. Some intervals last as long as 5 minutes. Also, the different scenarios have different statistics, e.g., Museum scenarios prompt more frequent engagement. All this variability indicates the difficulty of the task.

The new dataset is the first of its kind to explicitly define and thoroughly annotate ego-engagement. It is also substantially larger than datasets used in related areas—nearly 14 hours of video, with test videos over 30 minutes each. By contrast, clips in third-person saliency datasets are typically 20 seconds [23] to 2 minutes [42], since the interest is in gauging instantaneous gaze reactions.

3.3 Evaluating data consistency

How consistently do third-party annotators label engagement intervals? We analyze their consistency to verify the predictability and soundness of our definition.

Table 2 shows the analysis. We quantify label agreement in terms of the average F_1 score, whether at the frame or interval level (see Supp.). We consider two aspects of agreement: boundary (how well do annotators agree on the start and end points of a positive interval?) and presence (how well do they agree on the existence of a positive interval?).

First we compare how consistent each of the 10 annotators’ labels are with the consensus ground truth (see “Turker vs. Consensus”). They have reasonable agreement on the rough interval locations, which verifies the soundness of our definition. Still, the F_1 score is not perfect, which indicates that the task is non-trivial even for humans. Some discrepancies are due to the fact that even when two annotators agree on the presence of an interval, their annotations will not match exactly in terms of the start and end frame. For example, one annotator might mark the start when the recorder searches for items on the shelf, while another might consider it to be when the recorder grabs the item. Indeed, agreement on the presence criterion (right column) is even higher, 0.914. The “Random vs. Consensus” entry compares a prior-informed random guess to the ground truth.³ These two extremes give useful bounds of what we can expect

³ We randomly generate interval predictions 10 times based on the prior of interval length and temporal distribution and report the average.

		Frame F_1	Interval F_1	
			Boundary Presence	
Turker	vs. Consensus	0.818	0.837	0.914
	vs. Recorder	0.589	0.626	0.813
Random	vs. Consensus	0.426	0.339	0.481
	vs. Recorder	0.399	0.344	0.478

Table 2: Analysis of inter-annotator consistency.

from our computational model: a predictor should perform better than random, but will not exceed the inter-human agreement.

Next, we check how well the third-party labels match the experience of the first-person recorder (see “Turker vs. Recorder”). We collect 3 hours of self-annotation from 4 of the recorders, and compare them to the Turker annotations. Similar to above, we see the Turkers are considerably more consistent with the recorder labels compared to the prior-informed random guess, though not perfect. As one might expect, Turker annotations have higher recall, but lower (yet reasonable) precision against the first-person labels. Overall, the 0.813 F_1 score for Turker-Recorder presence agreement indicates our labels are fairly faithful to individuals’ subjective interpretation.

4 Approach

We propose to learn the motion patterns in first-person video that indicate engagement. Two key factors motivate our decision to focus on motion. First, camera motion often contains useful information about the recorder’s intention [10, 12, 13]. This is especially true in egocentric video, where the recorder’s head and body motion heavily influence the observed motion. Second, motion patterns stand to generalize better across different scenarios, as they are mostly independent of the appearance of the surrounding objects and scene.

Our approach has three main stages. First we compute *frame-wise* predictions (Sec. 4.1). Then we leverage those frame predictions to generate *interval* hypotheses (Sec. 4.2). Finally, we describe each interval as a whole and classify it with an interval-trained model (Sec. 4.3). By departing from traditional frame-based decisions [17, 26, 27], we capture long-term temporal dependencies. As we will see below, doing so is beneficial for detecting subtle periods of engagement and accounting for their variable length. Fig. 2 shows the workflow.

4.1 Initial frame-wise estimates

To first compute frame-wise predictions, we construct one motion descriptor per frame. We divide the frame into a grid of 16×12 uniform cells and compute the optical flow vector in each cell. Then we temporally smooth the grid motion with a Gaussian kernel. Since at this stage we want to capture attention within a granularity of a second, we set the width of the kernel to two seconds. As shown in [15], smoothing the flow is valuable to integrate out the regular unstable head bobbles by the recorder; it helps the descriptor focus on prominent scene and camera motion. The frame descriptor consists of the smoothed flow vectors

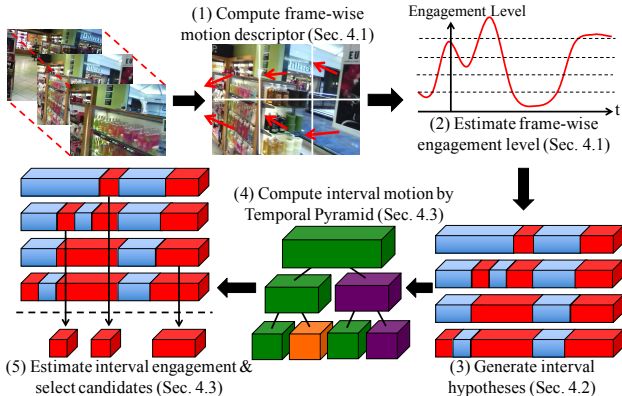


Fig. 2: Workflow for our approach.

concatenated across cells, together with the mean and standard deviation of all cells in the frame. It captures dominant egomotion and dynamic scene motion—both of which are relevant to first-person engagement.

We use these descriptors, together with the frame-level ground truth (cf. Sec. 3.2), to train an i.i.d. classifier. We use random forest classifiers due to their test-time efficiency and relative insensitivity to hyper-parameters, though of course other classifiers are possible. Given a test video, the confidence (posterior) output by the random forest is used as the initial frame-wise engagement estimate.

4.2 Generating interval proposals

After obtaining the preliminary estimate for each frame, we generate multiple hypotheses for engagement *intervals* using a level set method as follows. For a given threshold on the frame-based confidence, we obtain a set of positive intervals, where each positive interval consists of contiguous frames whose confidence exceeds the threshold. By sweeping through all possible thresholds (we use the decile), we generate multiple such sets of candidates. Candidates from all thresholds are pooled together to form a final set of *interval proposals*.

We apply this candidate generation process on both training data and test data. During training, it yields both positive and negative example intervals that we use to train an interval-level classifier (described next). During testing, it yields the hypotheses to which the classifier should be applied. This detection paradigm not only lets us avoid sliding temporal window search, but it also allows us to detect engagement intervals of variable length.

4.3 Describing and classifying intervals

For each interval proposal, we generate a motion descriptor that captures both the motion distribution and evolution over time. Motion evolution is important because a recorder usually performs multiple actions within an interval of engagement. For example, the recorder may stop, turn his head to stare at an object, reach out to touch it, then turn back to resume walking. Each action leads to a different motion pattern. Thus, unlike the temporally local frame-based descriptor above, here we aim to capture the statistics of the entire interval. We’d also

like the representation to be robust to time-scale variations (i.e., yielding similar descriptors for long and short instances of the same activity).

To this end, we use a temporal pyramid representation. For each level of the pyramid, we divide the interval from the previous level into two equal-length sub-intervals. For each sub-interval, we aggregate the frame motion computed in Sec. 4.1 by taking the dimension-wise mean and variance. So, the top level aggregates the motion of the entire interval, and its descendants aggregate increasingly finer time-scale intervals. The aggregated motion descriptors from all sub-intervals are concatenated to form a temporal pyramid descriptor. We use 3-level pyramids. To provide further context, we augment this descriptor with those of its temporal neighbor intervals (i.e., before and after). This captures the motion *change* from low engagement to high engagement and back.

We train a random forest classifier using this descriptor and the interval proposals from the training data, this time referring to the interval-level ground truth from Sec. 3.2. At test time, we apply this classifier to a test video’s interval proposals to score each one. If a frame is covered by multiple interval proposals, we take the highest confidence score as the final prediction per frame.

4.4 Discussion

Our method design is distinct from previous work in video *attention*, which typically operates per frame and uses temporally local measurements of motion [1, 2, 16, 17, 26, 27]. In contrast, we estimate engagement from interval hypotheses bootstrapped from initial frame estimates, and our representation captures motion changes over time at multiple scales. People often perform multiple actions during an engagement interval, which is well-captured by considering an interval together. For example, it is hard to tell whether the recorder is attracted by an object when we only know he glances at it, but it becomes clear if we know his following action is to turn to the object or to turn away quickly.

Simply flagging periods of low motion [15, 25, 27] is insufficient to detect all cases of heightened attention, since behaviors during the interval of engagement are often non-static and also exhibit learnable patterns. For example, shoppers move and handle objects they might buy; people sway while inspecting a painting; they look up and sweep their gaze downward when inspecting a skyscraper.

External sensors beyond the video stream could potentially provide cues useful to our task, such as inertial sensors to detect recorder motion and head orientation. However, such sensors are not always available, and they are quite noisy in practice. In fact, recent attempts to detect gazing behavior with inertial sensors alone yield false positive rates of 33% [25]. Similarly, although gaze could be informative for engagement, it requires greater instrumentation (i.e., eye tracker calibrated for each user) and will limit the applicability to generic egocentric video such as existing data on the web. This argues for the need for visual features for the challenging engagement detection task.

5 Experiments

We validate on two datasets and compare to many existing methods.



Fig. 3: Example engagement intervals detected by our method. Note the intra-interval variation: the recorder either performs multiple actions (Market), looks at an item from multiple views (Mall) or looks at multiple items (Museum). See videos on our website.

Baselines We compare with 9 existing methods, organized into four types:

Saliency Map: Following [17,21], we compute the saliency map for each frame and take the average saliency value. We apply the state-of-the-art learned video saliency model [1] and five others that were previously used for video summarization: [6, 7, 17, 19, 20]. We use the original authors’ code for [1, 6, 7, 19, 20] and implement [17]. Except [6], all these models use motion.

Motion Magnitude: Following [25,27], this baseline uses the inverse motion magnitude. Intuitively, the recorder becomes more still during his moments of high engagement as he inspects the object(s). We apply the same flow smoothing as in Sec. 4.1 and take the average.

Learned Appearance (CNN): This baseline predicts engagement based on the video content. We use state-of-the-art convolutional neural net (CNN) image descriptors, and train a random forest with the same frame-based ground truth our method uses. We use Caffe [43] and the provided pre-trained model (BVLC Reference CaffeNet).

Egocentric Important Region: This is the method of [28]. It is a learned metric designed for egocentric video that exploits hand detection, centrality in frame, etc. to predict the importance of regions for summarization. While the objective of “importance” is different than “engagement”, it is related and valuable as a comparison, particularly since it also targets egocentric data. We take the max importance per frame using the predictions shared by the authors.

Some of the baselines do not target our task specifically, a likely disadvantage. Nonetheless, their inclusion is useful to see if ego-engagement requires methods beyond existing saliency metrics. Besides, our baselines also include methods specialized for egocentric video [25,28], and one that targets exactly our task [25].

For the learned methods (ours, CNN and Important Regions), we use the classifier confidences to rate frames by their engagement level. Note that the CNN method has the benefit of training on the exact same data as our method. For the non-learned methods (saliency, motion), we use their magnitude. We evaluate two versions of our method: one with the interval proposals (Ours-interval) and one without (Ours-frame). The boundary agreement is used for interval prediction evaluation to favor methods with better localization of attention.

Datasets We evaluate on two datasets: our new UT Egocentric Engagement (UT EE) dataset and the public UT Egocentric dataset (UT Ego). We select all clips from UT Ego that contain browsing scenarios (mall, market), yielding 3 clips with total length of 58 minutes, and get them annotated with the same procedure in Sec. 3.2.

		Frame F_1	Interval F_1
GBVS	(Harel 2006 [19])	0.462	0.286
Self Resemblance	(Seo 2009 [20])	0.471	0.398
Bayesian Surprise	(Itti 2009 [7])	0.420	0.373
Salient Object	(Rahtu 2010 [6])	0.504	0.389
Video Attention	(Ejaz 2013 [17])	0.413	0.298
Video Saliency	(Rudoy 2013 [1])	0.435	0.396
Motion Mag.	(Rallapalli 2014 [25])	0.553	0.403
Cross Recorder	CNN Appearance	0.685	0.486
	Ours – frame	0.686	0.533
	Ours – interval	0.674	0.572
	Ours – GT interval	0.822	0.868
Cross Scenario	CNN Appearance	0.656	0.463
	Ours – frame	0.683	0.531
	Ours – interval	0.665	0.553
	Ours – GT interval	0.830	0.860
Cross Recorder AND Scenario	CNN Appearance	0.655	0.463
	Ours – frame	0.680	0.532
	Ours – interval	0.661	0.544
	Ours – GT interval	0.823	0.856

Table 3: F_1 -score accuracy of all methods on UT EE. (The cross-recorder/scenario distinctions are not relevant to the top block of methods, all of which do no learning.)

Implementation details We use the code of [44] for optical flow computation. Flow dominates our run-time, about 1.2 s per frame on 48 cores. The default settings are used for this and all the public saliency map codes. Using the scikit-learn package [45] for random forest, we train 2,400 trees in all results and leave all other parameters at default. The sample rate of video frames is 15 fps for optical flow and 1 fps for all other computation, including evaluation.

5.1 UT Egocentric Engagement (UT EE) dataset

We consider three strategies to form train-test data splits. The first is leave-one-recorder-out, denoted **cross-recorder**, in which we train a predictor for each recorder using exclusively video from *other* recorders. This setting tests the ability to generalize to new recorders (e.g., can we learn from John’s video to predict engagement in Mary’s video?). The second is leave-one-scenario-out, denoted as **cross-scenario**, in which we train a predictor for each scenario using exclusively video from other scenarios. This setting examines to what extent visual cues of engagement are independent of the specific activity or location the recorder (e.g., can we learn from a museum trip to predict engagement during a shopping trip?). The third strategy is the most stringent, disallowing any overlap in either the recorder or the scenario (**cross recorder AND scenario**).

Fig. 4(A)~(C) show the precision-recall curves for all methods and settings on the 14 hour UT EE dataset, and we summarize them in Table 3 using the F_1 scores; here we set the confidence threshold for each video such that 43.8% of its frames are positive, which is the ratio of positives in the entire dataset. Our method significantly outperforms all the existing methods. We also see our interval proposal idea has a clear positive impact on interval detection results. However, when evaluated with the frame classification metric (first column in Table 3), our interval method does not improve over our frame method. This is due to some inaccurate (too coarse) proposals, which may be helped by sampling

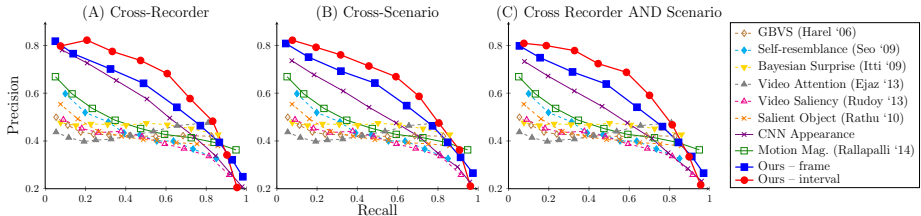


Fig. 4: Precision-recall accuracy on the UT EE dataset. Our approach detects engagement with much greater accuracy than an array of saliency and content-based methods, and our interval proposal idea improves the initial frame-wise predictions.

the level sets more densely. We also show an upper bound for the accuracy with perfect interval hypotheses (see Ours-GT interval), which emphasizes the need to go beyond frame-wise predictions as we propose.

Fig. 4 and Table 3 show our method performs similarly in all three train-test settings, meaning it generalizes to both new recorders and new scenarios. This is an interesting finding, since it is not obvious *a priori* that different people exhibit similar motion behavior when they become engaged, or that those behaviors translate between scenes and activities. This is important for applications, as it would be impractical to collect data for all recorders and scenarios.

The CNN baseline, which learns which video content corresponds to engagement, does the best of all the baselines. However, it is noticeably weaker than our motion-based approach. This result surprised us, as we did not expect the *appearance* of objects in the field of view during engagement intervals to be consistent enough to learn at all. However, there are some intra-scenario visual similarities in a subset of clips: four of the Museum videos are at the same museum (though the recorders focus on different parts), and five in the Mall contain long segments in clothing stores (albeit different ones). Overall we find the CNN baseline often fails to generate coherent predictions, and it predicts intervals much shorter than the ground truth. This suggests that appearance alone is a weaker signal than motion for the task.

Motion Magnitude (representative of [25,27]) is the next best baseline. While better than the saliency metrics, its short-term motion and lack of learning lead to substantially worse results than our approach. This also reveals that people often move while they engage with objects they want to learn more about.

Finally, despite their popularity in video summarization, Saliency Map methods [1,6,7,17,19,20] do not predict temporal ego-engagement well. In fact, they are weaker than the simpler motion magnitude baseline. This result accentuates the distinction between predicting gaze (the common saliency objective) and predicting first-person engagement. Clearly, spatial attention does not directly translate to the task. While all the Saliency Map methods (except [6]) incorporate motion cues, their reliance on temporally local motion, like flickers, makes them perform no better than the purely static image methods.

Fig. 3 shows example high engagement frames. **Please see the project webpage for video results.**

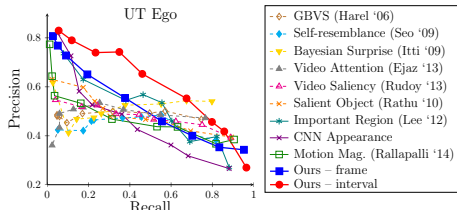


Fig. 5: Precision-recall accuracy on UT Ego dataset.

5.2 UT Egocentric dataset

Fig. 5 shows the results on the UT Ego dataset. The outcomes are consistent with those on UT EE above, and again our method performs the best. Whereas [28] is both trained and tested on UT Ego, our method does not do any training on the UT Ego data; rather, we use our model trained on UT EE. This ensures fairness to the baseline (and some disadvantage to our method).

Our method outperforms the Important Regions [28] method, which is specifically designed for first-person data. This result gives further evidence of our method’s cross-scenario generalizability. Important Regions [28] does outperform the Saliency Map methods on the whole, indicating that high-level semantic concepts are useful for detecting engagement, more so than low-level saliency. The CNN baseline does poorly, which reflects that its content-specific nature hinders generalization to a new data domain.

5.3 Start point correctness

Finally, Fig. 6 evaluates start point accuracy on UT EE. This setting is of interest to applications where it is essential to know the onset of engagement, but not necessarily its temporal extent. Here we run our method in a streaming fashion by using its frame-based predictions, without the benefit of hindsight on the entire intervals. Due to space limits, we defer to the Supp. for details.

6 Conclusion

We explore engagement detection in first-person video. By precisely defining the task and collecting a sizeable dataset, we offer the first systematic study of this problem. We introduced a learning-based approach that discovers the connection between first-person motion and engagement, together with an interval proposal approach to capture a recorder’s long-term motion. Results on two datasets show our method consistently outperforms a wide array of existing methods for visual attention. Our work provides the foundation for a new aspect of visual attention research. In future work, we will examine the role of external sensors (e.g., audio, gaze trackers, depth) that could assist in ego-engagement detection when they are available.

Acknowledgement

This research is supported in part by ONR YIP N00014-12-1-0754 and NSF IIS-1514118.

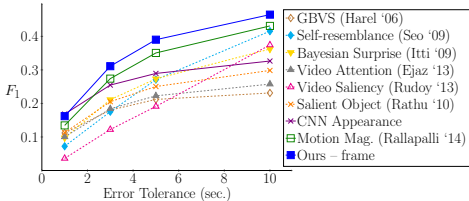


Fig. 6: Start-point accuracy on UT EE, measuring how well the onset of an engagement interval is detected in a streaming manner.

References

1. Rudoy, D., Goldman, D., Shechtman, E., Zelnik-Manor, L.: Learning video saliency from human gaze using candidate selection. In: CVPR. (2013)
2. Han, J., Sun, L., Hu, X., Han, J., Shao, L.: Spatial and temporal visual attention prediction in videos using eye movement data. *Neurocomputing* **145** (Dec 2014) 140–153
3. Lee, W., Huang, T., Yeh, S., Chen, H.: Learning-based prediction of visual attention for video signals. *IEEE TIP* **20**(11) (Nov 2011)
4. Abdollahian, G., Taskiran, C., Pizlo, Z., Delp, E.: Camera motion-based analysis of user generated video. *TMM* **12**(1) (Jan 2010) 28–41
5. Mahadevan, V., Vasconcelos, N.: Spatiotemporal saliency in dynamic scenes. *TPAMI* **32**(1) (Jan 2010) 171–177
6. Rahtu, E., Kannala, J., Salo, M., Heikkila, J.: Segmenting salient objects from images and videos. In: ECCV. (2010)
7. Itti, L., Baldi, P.: Bayesian surprise attracts human attention. *Vision Research* **49**(10) (2009) 1295–1306
8. Liu, H., Jiang, S., Huang, Q., Xu, C.: A generic virtual content insertion system based on visual attention analysis. In: ACM MM. (2008)
9. Li, Y., Fathi, A., Rehg, J.M.: Learning to predict gaze in egocentric video. In: ICCV. (2013)
10. Yamada, K., Sugano, Y., Okabe, T., Sato, Y., Sugimoto, A., Hiraki, K.: Attention prediction in egocentric video using motion and visual saliency. In: *Advances in Image and Video Technology*. (2012)
11. Yamada, K., Sugano, Y., Okabe, T., Sato, Y., Sugimoto, A., Hiraki, K.: Can saliency map models predict human egocentric visual attention? In: ACCV Workshop. (2011)
12. Kender, J., Yeo, B.L.: On the structure and analysis of home videos. In: ACCV. (2000)
13. Li, K., Oh, S., Perera, A., Fu, Y.: A videography analysis framework for video retrieval and summarization. In: BMVC. (2012)
14. Gygli, M., Grabner, H., Riemenschneider, H., Gool, L.V.: Creating summaries from user videos. In: ECCV. (2014)
15. Poleg, Y., Arora, C., Peleg, S.: Temporal segmentation of egocentric videos. In: CVPR. (2014)
16. Nguyen, T.V., Xu, M., Gao, G., Kankanhalli, M., Tian, Q., Yan, S.: Static saliency vs. dynamic saliency: a comparative study. In: ACM MM. (2013)
17. Ejaz, N., Mehmood, I., Baik, S.: Efficient visual attention based framework for extracting key frames from videos. *Image Communication* **28** (2013) 34–44
18. Itti, L., Dhavale, N., Pighin, F.: Realistic avatar eye and head animation using a neurobiological model of visual attention. In: Proc. SPIE 48th Annual International Symposium on Optical Science and Technology. Volume 5200. (Aug 2003) 64–78
19. Harel, J., Koch, C., Perona, P.: Graph-based visual saliency. In: NIPS. (2007)
20. Seo, H., Milanfar, P.: Static and space-time visual saliency detection by self-resemblance. *J. of Vision* **9**(7) (2009) 1–27
21. Ma, Y.F., Lu, L., Zhang, H.J., Li, M.: A user attention model for video summarization. In: ACM MM. (2002)
22. Kienzle, W., Schölkopf, B., Wichmann, F., Franz, M.: How to find interesting locations in video: a spatiotemporal interest point detector learned from human eye movements. In: DAGM. (2007)
23. Dorr, M., Martinetz, T., Gegenfurtner, K.R., Barth, E.: Variability of eye movements when viewing dynamic natural scenes. *J. of Vision* **10**(10) (2010) 1–17

24. Pilu, M.: On the use of attention clues for an autonomous wearable camera. Technical Report HPL-2002-195, HP Laboratories Bristol (2003)
25. Rallapalli, S., Ganesan, A., and V. Padmanabhan, K.C., Qiu, L.: Enabling physical analytics in retail stores using smart glasses. In: *MobiCom*. (2014)
26. Nakamura, Y., Ohde, J., Ohta, Y.: Structuring personal activity records based on attention-analyzing videos from head mounted camera. In: *ICPR*. (2000)
27. Cheatle, P.: Media content and type selection from always-on wearable video. In: *ICPR*. (2004)
28. Lee, Y.J., Ghosh, J., Grauman, K.: Discovering important people and objects for egocentric video summarization. In: *CVPR*. (2012)
29. Lu, Z., Grauman, K.: Story-driven summarization for egocentric video. In: *CVPR*. (2013)
30. Aghazadeh, O., Sullivan, J., Carlsson, S.: Novelty detection from an egocentric perspective. In: *CVPR*. (2011)
31. Hoshen, Y., Ben-Artzi, G., Peleg, S.: Wisdom of the crowd in egocentric video curation. In: *CVPR Workshop*. (2014)
32. Park, H.S., Jain, E., Sheikh, Y.: 3d gaze concurrences from head-mounted cameras. In: *NIPS*. (2012)
33. Fathi, A., Hodgins, J., Rehg, J.: Social interactions: A first-person perspective. In: *CVPR*. (2012)
34. Fathi, A., Farhadi, A., Rehg, J.: Understanding egocentric activities. In: *ICCV*. (2011)
35. Poleg, Y., Arora, C., Peleg, S.: Temporal segmentation of egocentric videos. In: *CVPR*. (2014)
36. Pirsiavash, H., Ramanan, D.: Detecting activities of daily living in first-person camera views. In: *CVPR*. (2012)
37. Damen, D., Leelasawassuk, T., Haines, O., Calway, A., Mayol-Cuevas, W.: You-do, i-learn: Discovering task relevant objects and their modes of interaction from multi-user egocentric video. In: 2014. (2014)
38. Soran, B., Farhadi, A., Shapiro, L.: Action recognition in the presence of one egocentric and multiple static cameras. In: *ACCV*. (2014)
39. Kitani, K., Okabe, T., Sato, Y., Sugimoto, A.: Fast unsupervised ego-action learning for first-person sports video. In: *CVPR*. (2011)
40. Spriggs, E., la Torre, F.D., Hebert, M.: Temporal segmentation and activity classification from first-person sensing. In: *CVPR Workshop on Egocentric Vision*. (2009)
41. Li, Y., Ye, Z., Rehg, J.: Delving into egocentric actions. In: *CVPR*. (2015)
42. Mital, P.K., Smith, T.J., Hill, R.L., Henderson, J.M.: Clustering of gaze during dynamic scene viewing is predicted by motion. *Cognitive Computation* **3**(1) (Mar 2011) 5–24
43. Jia, Y., Shelhamer, E., Donahue, J., Karayev, S., Long, J., Girshick, R., Sutskever, I., Darrell, T.: Caffe: Convolutional architecture for fast feature embedding. arXiv preprint arXiv:1408.5093 (2014)
44. Liu, C.: Beyond Pixels: Exploring New Representations and Applications for Motion Analysis. PhD thesis, Massachusetts Institute of Technology (May 2009)
45. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., Duchesnay, E.: Scikit-learn: Machine learning in Python. *JMLR* **12** (Nov 2011) 2825–2830