

Switch-a-View: View Selection Learned from Unlabeled In-the-wild Videos

Sagnik Majumder¹ Tushar Nagarajan¹ Ziad Al-Halah² Kristen Grauman¹
¹UT Austin ²University of Utah

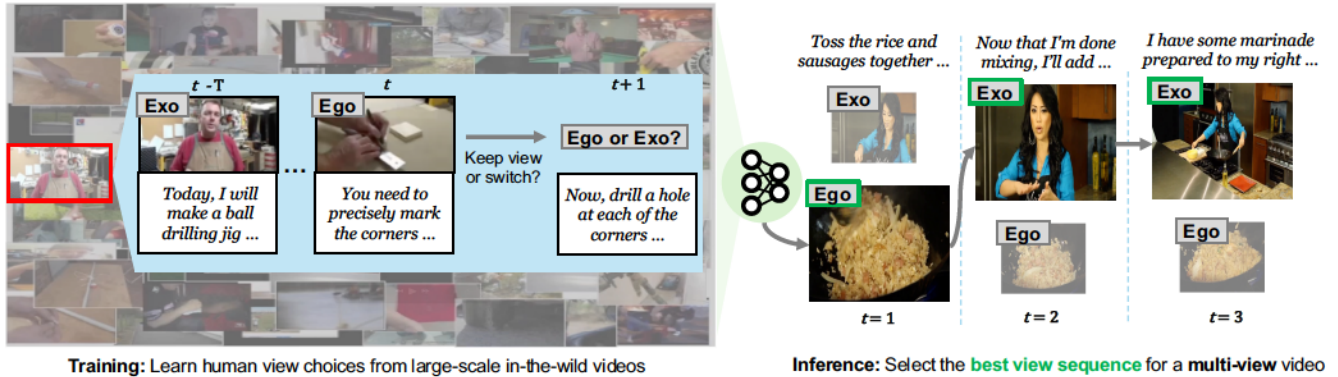


Figure 1. Given a multi-view narrated how-to video, can we select the sequence of camera viewpoints that best show the activity—automating the camerawork that is today done with manual editing? While direct supervision for this task is impractical, our SWITCH-A-VIEW approach shows how to learn *typical viewpoint choice patterns* from large-scale unlabeled in-the-wild instructional videos (left), then translate those patterns to novel multi-view videos (right), yielding an informative how-to that hops between the most useful ego/exo viewpoints.

Abstract

We introduce SWITCH-A-VIEW, a model that learns to automatically select the viewpoint to display at each timepoint when creating a how-to video. The key insight of our approach is how to train such a model from unlabeled—but human-edited—video samples. We pose a pretext task that pseudo-labels segments in the training videos for their primary viewpoint (egocentric or exocentric), and then discovers the patterns between the visual and spoken content in a how-to video on the one hand and its view-switch moments on the other hand. Armed with this predictor, our model can be applied to new multi-view video settings for orchestrating which viewpoint should be displayed when, even when such settings come with limited labels. We demonstrate our idea on a variety of real-world videos from *HowTo100M* and *Ego-Exo4D*, and rigorously validate its advantages. Project: https://vision.cs.utexas.edu/projects/switch_a_view/.

1. Introduction

Video is an amazing medium for communication, and today’s widely used Internet platforms make it easy to create and share content broadly. Instructional or “how-to” video is

particularly compelling in this setting: YouTube, TikTok, and similar sites have democratized the ability to share our talents with others, by both showing and telling how to perform some special skill. From how to plant a garden, how to make yogurt, how to fold origami, or how to give a dog a haircut, there is no shortage of how-to nuggets produced and consumed by users of many ages and backgrounds.

Creating an effective how-to video, however, is not trivial. From potentially hours of footage from multiple cameras capturing all aspects of the instructional activity, a creator needs to edit down to the essential steps of their demonstration *and* decide on the camera viewpoint (view) for each temporal segment that best reveals what they want to show. For example, when showing how to cut the dog’s hair, the instructor might first appear standing beside the dog—the camera more distant—then the camera may zoom close up to her using scissors and describing how to trim near the ear, then zoom back out while she shows progress across the dog’s body. How-to videos often exhibit this sequential mix of “exocentric” and “egocentric-like” viewpoints to effectively recap the procedure with clear visuals.

The status quo is to either orchestrate camerawork live while filming, or do post-recording editing among the multiple available cameras—both of which are labor intensive. Work in automatic cinematography [7, 16, 17, 23, 53, 61],

though inspiring, relies on heuristics or domain-specific models that are not equipped to address automatic editing of video demonstrations. How could we train an “AI how-to cameraman”, which, given a stream of 2 or more simultaneous camera views, could hop between them intelligently?

Supervising this learning task presents a problem. There are vast amounts of positive examples of well-edited how-to videos, but those edited results hide the “negatives”—the viewpoints that were *not* chosen for inclusion in the final video at any given time point. Those are left on the cutting room floor. This makes it unclear how to translate the editing patterns in in-the-wild edited video to new data.

To tackle this learning challenge, we design a pretext task for learning human view preferences from varying-view instructional videos on the Web. “Varying-view” means that the source training videos display an arbitrary number of view switches over the course of the video (e.g., from ego to exo and back as in our example above), and contain only *one* viewpoint at any time. We introduce a model called SWITCH-A-VIEW that learns from such data; it uses past frames together with the how-to narrations spoken by the demonstrator to learn a binary classifier indicating whether the viewpoint is going to switch or not at the current time step. Then, we deploy this pretext-trained model in *multi-view*, narrated video settings with limited best view labels, and decide how to orchestrate the view selection of such videos over time. In this way, our approach captures the view-switch patterns from widely diverse unlabeled in-the-wild videos, then translates those trends to automatically direct the camerawork in new instances. See Fig. 1.

We train and evaluate our approach on HowTo100M [37], an extensive repository of real-world how-to videos, and further show generalization to multi-view Ego-Exo4D [19] videos. Our data confirms that human judges exhibit substantial agreement on what constitutes a “best view” in a how-to video, establishing that it is possible to rigorously evaluate this task. Furthermore, our results show SWITCH-A-VIEW outperforms the state-of-the-art in multi-view video view selection [35] as well as multimodal retrieval [59] and other baselines.

2. Related work

Automatic cinematography. In automatic cinematography, systems automate the process of creating an effective video presentation given a video scene. Standard techniques include planning and controlling camera movements, angles and transitions. Prior works target classroom environments [17, 22, 61], group activities [2], or (pseudo-)panoramic recordings [6, 7, 9, 16, 53, 54, 60]. Different from all of the above, we tackle view selection in multi-view *instructional* scenarios. Moreover, we seek a lighter-weight supervision solution: whereas prior work uses supervised discriminative methods requiring large-scale best

view labels [7, 23, 53] or bootstraps view selector training using multi-view videos annotated with view-agnostic narrations [34], we aim to learn view selection from readily available in-the-wild *unlabeled* instructional videos. Furthermore, our model is multimodal, integrating both the video content as well as its transcribed speech.

View selection in active perception. More distant from our problem, work in active perception and robotics considers how agents can intelligently select their visual input stream. This includes next view selection, where an embodied agent learns to actively place a camera for tasks like visual recognition [1, 8, 12, 24, 43] and segmentation [47, 48]. Whereas the objective in such work is to spend less agent time or compute to see sufficient content, our goal is instead to choose an informative camera view for human consumption. In our setting, the cameras are placed at certain stationary locations, or worn by a person performing an activity.

Weak supervision from Web data. Large-scale instructional data from the Web has been shown to provide weak supervision for understanding instructional activities, by aligning frames [36] and narrations [30, 36] with their step descriptions from instructional Web articles (e.g., Wiki-How), or through modeling the temporal order and interdependence of steps [3, 63, 64]. Unlike any of these methods, we tackle a distinct problem of weakly supervised view-switch detection in instructional videos, with the end goal of using the detector for view selection.

Video summarization. Temporal video summarization [4, 21, 38, 40, 45] entails creating a short but informative summary of a long video by subsampling keyframes or clips from it. While early methods are largely unsupervised [25, 32, 41, 51], more recent works derive supervision from manual labels [18, 20, 28, 44, 52, 62]. Limited work explores summarization in the context of multiple input videos [10, 13, 40, 46]. Video summarization and viewpoint selection are two entirely distinct tasks. Video summarization aims to downsample the video in time to the essential parts, whereas our task essentially requires downsampling the video in *space* to isolate the most informative viewpoint.

3. Approach

Our goal is to train a model to predict the “best view sequence” for multi-camera instructional videos — the sequence of camera viewpoints (views) that a human would most likely select to demonstrate an instructional activity (e.g., a close-up view of ingredients in a cooking video, moving to a wide-shot view when the chef speaks and gestures). To tackle this, we train a model for the proxy task of detecting “view switches” in varying-view instructional videos, which we then bootstrap to form a view selection model.

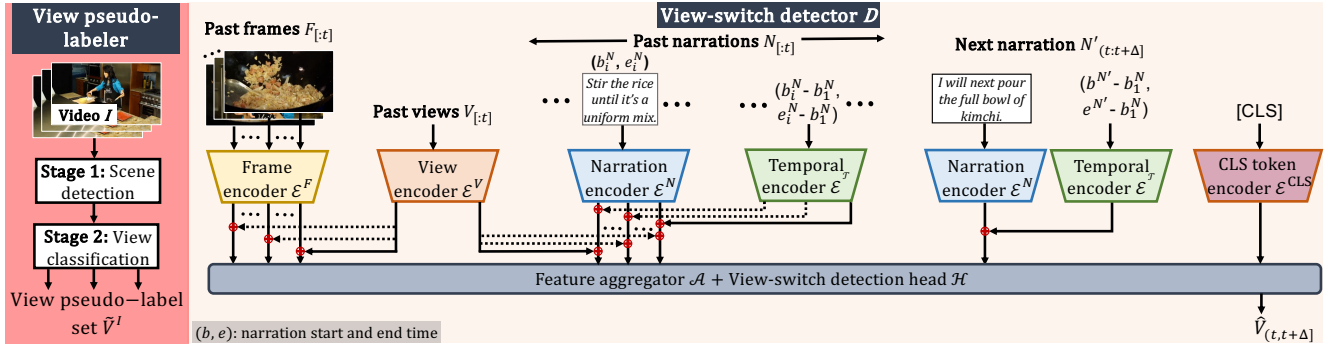


Figure 2. Given varying-view instructional videos—videos composed of a sequence of views chosen by human(s) to accurately show the instructional activity at all times—our goal is to train a view-switch detector D that can predict if the view should switch or not, at any time in a new video. Our hypothesis is that such a detector, when trained on large-scale and in-the-wild videos, can capture human view preferences and facilitate learning best view selection in multi-view settings with limited labels. However, such in-the-wild videos lack view labels. To train nevertheless, we propose an approach comprising (a) a view pseudo-labeler (left) that given a varying-view instructional video I , automatically classifies views in it and generates a pseudo-label set \tilde{V}^I , and (b) a view-switch detector D (right) that given the pseudo-labels \tilde{V}^I and any time t in I , learns to predict the next view. The prediction is conditioned on the past frames, past narrations, and the next narration, where narrations are naturally occurring spoken content from the how-to demonstrator.

First, we formally define our pretext task (Sec. 3.1). Next, we describe how to source pseudo-labels for our pretext task by automatically classifying views in varying-view videos (Sec. 3.2). We then describe our method and how to train it to predict view-switches (Sec. 3.3). Finally, we describe how our view-switch detector can bootstrap learning a view selection model (Sec. 3.4 and 3.5) with limited labels.

3.1. View-switch detection as a pretext task

We introduce our pretext task: view-switch detection in varying-view instructional videos. Consider a varying-view instructional video I , where the view changes back and forth over time between a close-up / *egocentric-like* (ego) view, and a wide shot / *exocentric-like* (exo) view.¹ This results in a sequence of varying views V . The instructional video also contains a sequence of narrations N , where each narration N_i has a start and end time, (b_i, e_i) , and provides commentary transcribed to text. These narrations are free-form spoken language from the demonstrator, which capture their actions (“hammer the nail in there”) as well as side comments (“sometimes I use my sander instead”, “thanks for watching!”).

We formulate the view-switch detection task as a two-class *view prediction* problem, where at any time t in the video, the model must detect if the view should be of type *ego* or *exo*, to best showcase the activity over the *next* Δ seconds. More specifically, we require a model D that predicts the human-preferred view $V_{(t,t+\Delta]}$ given the past video, narrations and views, as well as cues from the next narration. Formally,

$$D(F_{[t]}, N_{[t]}, V_{[t]}, N'_{(t,t+\Delta]}) = V_{(t,t+\Delta]},$$

¹We adopt this ego/exo view taxonomy given their importance and prevalence in instructional video datasets [19, 37, 55]

where $F_{[t]}$ is the past frames, $N_{[t]}$ is the past narrations and $V_{[t]}$ is the past views. $N'_{(t,t+\Delta]}$ is the next narration, if it overlaps with the prediction interval, and an empty string otherwise. This formulation provides a path from the *next-view prediction task* to the *view-switch task*: since the most recent past view is given, estimating the desired next view—and comparing it with the latest past view—is equivalent to predicting whether the view switched.

While past narrations provide high-level cues about past activity steps and what views were chosen to demonstrate those steps, past frames offer more fine-grained information about the same. They together form the past context that can help anticipate the next view. The next narration is essential to disambiguate between various potential actions that the demonstrator may do next, and the language directly hints at the appropriate views (e.g., the person says “next, let’s take a closer look at ...” suggesting an ego view). Thus, combining these inputs will offer valuable cues to our detector.²

Critically, we aim to train this detector on large-scale and in-the-wild instructional videos [37, 55]. We show that training for this pretext task can enable view selection models for multi-camera settings, with limited supervision. In short, representations developed to detect when to “switch view” can be repurposed with minimal modification to select the “best view” to switch to, since they contain rich knowledge of human-selected view-switch patterns in a large variety of in-the-wild scenarios. Next, we show how to source pseudo-labels to train such models.

²Note that next-step narrations are also available at inference time, when we have multi-view content and a full narration track, and we aim to perform view selection.

3.2. Sourcing “view-switch” pseudo-labels

Instructional videos [37, 55] are an ideal source of varying-view data, however they do not come paired with information about what camera viewpoint is chosen for each segment. We therefore design a strategy to automatically identify and pseudo-label their underlying view sequences. We do this in two stages (Fig. 2 left).

First, given a video I , we use an off-the-shelf scene detector (PySceneDetect [5]) to compute scene boundaries. Using this, we split the video into a sequence of contiguous shots. Next, we classify each frame in the video using a pre-trained ego vs. exo view classifier, and then aggregate the class predictions into a shot-level pseudo-label. Specifically, given a shot from I , we first split it into a sequence of fixed-length clips. Next, we feed each clip to the view classifier that produces the probability that the clip is from an ego vs. exo view. We then compute the pseudo-label for the whole shot by averaging the view probabilities across all its clips. We repeat these steps for all shots, and assign each frame in I the same pseudo-label as the shot it lies in, to finally obtain a pseudo-label set \tilde{V}^I . Combining the classifier with the scene detector reduces the overall noise in the pseudo-labels due to classification failures at scene boundaries. We use a learned model [29] for ego-exo view classification, trained on the Charades-Ego [50] dataset.

3.3. View-switch detector design

Given a video I and any time t in it, our view-switch detector D must successfully predict the view for the future time interval $(t, t + \Delta]$. It must do so using the frames, narrations and views from the past, and also the next narration, if it overlaps with the prediction interval (c.f. Sec. 3.1). See Fig. 2 right. In the following, we provide details on how our method extracts features from each input and then aggregates them for making a view prediction.

Frame encoding. We begin by using a frame encoder \mathcal{E}^F to embed the past frames $F_{[:t]}$ and produce a visual feature sequence f , where each frame F_i has a feature f_i . We further enhance each feature f_i by using a view encoder \mathcal{E}^V to embed the corresponding view V_i^F into a view feature v_i^F and add it to f_i . We also encode frame F_i ’s temporal position relative to the start time of the first past narration, denoted by \mathcal{T}_i^F , using a temporal encoder \mathcal{E}^T into a feature τ_i^F , and add it to the enhanced frame feature. Formally,

$$f_i = \mathcal{E}^F(F_i) + \mathcal{E}^V(V_i^F) + \mathcal{E}^T(\mathcal{T}_i^F). \quad (1)$$

Producing a feature per frame and augmenting it with view and temporal information helps us create a *fine-grained*, and *view-* and *temporally-aware* representation that is anticipative of the next view.

Narration encoding. Next, we encode each past narration from $N_{[:t]}$, and the next narration $N'_{(t,t+\Delta]}$ by using an LLM

encoder. This generates a text feature sequence n for the past narrations and a single text feature n' for the next narration.

Similar to our encoding of past frames, we also make the features for past narrations *view-aware*. To do so, we first produce a per-view count of the frames that lie in the interval of each past narration N_i . We then estimate a dominant view V_i^N for the narration—called narration view, henceforth—by setting it to the most frequent view according to the per-view frame count. Next, we use our view encoder \mathcal{E}^V to embed the narration view V_i^N into a view feature v_i^N . Finally, we update the narration feature n_i by adding it with v_i^N .

Moreover, for both past and next narrations, we provide their temporal information to our model so that it can infer the alignment between the frames and the narrations, and use it to improve its cross-modal reasoning. To this end, we first normalize the start and end time pair, (b, e) , for each past narration N_i and next narration N' , to be relative to the start time of the first past narration. We then compute the mean time of each pair: \mathcal{T}_i^N and $\mathcal{T}^{N'}$ for N_i and N' , respectively. These means convey the temporal locations of the narrations relative to each other. Next, we encode each relative mean with the temporal encoder \mathcal{E}^T and obtain a temporal feature: τ_i^N and $\tau^{N'}$ for N_i and N' , respectively. Finally, we update the narration features by adding them with their temporal features. Formally,

$$\begin{aligned} n_i &= \mathcal{E}^N(N_i) + \mathcal{E}^V(V_i^N) + \mathcal{E}^T(\mathcal{T}_i^N); \\ n' &= \mathcal{E}^N(N') + \mathcal{E}^T(\mathcal{T}^{N'}). \end{aligned} \quad (2)$$

Feature aggregation and view classification. To aggregate the visual and narration features, we first add modality features m^F and m^N to the frame features f , and narration features, n and n' , respectively. Formally,

$$f_i = f_i + m^F; \quad n_i = n_i + m^N; \quad n' = n' + m^N \quad (3)$$

These are learnable embeddings that enable our model to distinguish between the visual and text modalities, and successfully do cross-modal reasoning.

We also introduce a [CLS] token in our model, and embed it with an encoder \mathcal{E}^{CLS} to produce a feature c , so that the output of our feature aggregator, which corresponds to the [CLS] token, can be used to estimate the next view. Next, we feed the frame features f , the past narration features n , the next narration feature n' , and the [CLS]-token feature c into a feature aggregator \mathcal{A} . \mathcal{A} comprises a transformer [58] encoder that performs self-attention on all features and extracts multi-modal cues that are predictive of the next view. Finally, we take the output feature of \mathcal{A} , which corresponds to the [CLS] token, and pass it to a view classification head \mathcal{H} to get an estimate $\hat{V}_{(t,t+\Delta]}$ of the next view $V_{(t,t+\Delta]}$. Formally,

$$\hat{V}_{(t,t+\Delta]} = \mathcal{H}(\mathcal{A}(f, n, n', c)[j_{\text{CLS}}]), \quad (4)$$

where j_{CLS} is the feature index for the [CLS] token.

3.4. Repurposing switch detection for view selection

Recall that in view selection, given a multi-view instructional video I and any time t in it, the goal is to predict the view that is preferred by humans for showing the activity in an interval $[t, t+\Delta]$. We introduce a view selector S for tackling this task. S is a modification of our view-switch detector D , such that S additionally has access to the *frames from the simultaneously captured ego and exo views* during the prediction interval $[t, t+\Delta]$.

To this end, we first use our frame encoder \mathcal{E}^F to embed the ego frames $F_{[t,t+\Delta]}^G$ and exo frames $F_{[t,t+\Delta]}^X$ into visual features f^G and f^X , respectively. Next, we append f^G and f^X to the input sequence of our feature aggregator \mathcal{A} . Finally, we treat \mathcal{A} 's output feature for its [CLS] token input, as a representation of the best view for $[t, t+\Delta]$, and feed it to the detector's view classification head \mathcal{H} to get an estimate $\hat{V}_{[t,t+\Delta]}$ of the best view $V_{[t,t+\Delta]}$. Formally,

$$\hat{V}_{[t,t+\Delta]} = \mathcal{H}(\mathcal{A}(f, n, n, f^G, f^X, c)[j_{\text{CLS}}]), \quad (5)$$

where j_{CLS} is the feature index for the [CLS] token.

To learn view selection we initialize S with the parameters of our detector, trained on the view-switch detection task, and finetune it using a small set of samples labeled for view selection. Next, we provide details for training and finetuning.

3.5. Model training objective

We train our view-switch detector D with a view classification loss \mathcal{L}^D . We set \mathcal{L}^D to

$$\mathcal{L}^D = \mathcal{L}_{\text{CE}}(\hat{V}_{[t,t+\Delta]}, \tilde{V}_{[t,t+\Delta]}), \quad (6)$$

where $\hat{V}_{[t,t+\Delta]}$ is our estimated view (c.f. Sec. 3.3) and $\tilde{V}_{[t,t+\Delta]}$ is the pseudo-label from our view pseudo-labeler (c.f. Sec. 3.2).

To train our view selector S , we obtain a *small* train set of best view labels, B , such that $B = \{V_{[t_1,t_1+\Delta]}, \dots, V_{[t_W,t_W+\Delta]}\}$, and W is the label count in B . For each best view label $V_{[t_w,t_w+\Delta]} \in B$, and the corresponding view estimate $\hat{V}_{[t_w,t_w+\Delta]}$, per our view selector S (c.f. Sec. 3.4), we set our view selection loss \mathcal{L}^S to a cross-entropy loss, such that

$$\mathcal{L}^S = \mathcal{L}_{\text{CE}}(\hat{V}_{[t_w,t_w+\Delta]}, V_{[t_w,t_w+\Delta]}). \quad (7)$$

Once trained, our framework can accurately choose the preferred view in novel multi-view videos.

4. Datasets and annotations

Datasets. We use two datasets in our experiments. **HT100M** [37] is a large-scale dataset of narrated and in-the-wild instructional videos. These videos are view-varying

in nature, and the views can be broadly categorized as ego or exo. This, along with the diversity and realism of HT100M, makes it ideal for our view-switch detection task. **Ego-Exo4D** [19] contains multi-view videos, where each video is captured with five time-synced cameras—one is an ego camera worn by a human performing an instructional activity, and the other four are stationary exo cameras placed around the scene. Moreover, the narrate-and-act (N&A) subset of Ego-Exo4D has videos of humans narrating and performing an activity, where the narrations are free-form and match in style with HT100M, making it compatible with our task of view selection with limited labels.

Training data. To train the view-switch detector, we use 3,416 hours of HT100M videos spanning a diverse set of activities (cooking, DIY, household, etc.) and pseudo-label shots from these videos (c.f. Sec. 3.2). See Supp. for details.

Evaluation data. For evaluation, we use both HT100M and Ego-Exo4D [19], where the view-switch evaluation on Ego-Exo4D is *zero-shot*. While the training sets are automatically generated and pseudolabeled, we ensure a gold-standard test set free of noise by manually annotating videos for our tasks. To this end, we recruit trained annotators to manually annotate the view types for HT100M and the human-preferred views for Ego-Exo4D.

For HT100M, we identify 975 hours of videos that do not overlap with our train videos above. We segment 4,487 fixed-length clips, each with length set to the prediction interval Δ (c.f. Sec. 3.1). Next, we ask trained annotators to label these clips as either ego or exo. See Supp. for full annotation instructions and more details.

For Ego-Exo4D, we create a test set containing 2.7 hours of N&A videos spanning six activity categories (cooking, bike repair, rock climbing, dancing, soccer, basketball). For each video, we use its "best-exo-view" annotation from Ego-Exo4D to generate an ego-exo view pair comprising the single ego and the best exo view. As before, we create Δ length clips from each view. We then couple the pair with its closest atomic activity description (time-stamped manual descriptions of the camera wearer's activity [19]) and ask our annotators to label the *view* between the two that *best* demonstrates the activity described in the narration (see Supp. Fig. 3). Importantly, this means that annotators specifically select the "best" view as the one that most clearly illustrates the *current actions of the camera wearer*, consistent with our how-to video view selection goal.

Annotator agreement on best view. To ensure annotation quality for *both* datasets, in addition to providing detailed annotation guidelines and concrete examples (available in Supp.), we require annotators to take qualifiers with stringent passing criteria and we solicit 9 annotators' responses for each instance. We accept an annotation only if the inter-

annotator agreement is at least 78%, meaning at least 7 out of 9 annotators agree. This resulted in a Cohen’s kappa coefficient [11] of 0.65 for HT100M and 0.70 for Ego-Exo4D—both of which constitute “substantial” agreement [27]. This solid agreement assures the quality of our test set; despite there being some room for subjectivity in deciding the “best” view for a how-to, this data shows human judges are indeed able to substantially agree.

This results in a final total of 3,151 and 5,049 test instances for HT100M and Ego-Exo4D, respectively. In Supp. we filter with even higher agreement thresholds, yielding even more selective (but smaller) test sets; trends for our method vs. baselines remain consistent.

Data for view selection with limited labels. We train and evaluate our view selector on a small dataset comprising Ego-Exo4D [19] videos. For our train data, we follow our annotation protocol for evaluating view-switch detection on Ego-Exo4D, and collect view annotations for a total of 3.5 hours of train videos. This results in a total of 6,634 train instances. For evaluation, we use our test set from view-switch detection. This reuse is possible since a label indicates the type (ego/exo) of the desired next view for view-switch detection, or the desired current view for view selection. Train and test videos for this task are disjoint. See Supp. for details.

5. Experiments

Implementation. We set the durations of past frames to 8 seconds and past narrations to 32 seconds, and the prediction interval to $\Delta = 2$ seconds. We set the sample count for view selection to $W = 5000$. We evaluate view-switch detection on HowTo100M [37] by obtaining the views for the past frames (c.f. Sec. 3.3) from our pseudo-labeler. For Ego-Exo4D, we adopt a teacher-forcing setup and evaluate both tasks by using the ground-truth annotations for past frames and views. We implement our view-switch detector D and view selector S using the DINOv2 [39] encoder for our frame encoder \mathcal{E}^F , the Llama 2 [56] encoder for our narration encoder \mathcal{E}^N , a 8-layer transformer encoder [58] for our feature aggregator \mathcal{A} , a 2-layer MLP for the view classification head \mathcal{H} , and learnable embedding layers for our view encoder \mathcal{E}^V and temporal encoder \mathcal{E}^T . See Supp. for more details.

Baselines. We provide strong baselines comprising SOTA models and representations, as well as relevant heuristics. For *view-switch detection*, we compare against

- **InternVideo2 retrieval [59]:** a set of baselines that given the most recent past frame (**Retrieval [59]-F**), most recent past narration (**Retrieval [59]-N**), or next narration (**Retrieval [59]-N'**), first encodes [59] them into fine-grained features that capture multi-frame temporal contexts, then uses feature similarity to retrieve a nearest

Model	HowTo100M [37]			Ego-Exo4D [19]		
	Accuracy	AUC	AP	Accuracy	AUC	AP
All-ego/exo	50.0	50.0	50.0	50.0	50.0	50.0
Random	52.0	52.0	51.0	49.3	49.3	49.7
Last-frame	42.3	42.3	53.4	50.0	50.0	50.0
First-person pronoun detector	47.8	47.8	46.4	50.3	50.3	50.1
Retrieval [59]-F	<u>53.4</u>	<u>53.4</u>	<u>53.2</u>	52.6	<u>52.6</u>	<u>53.6</u>
Retrieval [59]-N	52.1	52.1	51.8	52.0	52.0	50.6
Retrieval [59]-N'	52.6	52.6	52.9	52.1	52.1	52.6
SWITCH-A-VIEW (Ours)	59.4	63.8	60.5	<u>51.2</u>	56.4	55.4

Table 1. View-switch detection results. Evaluation on Ego-Exo4D [19] is zero-shot. All values are in %, and higher is better.

Model	Accuracy	AUC	AP
All-ego/exo	50.0	50.0	50.0
Random	49.3	49.3	49.7
Last-frame	50.0	50.0	50.0
First-person pronoun detector	50.3	50.3	50.1
Retrieval [59]-F	52.3	52.3	53.6
Retrieval [59]-N	51.9	51.9	51.0
Retrieval [59]-N'	52.4	52.4	52.4
View-narration [59] Similarity	52.5	52.4	53.9
LangView [34]-smallData	52.1	52.6	53.2
LangView [34]-bigData (privileged)	<u>53.3</u>	<u>54.8</u>	<u>54.5</u>
Ours w/o pretraining	50.1	51.6	51.3
SWITCH-A-VIEW (Ours)	54.0	57.3	56.0

Table 2. Results and ablation for view selection with limited labels. All values are in %; higher is better. Significance $p \leq 0.05$.

neighbor of the same input type from the train set, and finally outputs the next view for F or N , or the corresponding view for N' , as its prediction.³

- **All-ego, All-exo, Random, Last-frame:** these are heuristics that use the ego view (**All-ego**), the exo view (**All-exo**), a randomly chosen (**Random**) view, or the view of the most recent past frame (**Last-frame**), as their prediction.
- **First-person pronoun detector:** a heuristic that predicts exo when it detects first-person pronouns like “I”, “We”, “My” or “Our” in the next narration, as human editors often use a wide shot that reveals their face or full body, when using such pronouns.

For *view selection with limited labels*, in addition to the baselines listed above, we compare against the following:

- **LangView [34]:** a SOTA view selector that uses narrations for weakly supervised pretraining. We finetune this model with our Ego-Exo4D labels (Sec. 4). We evaluate two versions of this baseline: LangView-*bigData* and LangView-*smallData*, which use large-scale Ego-Exo4D [34] videos, and our same small subset (Sec. 4), respectively, for pretraining. Note that the *bigData* variant enjoys access to **98x** more training samples than our method, an advantage

³See Supp. for parallel evaluation with CLIP [42]-style encoders, which generally underperformed InternVideo2 [59] encoders.

for the baseline.

- **View-narration [59] Similarity (VN-Sim):** separately computes the cosine similarity between the InternVideo2 features [59] for each view and the next narration, and picks the view most similar to the narration.

LangView evaluates how our model fares against SOTA view selection, while the retrieval and view-narration baselines analyze whether SOTA video-language embeddings are sufficient for this task. The heuristics verify the challenging nature of the tasks.

Evaluation metrics. We consider three metrics: 1) **Accuracy**, which directly measures the agreement between our predictions and labels; 2) **AUC**, the area under the ROC curve; and 3) **AP**, the average precision (AP) of the precision vs. recall curve. We use AUC and AP to account for the possible class imbalance in our collected annotations. Moreover, for each metric, we separately compute its value for the *same-view* and *view-switch* instances in our test sets, and report the mean. This lets us account for differences in the same-view and view-switch frequency, and obtain unbiased performance measures.

View-switch detection. In Table 1, we report our view-switch detection results. The heuristics generally perform the worst on both datasets, underlining the challenging nature of the task. The Retrieval [59] baselines improve over them, indicating that our model inputs do provide cues about the the view type. Among the Retrieval baselines, retrieving using the most recent past frame performs the best, showing that the past frames offer fine-grained task-relevant information beyond the narration words. Moreover, retrieving with the next narration is better than retrieving with the most recent past narration, revealing that the next narration carries more pertinent details about the desired view. This is likely because the next narration is better aligned with the time interval for which the view is being predicted.

Our method outperforms all baselines on both datasets, with the AUC margin over the best baseline, Retrieval [59]-*F*, being as high as 10.4% on HowTo100M (HT100M) [37] and 3.8% on Ego-Exo4D [19]. Our improvement over the Retrieval baselines show that computing feature [59]-level similarities are not enough for this task. Instead, learning it by leveraging complementary cues from both narrations and frames is critical. Moreover, our zero-shot results on Ego-Exo4D speak to our model’s efficacy vis-a-vis learning human view patterns from large-scale and in-the-wild videos, which generalize to different scenarios, without any training.

View selection. Table 2 shows our results on view selection with limited labels. For the heuristics and Retrieval [59] baselines, we observe the same performance trends as view-switch detection. The View-narration [59] Similarity (VN-Sim) baseline marginally improves over these methods, indi-

Model	HowTo100M [37]			Ego-Exo4D [19]		
	Accuracy	AUC	AP	Accuracy	AUC	AP
<i>N</i> -only	53.5	54.4	52.3	50.0	48.7	49.0
<i>N'</i> -only	55.4	57.8	56.2	49.8	49.8	50.0
<i>F</i> -only	53.3	54.5	54.7	51.0	53.4	53.2
(<i>F</i> , <i>N'</i>)-only	55.5	60.1	<u>58.1</u>	52.1	54.2	52.6
(<i>N</i> , <i>N'</i>)-only	<u>57.5</u>	59.3	56.6	50.0	53.0	52.6
(<i>F</i> , <i>N</i>)-only	56.0	<u>60.9</u>	57.4	51.8	<u>54.9</u>	<u>54.2</u>
Ours	59.4	63.8	60.5	<u>51.2</u>	56.4	55.4

Table 3. Ablation study for view-switch detection. All values are in %, and higher is better. Significance $p \leq 0.05$.

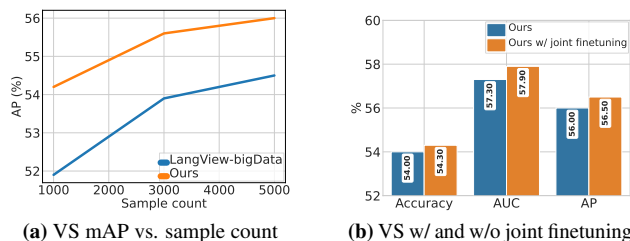


Figure 3. (a) Effect of sample count on our view selection (VS) performance; (b) Impact of joint finetuning with narration-based pseudo-labels [34] and best view labels on view selection (VS)

cating the frames from candidate views when combined with the corresponding narration (*N'*) provide direct cues about the preferred view. LangView [34]’s results benefit from its language-guided training, generally outperforming VN-Sim.

Our method significantly improves over all baselines, with the AUC margin over the best baseline, LangView [34]-bigData, being . 2.5%. Our gains over VN-Sim underscore that using feature similarity to match the activity described in the next narration with a candidate view does not suffice, and instead a model like ours, which can leverage multi-modal cues from the combination of both past and candidate frames, and past and next narrations, is valuable for this task. Training our model from scratch with only the small set of best view labels (“ours w/o pretraining”) is significantly weaker, showing that our view-switch pretraining idea is doing the heavy lifting.

Our gains over the SOTA LangView [35] show that learning view selection from language is less effective than that from large-scale human-edited videos, even when the videos and language are available at scale (bigData). Moreover, the insights of LangView and our work are complementary. We find if we fine-tune SWITCH-A-VIEW with LangView’s narration-based pseudo-labels, in *addition* to our labels (Sec. 4), we achieve further gains. See Fig. 3b, and Supp. for experiment details.

Ablations. Table 3 shows our ablation results for view-switch detection. Dropping any one input to our model degrades performance, indicating that each input plays a role. Dropping two inputs hurt the performance even more, showing that higher amounts of inputs is better in any com-

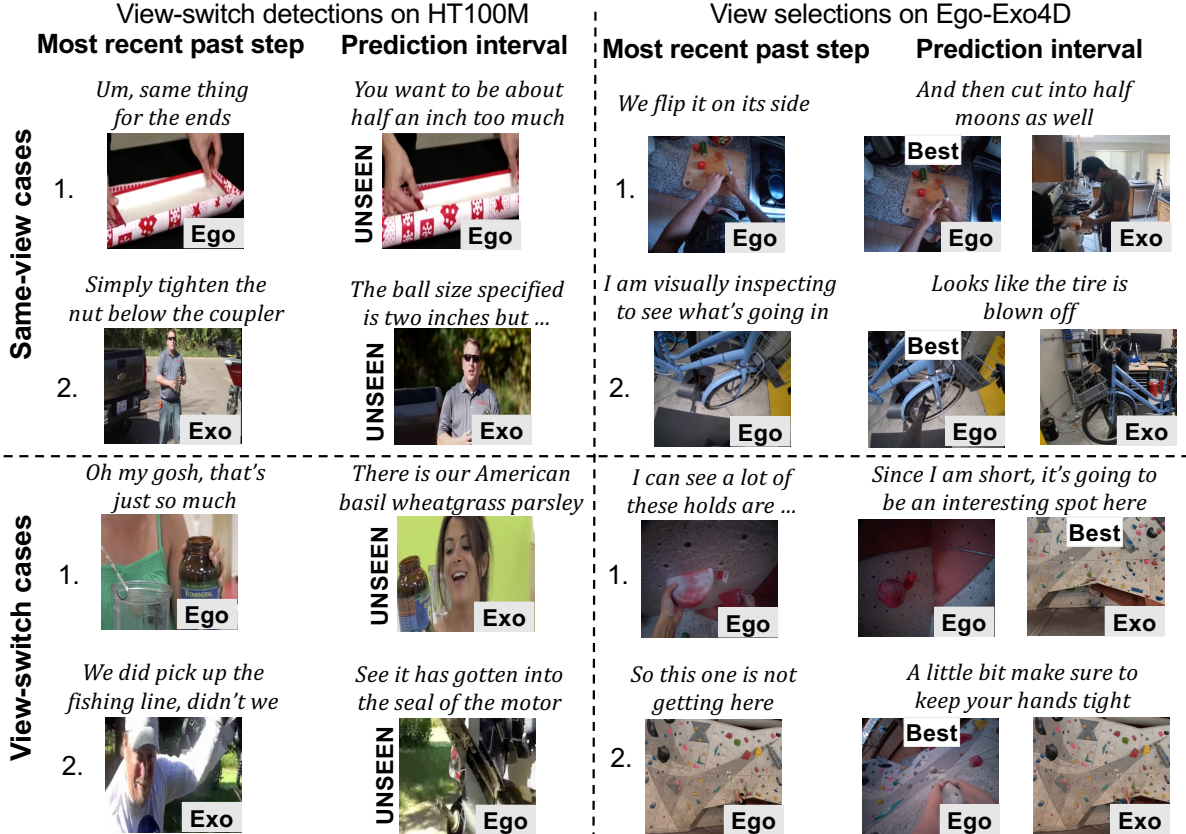


Figure 4. **Left:** successful view-switch detections by our model on same-view (**top**) and view-switch cases (**bottom**). Our model correctly detects view switches by popenopenotentially anticipating the next step using past frames (same-view sample 1, and view-switch sample 2) or leveraging the content of the next narration (same-view sample 2, and view-switch sample 1 and 2). **Right:** successful view selections by our model on same-view (**top**) and view-switch cases (**bottom**). For view selection as well, our model can predict the desired next view by relying on the next narration (same-view sample 1, and view-switch sample 1 and 2), or anticipate it using the past narrations (same-view sample 1 and 2), or the past frames (same-view sample 1). These examples show that all three inputs play a role in our model predictions.

bination, possibly because our model design enables extracting complementary cues from them in all configurations. Moreover, using past frames instead of narrations improves performance, re-affirming that vision provides fine-grained features necessary for high performance. Finally, using N' instead of N also helps improve performance in some cases, illustrating the next narration’s role.

In Fig. 3a, we study the effect of sample count on view selection. Our model already improves performance with as few as 1000 samples. This plot also highlights the low-shot success of our model vs. the best baseline, LangView-bigData.

See Supp. for more analysis, including the effect of the past frame and narration durations on model performance, and its scenario-level breakdown.

Qualitative examples. Fig. 4-left shows our model’s successful view-switch detections on both same-view (**top**) and view-switch cases (**bottom**); see caption for details. We also notice some common **failure modes** with our model.

For view-switch detection, our model sometimes fails when there is no next narration overlapping with the prediction interval, and neither the past frames nor narrations are predictive of the next view. In another failure type, the past views are wrongly categorized by our pseudo-labeler for HowTo100M [37] or by professional annotators for Ego-Exo4D [19]. This leads to our model getting confused and predicting the wrong next view. For view selection, in addition to these failures, our model can fail when both views look equally good. See Supp. for video examples.

6. Conclusion

We provided an approach for learning to select views from instructional video by bootstrapping already-edited, unlabeled in-the-wild content. Results show the method’s efficacy. In future work, we plan to generalize to *continuous* view selection, potentially by integrating ideas from new view synthesis, and will explore modeling user attention for personalized view selection.

References

- [1] Phil Ammirato, Patrick Poirson, Eunbyung Park, Jana Košecká, and Alexander C Berg. A dataset for developing and benchmarking active vision. In *2017 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1378–1385. IEEE, 2017. 2
- [2] Ido Arev, Hyun Soo Park, Yaser Sheikh, Jessica Hodgins, and Ariel Shamir. Automatic editing of footage from multiple social cameras. *ACM Trans. Graph.*, 33(4), 2014. 2
- [3] Kumar Ashutosh, Santhosh Kumar Ramakrishnan, Triantafyllos Afouras, and Kristen Grauman. Video-mined task graphs for keystep recognition in instructional videos. *Advances in Neural Information Processing Systems*, 36, 2024. 2
- [4] Taivanbat Badamdorj, Mrigank Rochan, Yang Wang, and Li Cheng. Contrastive learning for unsupervised video highlight detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 14042–14052, 2022. 2
- [5] Brandon Castellano. Pyscenedetect. <https://github.com/Breakthrough/PySceneDetect>. 4, 15
- [6] Seunghoon Cha, Jungjin Lee, Seunghwa Jeong, Younghui Kim, and Junyong Noh. Enhanced interactive 360° viewing via automatic guidance. *ACM Trans. Graph.*, 39(5), 2020. 2
- [7] Jianhui Chen, Keyu Lu, Sijia Tian, and Jim Little. Learning sports camera selection from internet videos. In *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1682–1691. IEEE, 2019. 1, 2
- [8] Ricson Cheng, Ziyang Wang, and Katerina Fragkiadaki. Geometry-aware recurrent neural networks for active visual recognition. *Advances in Neural Information Processing Systems*, 31, 2018. 2
- [9] Shih-Han Chou, Yi-Chun Chen, Kuo-Hao Zeng, Hou-Ning Hu, Jianlong Fu, and Min Sun. Self-view grounding given a narrated 360 {deg} video. *arXiv preprint arXiv:1711.08664*, 2017. 2
- [10] Wen-Sheng Chu, Yale Song, and Alejandro Jaimes. Video co-summarization: Video summarization by visual co-occurrence. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3584–3592, 2015. 2
- [11] J. Cohen. A coefficient of agreement for nominal scales. *Educational and Psychological Measurement*, 20(1):37–46, 1960. 6
- [12] Ruoyi Du, Wenqing Yu, Heqing Wang, Ting-En Lin, Dongliang Chang, and Zhanyu Ma. Multi-view active fine-grained visual recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 1568–1578, 2023. 2
- [13] Mohamed Elfeki, Liqiang Wang, and Ali Borji. Multi-stream dynamic video summarization. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 339–349, 2022. 2
- [14] 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. 16
- [15] Christoph Feichtenhofer, Yanghao Li, Kaiming He, et al. Masked autoencoders as spatiotemporal learners. *Advances in neural information processing systems*, 35:35946–35958, 2022. 16, 17
- [16] J. Foote and D. Kimber. Flycam: practical panoramic video and automatic camera control. In *2000 IEEE International Conference on Multimedia and Expo. ICME2000. Proceedings. Latest Advances in the Fast Changing World of Multimedia (Cat. No.00TH8532)*, pages 1419–1422 vol.3, 2000. 1, 2
- [17] Michael Gleicher and James Masanz. Towards virtual videography (poster session). In *Proceedings of the Eighth ACM International Conference on Multimedia*, page 375–378, New York, NY, USA, 2000. Association for Computing Machinery. 1, 2
- [18] Boqing Gong, Wei-Lun Chao, Kristen Grauman, and Fei Sha. Diverse sequential subset selection for supervised video summarization. In *Advances in Neural Information Processing Systems*. Curran Associates, Inc., 2014. 2
- [19] Kristen Grauman, Andrew Westbury, Lorenzo Torresani, Kris Kitani, Jitendra Malik, Triantafyllos Afouras, Kumar Ashutosh, Vijay Baiyya, Siddhant Bansal, Bikram Boote, et al. Ego-exo4d: Understanding skilled human activity from first-and third-person perspectives. *arXiv preprint arXiv:2311.18259*, 2023. 2, 3, 5, 6, 7, 8, 12, 14, 15, 16
- [20] Michael Gygli, Helmut Grabner, Hayko Riemenschneider, and Luc Van Gool. Creating summaries from user videos. In *Computer Vision – ECCV 2014*, pages 505–520, Cham, 2014. Springer International Publishing. 2
- [21] Bo He, Jun Wang, Jieli Qiu, Trung Bui, Abhinav Shrivastava, and Zhaowen Wang. Align and attend: Multimodal summarization with dual contrastive losses. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 14867–14878, 2023. 2
- [22] Rachel Heck, Michael Wallick, and Michael Gleicher. Virtual videography. In *Proceedings of the 14th ACM International Conference on Multimedia*, page 961–962, New York, NY, USA, 2006. Association for Computing Machinery. 2
- [23] Hou-Ning Hu, Yen-Chen Lin, Ming-Yu Liu, Hsien-Tzu Cheng, Yung-Ju Chang, and Min Sun. Deep 360 pilot: Learning a deep agent for piloting through 360deg sports videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3451–3460, 2017. 1, 2
- [24] Dinesh Jayaraman and Kristen Grauman. End-to-end policy learning for active visual categorization. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 41(7):1601–1614, 2019. 2
- [25] Hong-Wen Kang, Y. Matsushita, Xiaou Tang, and Xue-Quan Chen. Space-time video montage. In *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR’06)*, pages 1331–1338, 2006. 2
- [26] 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. 16
- [27] J Landis and G Koch. The measurement of observer agreement for categorical data. *Biometrics*, 1977. 6

- [28] Yandong Li, Liqiang Wang, Tianbao Yang, and Boqing Gong. How local is the local diversity? reinforcing sequential determinantal point processes with dynamic ground sets for supervised video summarization. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 151–167, 2018. 2
- [29] Yanghao Li, Tushar Nagarajan, Bo Xiong, and Kristen Grauman. Ego-exo: Transferring visual representations from third-person to first-person videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6943–6953, 2021. 4
- [30] Xudong Lin, Fabio Petroni, Gedas Bertasius, Marcus Rohrbach, Shih-Fu Chang, and Lorenzo Torresani. Learning to recognize procedural activities with distant supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13853–13863, 2022. 2
- [31] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 16, 17
- [32] Zheng Lu and Kristen Grauman. Story-driven summarization for egocentric video. In *Proceedings of the 2013 IEEE Conference on Computer Vision and Pattern Recognition*, page 2714–2721, USA, 2013. IEEE Computer Society. 2
- [33] Yiwei Ma, Guohai Xu, Xiaoshuai Sun, Ming Yan, Ji Zhang, and Rongrong Ji. X-clip: End-to-end multi-grained contrastive learning for video-text retrieval. In *Proceedings of the 30th ACM international conference on multimedia*, pages 638–647, 2022. 13, 14
- [34] Sagnik Majumder, Tushar Nagarajan, Ziad Al-Halah, Reina Pradhan, and Kristen Grauman. Which viewpoint shows it best? language for weakly supervising view selection in multi-view videos. *arXiv preprint arXiv:2411.08753*, 2024. 2, 6, 7, 12, 13
- [35] Sagnik Majumder, Tushar Nagarjan, Ziad Al-Halah, Reina Pradhan, and Kristen Grauman. Which viewpoint shows it best? language for weakly supervising view selection in multi-view videos. 2024. 2, 7
- [36] Effrosyni Mavroudi, Triantafyllos Afouras, and Lorenzo Torresani. Learning to ground instructional articles in videos through narrations. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 15201–15213, 2023. 2
- [37] 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 *Proceedings of the IEEE/CVF international conference on computer vision*, pages 2630–2640, 2019. 2, 3, 4, 5, 6, 7, 8, 12, 13, 14, 15
- [38] Medhini Narasimhan, Arsha Nagrani, Chen Sun, Michael Rubinstein, Trevor Darrell, Anna Rohrbach, and Cordelia Schmid. Tl; dw? summarizing instructional videos with task relevance and cross-modal saliency. In *European Conference on Computer Vision*, pages 540–557. Springer, 2022. 2
- [39] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023. 6, 16
- [40] Rameswar Panda and Amit K Roy-Chowdhury. Collaborative summarization of topic-related videos. In *Proceedings of the IEEE Conference on computer vision and pattern recognition*, pages 7083–7092, 2017. 2
- [41] Danila Potapov, Matthijs Douze, Zaid Harchaoui, and Cordelia Schmid. Category-specific video summarization. In *Computer Vision – ECCV 2014*, pages 540–555, Cham, 2014. Springer International Publishing. 2
- [42] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 6, 12, 13, 14
- [43] Santhosh K Ramakrishnan, Dinesh Jayaraman, and Kristen Grauman. Emergence of exploratory look-around behaviors through active observation completion. *Science Robotics*, 4(30):eaaw6326, 2019. 2
- [44] Mrigank Rochan, Linwei Ye, and Yang Wang. Video summarization using fully convolutional sequence networks. In *Proceedings of the European conference on computer vision (ECCV)*, pages 347–363, 2018. 2
- [45] Mrigank Rochan, Mahesh Kumar Krishna Reddy, Linwei Ye, and Yang Wang. Adaptive video highlight detection by learning from user history. In *Computer Vision – ECCV 2020*, pages 261–278, Cham, 2020. Springer International Publishing. 2
- [46] Abhimanyu Sahu and Ananda S. Chowdhury. Shot level egocentric video co-summarization. In *2018 24th International Conference on Pattern Recognition (ICPR)*, pages 2887–2892, 2018. 2
- [47] Soroush Seifi and Tinne Tuytelaars. Attend and segment: Attention guided active semantic segmentation. In *European Conference on Computer Vision*, pages 305–321. Springer, 2020. 2
- [48] Soroush Seifi, Abhishek Jha, and Tinne Tuytelaars. Glimpse-attend-and-explore: Self-attention for active visual exploration. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 16137–16146, 2021. 2
- [49] Gunnar A Sigurdsson, Abhinav Gupta, Cordelia Schmid, Ali Farhadi, and Karteek Alahari. Actor and observer: Joint modeling of first and third-person videos. In *proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7396–7404, 2018. 14
- [50] Gunnar A Sigurdsson, Abhinav Gupta, Cordelia Schmid, Ali Farhadi, and Karteek Alahari. Charades-ego: A large-scale dataset of paired third and first person videos. *arXiv preprint arXiv:1804.09626*, 2018. 4, 14
- [51] Michael Smith and Takeo Kanade. Video skimming for quick browsing based on audio and image characterization. Technical Report CMU-CS-95-186, Pittsburgh, PA, 1995. 2
- [52] Yale Song, Jordi Vallmitjana, Amanda Stent, and Alejandro Jaimes. Tvsum: Summarizing web videos using titles. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5179–5187, 2015. 2
- [53] Yu-Chuan Su, Dinesh Jayaraman, and Kristen Grauman. Pano2vid: Automatic cinematography for watching 360

- videos. In *Asian Conference on Computer Vision*, pages 154–171. Springer, 2016. [1](#), [2](#)
- [54] Xinding Sun, J. Foote, D. Kimber, and B. S. Manjunath. Region of interest extraction and virtual camera control based on panoramic video capturing. *Trans. Multi.*, 7(5):981–990, 2005. [2](#)
- [55] Yansong Tang, Dajun Ding, Yongming Rao, Yu Zheng, Danyang Zhang, Lili Zhao, Jiwen Lu, and Jie Zhou. Coin: A large-scale dataset for comprehensive instructional video analysis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1207–1216, 2019. [3](#), [4](#)
- [56] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023. [6](#), [16](#)
- [57] Du Tran, Heng Wang, Lorenzo Torresani, Jamie Ray, Yann LeCun, and Manohar Paluri. A closer look at spatiotemporal convolutions for action recognition. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 6450–6459, 2018. [16](#)
- [58] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. [4](#), [6](#), [17](#)
- [59] Yi Wang, Kunchang Li, Xinhao Li, Jiashuo Yu, Yan He, Guo Chen, Baoqi Pei, Rongkun Zheng, Zun Wang, Yansong Shi, et al. Internvideo2: Scaling foundation models for multimodal video understanding. In *European Conference on Computer Vision*, pages 396–416. Springer, 2024. [2](#), [6](#), [7](#), [12](#), [13](#)
- [60] Bo Xiong and Kristen Grauman. Snap angle prediction for 360° panoramas. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018. [2](#)
- [61] Cha Zhang, Yong Rui, Jim Crawford, and Li-Wei He. An automated end-to-end lecture capture and broadcasting system. *ACM Trans. Multimedia Comput. Commun. Appl.*, 4(1), 2008. [1](#), [2](#)
- [62] Ke Zhang, Wei-Lun Chao, Fei Sha, and Kristen Grauman. Summary transfer: Exemplar-based subset selection for video summarization. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1059–1067, 2016. [2](#)
- [63] Yiwu Zhong, Licheng Yu, Yang Bai, Shangwen Li, Xueting Yan, and Yin Li. Learning procedure-aware video representation from instructional videos and their narrations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14825–14835, 2023. [2](#)
- [64] Honglu Zhou, Roberto Martín-Martín, Mubbasir Kapadia, Silvio Savarese, and Juan Carlos Niebles. Procedure-aware pretraining for instructional video understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10727–10738, 2023. [2](#)

7. Supplementary material

In this supplementary material we provide additional details about:

- Video for qualitatively illustrating of our main idea and also qualitatively evaluating of our view-switch detections and view selections (Sec. 7.1), as mentioned in ‘Qualitative examples’ in Sec. 5 in main
- Analysis of the impact of our shot-level pseudo-labeling on view-switch detection performance (Sec. 7.2), as referenced in Sec. 3.2 in main
- Annotation filtering and model evaluation with higher inter-annotator agreement thresholds (Sec. 7.3), as noted in ‘Annotator agreement on best view’ in Sec. 4 in in main
- View-selection results upon finetuning our view-switch detector jointly with narration-based pseudo-labels and our best view labels (Sec. 7.4), as mentioned in ‘View selection’ in Sec. 5 in main
- Analysis of the impact of the duration of our past frames on view-switch detection performance (Sec. 7.5), as noted in ‘Ablations’ in Sec. 5 in main
- Analysis of the impact of the duration of our past narrations on view-switch detection performance (Sec. 7.6), as referenced in ‘Ablations’ in Sec. 5 in main
- Scenario-level breakdown of view selection performance (Sec. 7.7), as noted in ‘Ablations’ in Sec. 5 in main
- Feature similarity baseline evaluation with CLIP [42]-style encoders (Sec. 7.8), as mentioned in ‘Baselines’ in Sec. 5 in main
- Dataset details (Sec. 7.9) in addition to the ones provided in Sec. 3.2, and ‘Training data’, ‘Evaluation data’ and ‘Data for view selection with limited labels’ in Sec. 4 in main
- Annotation details (Sec. 7.10) in addition to the ones provided in ‘Evaluation data’ and ‘Annotator agreement on best view’ in Sec. 4 in main
- Additional implementation details (Sec. 7.11), as referenced in ‘Implementation’ in Sec. 5 in main

7.1. Supplementary video

The supplementary video, available at https://vision.cs.utexas.edu/projects/switch_a_view/, qualitatively illustrates our task, View Selection with Limited Labels, and our main idea towards tackling that task, Weakly-Supervised Learning from Unlabeled In-the-wild Videos. We also show successful predictions by our model for both view-switch detection and view selection. For view selection, we additionally provide multi-step selection examples, where our model selects the best view over multiple consecutive steps. Finally, we illustrate our model’s common failure modes (‘Qualitative examples’ in Sec. 5 in main) with qualitative examples.

Model	Accuracy	AUC	AP
Ours w/o shot-level pseudo-labeling	51.5	51.9	52.3
Ours	59.4	63.8	60.5

Table 4. Analysis of the impact of our shot-level pseudo-labeling strategy on view-switch detection performance on the HowTo100M [37] dataset. All values are in %, and higher is better. Significance $p \leq 0.05$.

7.2. Shot-level pseudo-labeling

In Table 4, we report the results for an additional ablation study, in which we analyze the impact of our shot-level pseudo-labeling strategy (Sec. 3.2 in main) on view-switch detection with the HowTo100M [37] dataset. Upon replacing our shot-level pseudo-labeling strategy with a clip-level pseudo-labeling strategy (Sec 3.2 in main), we observe a drastic drop in model performance. This demonstrates that our pseudo-labeler is able to mitigate noisy clip-level predictions, particularly at scene boundaries.

7.3. Inter-annotator agreement threshold

In main, we evaluated our models with an inter-annotator agreement thresholds of 78%, meaning at least 7 out of 9 agree for each annotation instance (‘Annotator agreement on best view’ in Sec. 4 in main). Here, we evaluate even higher agreement thresholds of 89%—at least 8 out of 9 annotators agree, and 90%—all annotators agree. For HT100M [37], the number of samples drops from 3,151 to 1,840 at 80%, and 1,345 at 90%. For Ego-Exo4D [19], the same goes down from 5,049 to 3,421 at 80%, and 1,887 at 90%, respectively. Table 5 reports the results. Even at these higher and more challenging inter-annotator agreement thresholds, our model outperforms the strongest baseline—Retrieval [59]- F for view-switch detection, and LangView [34]-bigData for view selection—on both tasks.

7.4. Finetuning jointly with narration-based pseudo-labels and best view labels, for view selection

Here, we provide details on how we finetune our view-switch detector jointly with narration-based pseudo-labels [34] and our best view labels (‘Data for view selection with limited labels’ in Sec. 4 in main), for doing view selection. Essentially, we modify our view selector training loss \mathcal{L}^S (Sec. 3.5 in main) as follows:

$$\mathcal{L}^S = \mathcal{L}_{\text{CE}}(\check{V}_{(t_w, t_w + \Delta)}, V_{(t_w, t_w + \Delta)}) + \alpha * \mathcal{L}^{N'}, \quad (8)$$

where $\mathcal{L}^{N'} = \mathcal{L}_{\text{CE}}(\check{V}_{(t_w, t_w + \Delta)}, \tilde{V}_{(t_w, t_w + \Delta)}^{N'})$ is the cross-entropy loss between the predicted views $\check{V}_{(t_w, t_w + \Delta)}$ (Sec. 3.4 in main) and narration-based pseudo-labels [34] $\tilde{V}_{(t_w, t_w + \Delta)}^{N'}$, generated using next narrations N' (Sec. 3.1 in

Model	View-switch detection on HT100M			View selection on Ego-Exo4D		
	78%	89%	100%	78%	89%	100%
Retrieval [59]- F	53.2	53.6	53.6	—	—	—
View-narration [59] Similarity	—	—	—	53.9	54.2	53.7
LangView [34]-bigData	—	—	—	54.5	54.9	54.1
Ours	60.5	60.8	60.7	56.0	56.2	55.3

Table 5. Model performance (AP) vs. inter-annotator agreement threshold. — indicates that the baseline is not applicable for the particular task. All values are in %, and higher is better. Significance $p \leq 0.05$.

Model	Accuracy	AUC	AP
$T^F = 2$	59.4	63.1	59.6
$T^F = 4$	<u>59.0</u>	<u>63.4</u>	<u>60.2</u>
$T^F = 8$ (Ours)	59.4	63.8	60.5
$T^F = 16$	55.4	59.1	57.0
$T^F = 32$	52.8	55.0	53.6

Table 6. Analysis of the impact of the duration of past frames, denoted with T^F , on view-switch detection performance on the HowTo100M [37] dataset. All values are in %, and higher is better. Significance $p \leq 0.05$.

Model	Accuracy	AUC	AP
$T^N = 2$	56.1	55.9	56.2
$T^N = 4$	52.4	53.9	53.4
$T^N = 8$	55.5	60.2	58.0
$T^N = 16$	<u>56.1</u>	<u>60.2</u>	<u>58.0</u>
$T^N = 32$ (Ours)	59.4	63.8	60.5

Table 7. Analysis of the impact of the duration of past narrations, denoted with T^N , on view-switch detection performance on the HowTo100M [37] dataset. All values are in %, and higher is better. Significance $p \leq 0.05$.

main), and α is the weight on $\mathcal{L}^{N'}$, which we set α to 0.3 on the basis of validation.

7.5. Duration of past frames

In Table 6, we report our view-switch detection performance numbers for different durations of past frames, denoted by T^F , using the HowTo100M [37] dataset. We notice that our model performance declines monotonically as we move from our choice of $T^F = 8$ seconds (‘Implementation’ in Sec. 5 in main) to both smaller and larger values. While very short visual contexts fail to capture long-range temporal patterns in human-preferred view sequences, very long visual contexts might contain spurious signals that affect model performance. $T^F = 8$ seconds balances this trade-off and leads to the best model performance, per this study.

7.6. Duration of past narrations

In Table 7, we report our view-switch detection results for different durations of past narrations, denoted by T^N , with HowTo100M [37]. Upon reducing T^N to values lower than our choice of 32 seconds (‘Implementation’ in Sec. 5 in

main), our model performance declines monotonically. This shows that a longer past narration context helps better learn correlations between the text in the narrations and desired view types.

7.7. Scenario-level analysis of view-selection performance

Fig. 5 shows the breakdown of view selection performance per scenario, where only the scenarios with a minimum of 10 instances after filtering low-quality annotations (‘Data for view selection with limited labels’ in Sec. 4 in main) are shown. Compared to the best-performing baseline, LangView-bigData, our model’s performance goes up both in absolute and relative terms, from the procedural scenarios like Cooking or Bike Repair, to physical scenarios like Rock climbing. This demonstrates that our model is better able to handle more scenarios with more physical activity, and consequently, more view changes, than the best baseline.

7.8. Feature similarity baselines with CLIP [42] encoders

In ‘View-switch detection’ and ‘View selection’ in Sec. 5 in main, we evaluated feature similarity baselines with InternVideo2 [59] encoders. Here, we provide the parallel experiment with CLIP [42]-style encoders in Table 8 for view-switch detection, and Table 9 for view selection. Specifically, while for the retrieval baselines, we use the unmodified CLIP encoders, we use X-CLIP [33] encoders in the View-narration Similarity baseline for encoding multiple frames

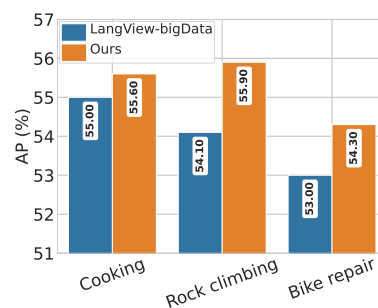


Figure 5. Per-scenario breakdown of our and the strongest baseline, LangView-bigData’s view selection performance, measured with AP (%).

Model	HowTo100M [37]			Ego-Exo4D [19]		
	Accuracy	AUC	AP	Accuracy	AUC	AP
Retrieval [42]- F	<u>53.0</u>	<u>53.0</u>	52.7	52.1	<u>52.1</u>	<u>53.4</u>
Retrieval [42]- N	52.3	52.3	51.8	51.8	51.8	50.7
Retrieval [42]- N'	52.6	52.6	<u>53.2</u>	52.0	52.0	52.5
Ours	59.4	63.8	60.5	<u>51.2</u>	56.4	55.4

Table 8. View-switch detection results. Evaluation on Ego-Exo4D [19] is zero-shot. All values are in %, and higher is better. Significance $p \leq 0.05$.

Model	Accuracy	AUC	AP
Retrieval [42]- F	52.1	52.1	<u>53.4</u>
Retrieval [42]- N	51.8	51.8	50.7
Retrieval [42]- N'	52.0	52.0	52.5
View-narration [33] Similarity	<u>52.5</u>	<u>52.2</u>	<u>53.4</u>
Ours	54.0	57.3	56.0

Table 9. Results and ablation study for view selection with limited labels. All values are in %, and higher is better. Significance $p \leq 0.05$.

in the ego and exo views (Sec. 3.4 in main). With CLIP, the similarity baselines generally perform worse. This happens possibly because, unlike InternVideo2, CLIP features are not very fine-grained and/or do not capture temporal context in the case of retrieval baselines. Furthermore, our model outperforms all feature similarity baselines, even when implemented with CLIP, highlighting its advantages over the baselines across different encoder choices.

7.9. Additional dataset details

Here, we give further dataset details, in addition to the ones provided in ‘Training data’, ‘Evaluation data’ and ‘Data for view selection with limited labels’ in Sec. 4 in main.

Charades-Ego [49] datasets for training view classifier in our pseudo-labeler. As mentioned in Sec. 3.2 in main in main, we train the view classifier of our pseudo-labeler on the Charades-Ego [50] dataset. To do so, we create a dataset containing 5,551 train, 615 val, and 1,597 test videos, where all videos are randomly sampled and the splits are completely disjoint. Moreover, we train and test our model by sampling fixed-length clips, where the clip length is set to 2 seconds.

HowTo100M [37] datasets for view-switch detection. As noted in ‘Training data’ and ‘Evaluation data’ in Sec. 4 in main in main, we train our view-switch detector on HowTo100M (HT100M) [37] and also use it as a dataset for evaluating view-switch detection. To do so, we sample a maximum of 500 videos from each category in the second level of the HT100M video classification hierarchy. This results in a total of 4,391 hours of HT100M videos.

In addition to the details provided in ‘Evaluation data’ in Sec. 4 in main, for creating the HT100M test set for evaluating view-switch detection, we also include the clips right after all view-switch boundaries, as identified by our pseudo-labeler, in the test set. This ensures that the test set is not totally dominated by the more frequently-occurring same-view instances, which can affect the estimation of our unbiased mean performance (‘Evaluation metrics’ in Sec. 5 in main). Finally, we provide the clip just before each clip being labeled, to the annotators in order to identify instances where the *ground-truth* view stays the same (same-view) and where it switches (view-switch). This allows us to separately evaluate these two alternate but important scenarios, and report *unbiased* mean performance (‘Evaluation metrics’ in Sec. 5 in main).

Ego-Exo4D [19] datasets. We create our datasets for Ego-Exo4D [19] (‘Evaluation data’ and ‘Data for view selection with limited labels’ in Sec. 4 in main) for evaluating view-switch detection, and training and evaluating view selection, by sampling clips from each video at a regular interval of 1 second.

7.10. Additional annotation details

Here, we provide further annotation details, in addition to what are provided in (‘Evaluation data’ and ‘Annotator agreement on best view’ in Sec. 4 in main). We start with the details that are common for both both HT100M [37] and Ego-Exo4D [19].

We use Amazon Mechanical Turk (MTurk) to collect annotations for both datasets. Before assigning an MTurk worker our job, we ensure that their prior annotation approval rate is more than 98%. We also require them to take pass short qualifiers (‘Annotator agreement on best view’ in Sec. 4 in main), each of which contains 10 annotation instances. We design these qualifiers such that they are very similar to our main jobs. Furthermore, we handpick the annotation instances in the qualifiers such we that know the ground-truth for each of them. This lets us easily compare an annotator’s choices against the ground-truths in the qualifiers, and consequently, gauge the their reliability. We only accept annotators who pass these qualifiers with a success rate of at least 90%.

Next, we provide dataset-specific annotation details.

HT100M. In Fig. 6, we show our annotation interface for HT100M (‘Evaluation data’ and ‘Annotator agreement on best view’ in Sec. 4 in main). In short, we provide a set of detailed instructions, which lists the different characteristics of both ego and exo clips⁴, and give examples for each characteristic. Additionally, we provide a more concise summary

⁴We refer to ego and exo views as “closeup” and “wide” shots, respectively, in order to easily explain the annotation process to the workers.

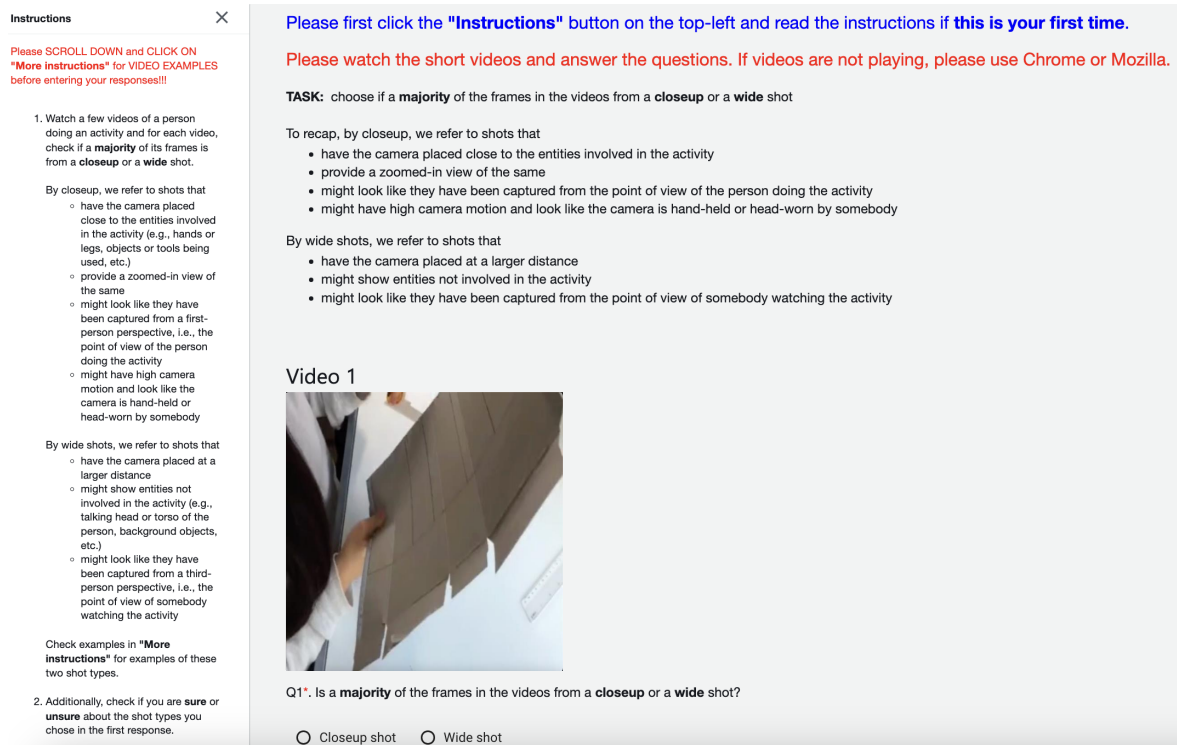


Figure 6. Sample interface for collecting HT100M [37] annotations (‘Evaluation data’ and ‘Annotator agreement on best view’ in Sec. 4 in main). Additionally, we also provide video examples for both ego (closeup shot) and exo (wide shot) clips, to help the annotators.

of the lists of per-view attributes on each annotation page to give a quick recap of the annotation task, to the workers. Finally, we provide video examples of both ego and exo clips, to further guide the annotation process.

Ego-Exo4D. In Fig. 7, we show our annotation interface for Ego-Exo4D (‘Evaluation data’ and ‘Annotator agreement on best view’ in Sec. 4 in main). In summary, we provide a set of detailed instructions on each annotation page, which describes the kind of information captured by the two views (ego and exo) and the role of the associated atomic description, and also specify that we expect the annotator to pick a view that best shows the activity mentioned in the atomic description and hence, useful for instructional purposes. To further assist with the annotation process, we provide examples showing pairs of clips from both ego and exo views, their associated atomic descriptions, and how an annotator should reason about which view is better for viewing the activity, in the context of its corresponding atomic description.

7.11. Additional implementation details

Here, we provide additional implementation details.

View assignment to past frames and narrations for Ego-Exo4D [19]. In Sec. 3.2 in main, we provide details about

how we assign our view pseudo-labels to past frames and past narrations, for using our model (view-switch detector or view selector) with HowTo100M [37]. Here, we describe our process for assigning views to the past frames and narrations for Ego-Exo4D [19].

Note that all frames or narrations in an Ego-Exo4D video might not be assigned a ground-truth best view during our annotation process, because the annotations for some couplets of pairs of clips from the two views, and their associated atomic descriptions (‘Evaluation data’ in Sec. 4 in main), might get discarded due to our annotation quality control measures (‘Annotator agreement on best view’ in Sec. 4 in main, and above). To tackle this, we set the best view for each past frame and past narration to the best view ground-truth for the nearest couplet of clip pair and its associated atomic description, for which the ground-truth has not been discarded.

7.11.1. View pseudo-labeler

Scene detector. As mentioned in Sec. 3.2 in main, we use PySceneDetect [5], an off-the-shelf scene detector, for detecting scene boundaries in HowTo100M [37] videos. Specifically, we use the image-content-based detector. Moreover, we set the weights for pixel colors to 1.0 and the same for object edges to 0.0, and the minimum shot length to 2 seconds.

Please first click the **"Instructions"** button on the top-left and read the instructions if **this is your first time**.

Please watch the short videos and answer the questions. If videos are not playing, please use Chrome or Mozilla.

TASK: Watch pairs of videos and choose which video--**Left** or **Right**--**more accurately** shows the **activity** mentioned in the **text** written below each pair, and is **also more useful** for **teaching** somebody **how to** do the activity ?

For each pair

- a person performs an activity, which is shown from two viewpoints
- the video on the **left** is captured from the **perspective of the person doing the activity**, and provides a **first-person view** of the activity
- the video on the **right** is captured with a camera placed next to the site of the activity, and provides a **third-person view** of the activity
- the **text** below each pair provides a **fine-grained description** of the activity.

By **'choose which video more accurately shows the activity in the text'**, we ask you to choose the video that more clearly shows the body parts and other objects involved in the activity, and these entities interact with each other, as described in the text, and is **also more useful** for **teaching** somebody **how to** do the activity.

Note that

- the letter **'C'**, when used as the subject in the text, denotes the person doing the activity
- other uppercase letters, like **'X'**, **'O'**, etc., when used in the text, denote other persons involved in the activity

Pair 1



Text: C places the garlic in the bowl on the countertop, using his right hand

Q1*. Which video--**Left** or **Right**--do you think more accurately shows the activity mentioned in the text?

By accurate, we mean that the video you choose should more clearly show the human body parts and objects, and their interactions as mentioned in the text. (refer to instructions for examples).

Left Right

Figure 7. Sample interface for collecting Ego-Exo4D [19] annotation (‘Evaluation data’ and ‘Annotator agreement on best view’ in Sec. 4 in main) . Additionally, we also provide examples showing pairs of clips from both ego and exo views, their associated atomic descriptions, and how to reason about which view is better for viewing an activity, to help the annotators.

View classifier. For our view classifier (Sec. 3.2 in main), we use the slow branch of a SlowFast [14] model that has a ResNet3D [57]-50 architecture and is pretrained on the Kinetics-400 [26] dataset, as the visual encoder. This encoder takes 8 uniformly sampled frames from each 2-second clip (c.f. Sec. 7.10), embeds them into a visual feature, and passes the visual feature to a linear classification head, which is implemented as a single linear layer. We initialize the parameters for all model components that are trained from scratch, using a parameter initialization strategy for masked auto-encoders for videos [15]. We train the model until convergence by using an AdamW [31] optimizer, a batch size of 32, and initial learning rates of 10^{-5} and 10^{-4} for the visual

encoder, and the classification head, respectively.

7.11.2. View-switch detector

Here, we provide more implementation details of view-switch detector, in addition to those provided in ‘Implementation’ in Sec. 5 in main. Our detector’s DINOv2 [39] frame encoder has 12 layers and takes in frames sampled at 4 fps, and produces a 768-dimensional feature for each frame. Our detector’s Llama 2 [56] text encoder begins by producing a token sequence for each input narration by tokenizing its text, and padding the tokenizer output to match the length of the longest token sequence in a batch, or truncating it to reduce its length to 512, depending on which length is shorter. Moreover, for the past narrations, it ensures that

their total token count does not cross 1024, by truncating wherever necessary. It then encodes each token from the past narration sequence, or the next narration, into a 4096-dimensional features. Next, it projects each such feature into a 768-dimensional feature using a linear layer. We implement our modules for encoding views for both frames and past narrations, and their relative temporal positions, as learnable embeddings that produce 768-dimensional features. Specifically, for encoding views, we use a learnable feature dictionary with 0 (ego) and 1 (exo) as keys. For encoding relative temporal positions, we first discretize the durations in seconds, by using bins of size 0.1 second, and then encode them with learnable feature dictionaries, in which the number of keys is set to the maximum number of bins. To aggregate all the above-mentioned 768-dimensional features, we use a 8-layer and 8-head transformer encoder [58] that adds a positional embedding [58] to each feature and performs self-attention on them. Finally, the 2-layer MLP for our view classification is a stack of two hidden linear layers with the output feature size of 256 and 64, respectively, and a final linear layer that estimates the next view. We initialize the parameters for all model components that are trained from scratch, using a parameter initialization strategy for masked auto-encoders for videos [15]. We freeze the pre-trained components of our view-switch detector and train the rest of the model until convergence with the AdamW [31] optimizer, a batch size of 48, and a learning rate of 10^{-6} .

7.11.3. View selector

Our view selector has the exact same architecture as our view-switch detector. For training it, we first initialize its parameters with those of our pretrained view-switch detector. We then freeze the frame encoder and the text encoder, and finetune the rest of the model until convergence with the exact same optimizer and hyperparameters as the ones used in our view-switch detector training.